

Baoxiang Liu
Maode Ma
Jincai Chang (Eds.)

LNCS 7473

Information Computing and Applications

Third International Conference, ICICA 2012
Chengde, China, September 2012
Proceedings



Springer

Commenced Publication in 1973

Founding and Former Series Editors:

Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen

Editorial Board

David Hutchison

Lancaster University, UK

Takeo Kanade

Carnegie Mellon University, Pittsburgh, PA, USA

Josef Kittler

University of Surrey, Guildford, UK

Jon M. Kleinberg

Cornell University, Ithaca, NY, USA

Alfred Kobsa

University of California, Irvine, CA, USA

Friedemann Mattern

ETH Zurich, Switzerland

John C. Mitchell

Stanford University, CA, USA

Moni Naor

Weizmann Institute of Science, Rehovot, Israel

Oscar Nierstrasz

University of Bern, Switzerland

C. Pandu Rangan

Indian Institute of Technology, Madras, India

Bernhard Steffen

TU Dortmund University, Germany

Madhu Sudan

Microsoft Research, Cambridge, MA, USA

Demetri Terzopoulos

University of California, Los Angeles, CA, USA

Doug Tygar

University of California, Berkeley, CA, USA

Gerhard Weikum

Max Planck Institute for Informatics, Saarbruecken, Germany

Baoxiang Liu Maode Ma Jincai Chang (Eds.)

Information Computing and Applications

Third International Conference, ICICA 2012
Chengde, China, September 14-16, 2012
Proceedings



Springer

Volume Editors

Baoxiang Liu
Jincai Chang
Hebei United University, College of Science
Tangshan 063000, Hebei, China
E-mail: 251983480@qq.com;1714064990@qq.com

Maode Ma
Nanyang Technological University, Singapore
E-mail: maode_ma@pmail.ntu.edu.sg

ISSN 0302-9743
ISBN 978-3-642-34061-1
DOI 10.1007/978-3-642-34062-8
Springer Heidelberg Dordrecht London New York

e-ISSN 1611-3349
e-ISBN 978-3-642-34062-8

Library of Congress Control Number: 2012948398

CR Subject Classification (1998): C.2, D.2, C.2.4, I.2.11, C.1.4, D.2.7

LNCS Sublibrary: SL 3 – Information Systems and Application, incl. Internet/Web and HCI

© Springer-Verlag Berlin Heidelberg 2012

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer. Violations are liable to prosecution under the German Copyright Law.

The use of general descriptive names, registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

Typesetting: Camera-ready by author, data conversion by Scientific Publishing Services, Chennai, India

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

Preface

Welcome to the proceedings of the 3rd International Conference on Information Computing and Applications (ICICA 2012), which was held during September 14–16, 2012, in Mountain Resort, Chengde, China.

As information technology advances, information computing and applications are becoming increasingly specialized. Information computing and applications including hardware, software, communications, and networks are growing with ever increasing scale and heterogeneity, and becoming overly complex. The complexity is getting more critical along with the growing number of applications. To cope with this complexity, information computing and applications focus on intelligent, selfmanageable, scalable computing systems and applications to the maximum extent possible without human intervention or guidance.

With the rapid development of information science and technology, information computing has become the third approach of science research. Information computing and applications is the field of study concerned with constructing intelligent computing, mathematical models, numerical solution techniques, and using computers to analyze and solve natural scientific, social scientific, and engineering problems. In practical use, it is typically the application of computer simulation, intelligent computing, internet computing, pervasive computing, scalable computing, trusted computing, autonomy-oriented computing, evolutionary computing, mobile computing, applications, and other forms of computation to problems in various scientific disciplines and engineering. Information computing and applications is an important underpinning for techniques used in information and computational science and there are many unresolved problems that are worth studying.

The ICICA 2012 conference provided a forum for engineers and scientists from academia, industry, and government to address the most innovative research and development including technical challenges and social, legal, political, and economic issues, and to present and discuss their ideas, results, work in progress, and experience on all aspects of information computing and applications.

There was a very large number of paper submissions (1089). All submissions were reviewed by at least three Program or Technical Committee members or external reviewers. It was extremely difficult to select the presentations for the conference because there were so many excellent and interesting submissions. In order to allocate as many papers as possible and keep the high quality of the conference, we finally decided to accept 330 papers for presentation, reflecting a 30.3% acceptance rate. A total of 100 papers were included in this volume. We believe that all of these papers and topics not only provided novel ideas, new results, work in progress, and state-of-the-art techniques in this field, but will also stimulate future research activities in the area of information computing and applications.

The exciting program for this conference was the result of the hard and excellent work of many others, such as Program and Technical Committee members, external reviewers, and Publication Chairs under a very tight schedule. We are also grateful to the members of the Local Organizing Committee for supporting us in handling so many organizational tasks, and to the keynote speakers for accepting to come to the conference with enthusiasm. Last but not least, we hope you enjoyed the conference program, and the beautiful attractions of Chengde, China.

September 2012

Yanchun Zhang
Baoliang Liu
Chunfeng Liu

Organization

ICICA 2012 was organized by Hebei United University, Hebei Scene Statistical Society, and sponsored by the National Science Foundation of China, Hunan Institute of Engineering, Yanshan University, Northeastern University at Qinhuangdao, and Chengde Petroleum College.

Executive Committee

Honourary Chair

Jun Li Hebei Polytechnic University, China

General Chairs

Yanchun Zhang University of Victoria, Australia
Baoxiang Liu Hebei Polytechnic University, China

Program Chairs

Chunfeng Liu Hebei Polytechnic University, China
Fengbo Hou Chengde Petroleum College, China
Wenjiang Du Chongqing Normal University, China

Local Arrangements Chairs

Jincai Chang Hebei Polytechnic University, China
Aimin Yang Hebei Polytechnic University, China

Steering Committee

Qun Lin Chinese Academy of Sciences, China
Maode Ma Nanyang Technological University, Singapore
Nadia Nedjah State University of Rio de Janeiro, Brazil
Lorna Uden Staffordshire University, UK
Yiming Chen Yanshan University, China
Changcun Li Hebei Polytechnic University, China
Zhijiang Wang Hebei Polytechnic University, China
Rongbo Zhu South-Central University for Nationalities,
China
Guohuan Lou Hebei Polytechnic University, China
Jixian Xiao Hebei Polytechnic University, China
Xinghuo Wan Hebei Polytechnic University, China
Chunying Zhang Hebei Polytechnic University, China
Dianchuan Jin Hebei Polytechnic University, China

Publicity Chairs

Aimin Yang	Hebei Polytechnic University, China
Xilong Qu	Hunan Institute of Engineering, China

Publication Chairs

Yuhang Yang	Shanghai Jiao Tong University, China
-------------	--------------------------------------

Financial Chair

Jincai Chang	Hebei Polytechnic University, China
--------------	-------------------------------------

Local Arrangements Committee

Lihong Li	Hebei Polytechnic University, China
Shaohong Yan	Hebei Polytechnic University, China
Yamian Peng	Hebei Polytechnic University, China
Lichao Feng	Hebei Polytechnic University, China
Yuhuan Cui	Hebei Polytechnic University, China

Secretaries

Kaili Wang	Hebei Polytechnic University, China
Jingguo Qu	Hebei Polytechnic University, China
Yafeng Yang	Hebei Polytechnic University, China

Program/Technical Committee

Yuan Lin	Norwegian University of Science and Technology, Norway
Yajun Li	Shanghai Jiao Tong University, China
Yanliang Jin	Shanghai University, China
Mingyi Gao	National Institute of AIST, Japan
Yajun Guo	Huazhong Normal University, China
Haibing Yin	Peking University, China
Jianxin Chen	University of Vigo, Spain
Miche Rossi	University of Padua, Italy
Ven Prasad	Delft University of Technology, The Netherlands
Mina Gui	Texas State University, USA
Nils Asc	University of Bonn, Germany
Ragip Kur	Nokia Research, USA
On Altintas	Toyota InfoTechnology Center, Japan
Suresh Subra	George Washington University, USA
Xiyin Wang	Hebei Polytechnic University, China
Dianxuan Gong	Hebei Polytechnic University, China
Chunxiao Yu	Yanshan University, China

Yanbin Sun	Beijing University of Posts and Telecommunications, China
Guofu Gui	CMC Corporation, China
Haiyong Bao	NTT Co., Ltd., Japan
Xiwen Hu	Wuhan University of Technology, China
Mengze Liao	Cisco China R&D Center, China
Yangwen Zou	Apple China Co., Ltd., China
Liang Zhou	ENSTA-ParisTech, France
Zhanguo Wei	Beijing Forestry University, China
Hao Chen	Hu'nan University, China
Lilei Wang	Beijing University of Posts and Telecommunications, China
Xilong Qu	Hunan Institute of Engineering, China
Duolin Liu	ShenYang Ligong University, China
Xiaozhu Liu	Wuhan University, China
Yanbing Sun	Beijing University of Posts and Telecommunications, China
Yiming Chen	Yanshan University, China
Hui Wang	University of Evry in France, France
Shuang Cong	University of Science and Technology of China, China
Haining Wang	College of William and Mary, USA
Zengqiang Chen	Nankai University, China
Dumisa Wellington Ngwenya	Illinois State University, USA
Hu Changhua	Xi'an Research Inst. of Hi-Tech, China
Juntao Fei	Hohai University, China
Zhao-Hui Jiang	Hiroshima Institute of Technology, Japan
Michael Watts	Lincoln University, New Zealand
Tai-hon Kim	Defense Security Command, Korea
Muhammad Khan	Southwest Jiaotong University, China
Seong Kong	The University of Tennessee, USA
Worap Kreesuradej	King Mongkut's Institute of Technology Ladkrabang, Thailand
Uwe Kuger	Queen's University of Belfast, UK
Xiao Li	CINVESTAV-IPN, Mexico
Stefanie Lindstaedt	Know-Center, Austria
Paolo Li	Polytechnic of Bari, Italy
Tashi Kuremoto	Yamaguchi University, Japan
Chun Lee	Howon University, South Korea
Zheng Liu	Nagasaki Institute of Applied Science, Japan
Michiharu Kurume	National College of Technology, Japan
Sean McLoe	National University of Ireland, Ireland
R. McMenemy	Queen's University Belfast, UK
Xiang Mei	The University of Leeds, UK
Cheol Moon	Gwangju University, South Korea
Veli Mumcu	Technical University of Yildiz, Turkey

Nin Pang	Auckland University of Technology, New Zealand
Jian-Xin Peng	Queen's University of Belfast, UK
Lui Piroddi	Technical University of Milan, Italy
Girij Prasad	University of Ulster, UK
Cent Leung	Victoria University of Technology, Australia
Jams Li	University of Birmingham, UK
Liang Li	University of Sheffield, UK
Hai Qi	University of Tennessee, USA
Wi Richert	University of Paderborn, Germany
Meh shafiei	Dalhousie University, Canada
Sa Sharma	University of Plymouth, UK
Dong Yue	Huazhong University of Science and Technology, China
YongSheng Ding	Donghua University, China
Yuezhi Zhou	Tsinghua University, China
Yongning Tang	Illinois State University, USA
Jun Cai	University of Manitoba, Canada
Sunil Maharaj Sentech	University of Pretoria, South Africa
Mei Yu	Simula Research Laboratory, Norway
Gui-Rong Xue	Shanghai Jiao Tong University, China
Zhichun Li	Northwestern University, China
Lisong Xu	University of Nebraska-Lincoln, USA
Wang Bin	Chinese Academy of Sciences, China
Yan Zhang	Simula Research Laboratory and University of Oslo, Norway
Ruichun Tang	Ocean University of China, China
Wenbin Jiang	Huazhong University of Science and Technology, China
Xingang Zhang	Nanyang Normal University, China
Qishi Wu	University of Memphis, USA
Jalel Ben-Othman	University of Versailles, France

Table of Contents

Internet Computing and Applications

A Novel Event Network Matching Algorithm	1
<i>Shan Jianfang and Liu Zongtian</i>	
Static Patterns Matching for High Speed Networks	15
<i>Kunpeng Jiang, Huifang Guo, Shengping Zhu, and Julong Lan</i>	
Information Propagation in Online Social Networks Based on User Behavior	23
<i>Niu Li and Han Xiaoting</i>	
Research on Operation and Management of Railway Transport of Dangerous Goods in Third-Party Logistics Enterprises	31
<i>Xin Li and Yue-fang Yang</i>	
A Password Authentication Scheme against Smart Card Security Breach	37
<i>Jing Shen and Yusong Du</i>	
A Vulnerability Attack Graph Generation Method Based on Scripts	45
<i>Bo Han, Qing Wang, Fajiang Yu, and Xianda Zhang</i>	
DFP-Growth: An Efficient Algorithm for Mining Frequent Patterns in Dynamic Database	51
<i>Zailani Abdullah, Tutut Herawan, A. Noraziah, and Mustafa Mat Deris</i>	
Analysis on Key Nodes Behavior for Complex Software Network	59
<i>Xizhe Zhang, Guolong Zhao, Tianyang Lv, Ying Yin, and Bin Zhang</i>	
Webpage Information Hiding Algorithm Based on Integration of Tags and Data	67
<i>Junling Ren and Li Zhang</i>	
Location Method of Underground Pipeline Monitoring Point Based on Cavity Detection	75
<i>Wei Zhu, Ping Sun, Ying nan Ma, Rui Song, Shi wei He, and Ke Hui Liu</i>	

Research on Purified Internet Environment for College Students 82
Qichun Zhong and Jinghong Hu

A Mobile-Certificate Security Method of Satellite-Earth Integration
 Networks 88
Qianmu Li, Qiugan Shi, Jun Hou, Yong Qi, and Hong Zhang

Multimedia Networking and Computing

Duration Modeling for Emotional Speech 98
Wen-Hsing Lai and Siou-Lin Wang

Research on Image Retrieval Based on Color and Shape Features 104
*Hongwei Zhao, Xiao Chen, Wei Huang, Pingping Liu, and
 Lingjiao Ma*

Existence and Simulations of Periodic Solution for Impulsive
 Predator-Prey System with Stage Structure for the Predator 112
Kaihua Wang, Wenxiang Zhang, and Zhanji Gui

Dynamics and Simulations of Multi-species Competition-Predator
 System with Impulsive 120
Yan Yan, Kaihua Wang, and Zhanji Gui

Research on the Distal Supervised Learning Model of Speech
 Inversion 128
Ying Chen and Shaobai Zhang

Low Power Pulse Width Modulation Design for Class D Audio
 Amplifier Systems 136
Ruei-Chang Chen, Shih-Fong Lee, and Yeong-Chau Kuo

A New Logic Method for Education Resource Software Guarantee 144
Guan Wei and Lv Yuanhai

Research on 3D Object Rounding Photography Systems and
 Technology 152
Zhenjie Hou, Junsheng Huang, and Jianhua Zhang

Facial Expression Feature Selection Based on Rough Set 159
Dong Li, Yantao Tian, Chuan Wan, and ShuaiShi Liu

Intelligent Computing and Applications

Concise Representations for State Spaces in Conformant Planning
 Tasks 167
Weisheng Li, Jiao Du, and Lifang Zhou

On Fast Division Algorithm for Polynomials Using Newton Iteration ... <i>Zhengjun Cao and Hanyue Cao</i>	175
On the Security of an Improved Password Authentication Scheme Based on ECC..... <i>Ding Wang, Chun-guang Ma, Lan Shi, and Yu-heng Wang</i>	181
Strategies to Develop Personal Knowledge Management Ability Based on M-Learning..... <i>Lin Hu</i>	189
Improved Detection Algorithm for MIMO Wireless Communication System Based on Chase Detector <i>Li Liu, Jinkuan Wang, Dongmei Yan, and Fulai Liu</i>	196
Improving Recommendation Performance through Ontology-Based Semantic Similarity <i>Mingxin Gan, Xue Dou, and Rui Jiang</i>	203
Formal Modeling and Model Checking Analysis of the Wishbone System-on-Chip Bus Protocol <i>Ricai Luo and Hua Tan</i>	211
A Time Synchronization Method for Wireless Sensor Networks <i>Chao Zou and Yueming Lu</i>	221
Real-Valued Negative Selection Algorithm with Variable-Sized Self Radius <i>Jinquan Zeng, Weiwen Tang, Caiming Liu, Jianbin Hu, and Lingxi Peng</i>	229
Scale Effect on Soil Attribute Prediction in a Complex Landscape Region <i>Zhenfu Wu, Yanfeng Zhao, Li Qi, and Jie Chen</i>	236
New Machine Learning Algorithm: Random Forest <i>Yanli Liu, Yourong Wang, and Jian Zhang</i>	246
Efficient Method of Formal Event Analysis <i>Ying Liu and Zongtian Liu</i>	253

Computational Statistics and Applications

Application of Monte Carlo Simulation in Reliability and Validity Evaluation of Two-Stage Cluster Sampling on Multinomial Sensitive Question..... <i>Qiao-qiao Du, Ge Gao, Zong-da Jin, Wei Li, and Xiang-yu Chen</i>	261
--	-----

New Lax-Friedrichs Scheme for Convective-Diffusion Equation	269
<i>Haixin Jiang and Wei Tong</i>	
Span of T-Colorings Multigraphs	277
<i>Juan Du</i>	
Dual-Scaled Method for the Rheology of Non-newtonian Boundary Layer and Its High Performance FEM	284
<i>Lei Hou, Hanling Li, Ming Zhang, Weijia Wang, Dezhi Lin, and Lin Qiu</i>	
Factor Model Averaging Quantile Regression and Simulation Study	291
<i>Zhimeng Sun</i>	
Analytical Studies and Experimental Examines for Flooding-Based Search Algorithms	299
<i>Hassan Barjini, Mohamed Othman, Hamidah Ibrahim, and Nur Izura Udzir</i>	
Innovative Study on the Multivariate Statistical Analysis Method of Chromatography Economy Analysis	307
<i>Shibing You, Yuan Hui, Xue Yu, and Lili Bao</i>	
Optimization of Lifting Points of Large-Span Steel Structure Based on Evolutionary Programming	315
<i>Xin Wang, Xu Lei, Xuyang Cao, Yang Zhou, and Shunde Gao</i>	
Modifying Feasible SQP Method for Inequality Constrained Optimization	323
<i>Zhijun Luo, Zhibin Zhu, and Guohua Chen</i>	
On the Solvable n -Lie Algebras	331
<i>Liu Jianbo, Zhang Yanyan, Men Yafeng, and Chen Wenying</i>	
Modules of Lie Algebra $G(A)$	337
<i>Zhang Yanyan, Liu Jianbo, Tao Wen, and Zhu Qin</i>	
Detection Performance Analysis of tests for Spread Targets in Compound-Gaussian Clutter	343
<i>Xiandong Meng, Zhiming He, Xiaowei Niu, and Ganzhong Feng</i>	
Rough Differential Equations in Rough Function Model	350
<i>Yun Wang, Xiaojing Xu, and Zhiqin Huang</i>	
 Cloud and Evolutionary Computing	
Cloud Computing Infrastructure and Application Study	358
<i>Ming Ye and ZeHui Qu</i>	

An Ant Colony Algorithm for Solving the Sky Luminance Model Parameters	365
<i>Ping Guo, Lin Zhu, Zhujin Liu, and Ying He</i>	
Tugboat Scheduling Problem Based on Trust-Based Ant Colony Optimization	373
<i>Su Wang, Min Zhu, Jun Zheng, and Kai Zheng</i>	
A Cloud Architecture with an Efficient Scheduling Technique	381
<i>Nawsher Khan, A. Noraziah, and Tutut Herawan</i>	
A Local Search Particle Swarm Optimization with Dual Species Conservation for Multimodal Optimization	389
<i>Dingcai Shen and Xuewen Xia</i>	
Cloud Computing: Analysis of Various Services	397
<i>Nawsher Khan, A. Noraziah, Tutut Herawan, and Mustafa Mat Deris</i>	
Quantum Ant Colony Algorithm Based on Bloch Coordinates	405
<i>Xiaofeng Chen, Xingyou Xia, and Ruiyun Yu</i>	
Energy Efficient VM Placement Heuristic Algorithms Comparison for Cloud with Multidimensional Resources	413
<i>Dailin Jiang, Peijie Huang, Piyuan Lin, and Jiacheng Jiang</i>	
A Python Based 4D Visualization Environment	421
<i>Lin Jing, Xipei Huang, Yiwen Zhong, Yin Wu, and Hui Zhang</i>	
Study of Trustworthiness Measurement and Kernel Modules Accessing Address Space of Any Process	429
<i>Ce Zhang, Gang Cui, Bin Jin, and Liang Wang</i>	
Human Resource Management System Based on Factory Method Design Pattern	437
<i>Xing Xu, Hao Hu, Na Hu, Lin Xiao, and Weiqin Ying</i>	
Ant Colony Optimization with Multi-Agent Evolution for Detecting Functional Modules in Protein-Protein Interaction Networks	445
<i>Junzhong Ji, Zhijun Liu, Aidong Zhang, Lang Jiao, and Chunnian Liu</i>	
Research on Genetic Segmentation and Recognition Algorithms	454
<i>Zhenjie Hou and Jianhua Zhang</i>	

Computer Engineering and Applications

Blog-Based Distributed Computation: Implementation of Software Verification System	461
<i>Takayuki Sasajima and Shin-ya Nishizaki</i>	

Model Transformation Method for Compensation Events and Tasks from Business Process Model to Flowchart	468
<i>Jian Deng, Bo Chen, and Jiazhi Zeng</i>	
Towards Efficient Replication of Documents in Chord: Case (r,s) Erasure Codes	477
<i>Rafał Kapelko</i>	
Path Planning for Crawler Crane Using RRT*	484
<i>Yuanshan Lin, Di Wu, Xin Wang, Xiukun Wang, and Shunde Gao</i>	
Study on Data Preprocessing for Daylight Climate Data	492
<i>Ping Guo, Shuai-Shuai Chen, and Ying He</i>	
Scalable Technique to Discover Items Support from Trie Data Structure	500
<i>A. Noraziah, Zailani Abdullah, Tutut Herawan, and Mustafa Mat Deris</i>	
On Soft Partition Attribute Selection	508
<i>Rabiei Mamat, Tutut Herawan, Noraziah Ahmad, and Mustafa Mat Deris</i>	
Effect of Vocabulary Preparation on Students Vocabulary and Listening Comprehension	516
<i>Lijun Li, Kaida He, and Qiudong He</i>	
3D Parametric Design on Trough Type Liquid Distributor Based on AutoCAD VBA	524
<i>Pengfei Zhang and LuoJia Wan</i>	
A Novel Differential Evolution Algorithm with Adaptive of Population Topology	531
<i>Yu Sun, Yuanxiang Li, Gang Liu, and Jun Liu</i>	
Empirical Study on the Relationship between Money Supply and Stock Market in Europe	539
<i>Yijun Li</i>	
Evaluation of Agricultural Information Service System	545
<i>Yanxia Wang and Xingjie Hui</i>	
Six-Mode Truncation and Chaotic Characteristics of Atmospheric Convection System	553
<i>Li Zhen</i>	
Security Uniform Office Format Specification and API Design Based on the Java Platform	560
<i>Ying Cai, Ning Li, and Chengxia Liu</i>	

Bioinformatics Analysis of the Complete Nucleotide Sequence of Duck Plague Virus UL22 Gene	569
<i>Li-Sha Yang, An-Chun Cheng, Ming-Shu Wang, De-Kang Zhu, Shun Chen, Ren-Yong Jia, and Xiao-Yue Chen</i>	

Knowledge Management and Applications

Evidence Conflict Analysis Approach to Obtain an Optimal Feature Set for Bayesian Tutoring Systems	576
<i>Choo-Yee Ting, Kok-Chin Khor, and Yok-Cheng Sam</i>	

Binary Vote Assignment Grid Quorum for Managing Fragmented Database	584
<i>A. Noraziah, Ainul Azila Che Fauzi, Noriyani Mohd Zin, and Tutut Herawan</i>	

WLAR-Viz: Weighted Least Association Rules Visualization	592
<i>A. Noraziah, Zailani Abdullah, Tutut Herawan, and Mustafa Mat Deris</i>	

Ontology-Based Genes Similarity Calculation with TF-IDF	600
<i>Yue Huang, Mingxin Gan, and Rui Jiang</i>	

Quantitative Study of Oilfield Casing Damage	608
<i>Deng Rui, Zhang Liang, and Guo Haimin</i>	

Complete SAT Solver Based on Set Theory	616
<i>Wensheng Guo, Guowu Yang, Qianqi Le, and William N.N. Hung</i>	

Application of XML Data Mining in GUI Run-Time State Clustering . . .	624
<i>Jing Feng and Tingjie ShangGuan</i>	

Strong Reduction for Typed Lambda Calculus with First-Class Environments	632
<i>Shin-ya Nishizaki and Mizuki Fujii</i>	

Reliable NoC Mapping Based on Scatter Search	640
<i>Qianqi Le, Guowu Yang, William N.N. Hung, and Wensheng Guo</i>	

Towards the Applied Hybrid Model in Decision Making: Support the Early Diagnosis of Type 2 Diabetes	648
<i>Andrea Carvalho Menezes, Placido Rogerio Pinheiro, Mirian Caliope Dantas Pinheiro, and Tarcísio Pequeno Cavalcante</i>	

Replacement Study of the Retention Time in Chromatography Economic Analysis	656
<i>Ping Shen and Shibing You</i>	

Study on Frustration Tolerance and Training Method of College Students	663
<i>Naisheng Wang</i>	
On Statistical Analysis and Management Countermeasures of Occupation Burnout for College Teachers	669
<i>Youcai Xue and Jianhui Pan</i>	
Development of a RFID Multi-point Positioning and Attendance System Based on Data Comparison Algorithm	677
<i>Wenyu Zhao, Jun Gao, Xiaotian Liu, and Yaping Wu</i>	

Communication Technology and Applications

High Throughput Constraint Repetition for Regular Expression Matching Algorithm	684
<i>Kunpeng Jiang, Julong Lan, and Youjun Bu</i>	
Interference Mitigation Based on Multiple SINR Thresholds in 60GHz Wireless Networks	692
<i>Weixia Zou, Fang Zhang, Guanglong Du, and Bin Li</i>	
Analysis of the Coupling Action in Nonlinear Harmonic Vibration Synchronization System	700
<i>Xiaohao Li and Zhenwei Zhang</i>	
Performance Analysis of Hierarchical Modulation with Higher Spectrum Efficiency for a Higher Data Rate T-DMB System	707
<i>Linlin Dong, Lixin Sun, Xiaoming Jiang, and Na Zhu</i>	
Research on the Controlling Technique of DFBLD in the Spectrum Absorption Optical Fiber Gas Detecting System	715
<i>Shutao Wang, Pengwei Zhang, and Xiaoqing Shao</i>	
An Efficient and Provable Secure PAKE Scheme with Robust Anonymity	722
<i>Cong Liu and Chuan-gui Ma</i>	
Study on QoS of Video Communication over VANET	730
<i>Shouzhi Xu, Pengfei Guo, Bo Xu, and Huan Zhou</i>	
Phase Noise Estimation and Mitigation for Cognitive OFDM Systems	739
<i>Yuan Jing, Haoyu Li, Xiaofeng Yang, Li Ma, Ji Ma, and Bin Niu</i>	
An Effective Multipath Approach to Reducing Congestion in WSNs	746
<i>Laomo Zhang, Ying Ma, and Guodong Wang</i>	

Immunity-Based Gravitational Search Algorithm	754
<i>Yu Zhang, Yana Li, Feng Xia, and Ziqiang Luo</i>	
A Clustering Method Based on Time Heat Map in Mobile Social Network	762
<i>Wang Ye, Wang Jian, and Yuan Jian</i>	
TNC-eSA: An Enhanced Security Access Solution to Office Networks . . .	770
<i>Jun Ma, Yuan-bo Guo, and Jinbo Xiong</i>	
Author Index	779

A Novel Event Network Matching Algorithm

Shan Jianfang and Liu Zongtian

School of Computer Engineering and Science, Shanghai University, Yanchang Road 149,
Zhabei, Shanghai, China
sjfshan@163.com, ztliu@shu.edu.cn

Abstract. Event network is a new semantic-based and event-oriented text representation model, the operation on event network is a good form of semantic computation, which can provide support for text semantic information processing. The paper proposes a new three-step matching algorithm for event network: event matching based on maximum similarity priority, relation matching based on isotropic-relational-distance matrix, and event network matching by integrating event matching and relation matching. The paper's experimental results show that the method is feasible and reasonable.

Keywords: Event, Relation, Maximum similarity priority, Event network, Isotropic-matrix.

1 Introduction

Text representation is an important issue in Natural Language Processing, such as information retrieval, text classification and so on. In recent years, there is a tendency to use richer text representations than just keywords and concepts. Under the circumstances, it is necessary to investigate the appropriate methods for comparison of two texts in any of these text representations. In the paper, text is represented with event network and a new method for matching two event networks is proposed.

Event network is a new text representation model, could be regarded as graph but with many extra features. Graph matching theory is of good use for event network matching. Most of graph matching come from graph theory and information retrieval. Matching theory is a central part of graph theory [1], Graph matching methods can be divided into two groups: the biggest common substructure and distance transform. In information retrieval, comparison of conceptual graph at its core has been studied and applied widely [2, 3]. Jonathan Poole [4] defined three similarities for conceptual graph matching: surface, structure and thematic similarity, surface similarity is similarity of matching conception, structure similarity is similarity of matching relationship, and thematic similarity takes into account pattern of conception and relationship. On the basis, H.P.Zhu [5] defined conception similarity, relationship similarity and graph similarity for conceptual graph semantic similarity. However, there is a significant difference between event network and graph and conceptual graph both in structure and meaning. The existing matching methods can not be used directly for event network matching. Therefore, it is necessary to study new method for matching event network.

The paper proposes a new method for event network matching: event matching based on maximum similarity priority, relation matching based on isotropic-relational-distance matrix, or isotropic-distance matrix for short, the third and final step is event network matching by integrating event matching and relation matching. The corresponding matching is calculated in each step, and choose the desirable weighting factors to decide emphasis on event or relation a little more. Our experimental results show that the method is feasible and reasonable.

2 Preliminary Work

2.1 Event and Event Similarity

Event is the basic unit of human cognition, it originated from cognitive science, often appears in the texts of philosophy, cognitive science, linguistics and Artificial Intelligence [6]. Different applications define event in very different ways[7] most of them emphasize two categories of event attributes, action (verb or action- noun) and action characteristics(participant, location, time, etc.). These attributes are called event element or element for short.

Event identification and event clustering are essential for event-based information query,automatic question answering and automatic summarization, etc. All these are directly bound up with event similarity. Shan Jianfang introduced how to calculate the similarity between two events ,which combines syntax, semantic, word sequences and time relation according to the feature of different event elements. The similarity between corresponding event elements is calculated firstly, then the similarity between the two events is calculated by the weighted summation of elements similarities. Given two events e_1 and e_2 , that can be expressed as the following: $e_1=(H_{11},H_{12},H_{13},\dots,H_{1n})$, $e_2=(H_{21},H_{22},H_{23},\dots,H_{2n})$, the formula for calculating similarity between e_1 and e_2 is $Sim(e_1,e_2)=\sum_{i=1}^n w_i Sim(H_{1i},H_{2i})$,where H_{1i} and H_{2i} are i-th event element of e_1 and e_2 , such as participant, location, time,etc., w_i is weighting factor of i-th event elements, subject to $\sum_{i=1}^n w_i =1$.

2.2 Event Network

Event network is a novel semantic-based and event-oriented text representation model, can preserve semantic information of text, show the relations between events, and also reflect the importance of events, dynamic behaviors of events. The operations on event network are good form of semantic computation, can provide support for text semantic information processing.

Shan Jianfang introduced how to represent text by event network,the nodes represent events, and the edges represent the correlations between events, event-relation diagrams

is constructed by extracting events and event relations from text, this is co-occurrence event network, or event network, it indicates that these events co-occur in the text.

Definition 1. Event network

An event network $EN = (E, R)$ is a two-tuples that meet the following conditions:

- (1) E is nonempty vertex set, is called event set, express as $E = \{e\}$.
- (2) R is edge (event relation) set: including taxonomic relationship and non-taxonomic relationship, $R \subseteq E \times E$. If there were a relation r between e_1 and e_2 , then the two events are linked by an edge labeled with “ r ”, express as $(e_1, e_2)_r$ or $r(e_1, e_2)$.

3 Event Set Matching

Text T_1 is represented as event network $EN_1 = (E_1, R_1)$, and T_2 is represented as $EN_2 = (E_2, R_2)$, where $E_1 = \{e_{11}, e_{12}, \dots, e_{1m}\}$ and $E_2 = \{e_{21}, e_{22}, \dots, e_{2n}\}$. Event similarity $Sim(e_{1i}, e_{2j}) (1 \leq i \leq m, 1 \leq j \leq n)$ is calculated according to the method mentioned in section 2.1. The results were tabulated as shown in Table 1.

Table 1. Event set similarity matrix $Sim(E_1, E_2)$

e_{1i}	$sim(e_{1i}, e_{2j})$	e_{2j}	e_{21}	e_{22}	...	e_{2j}	...	e_{2n}
e_{11}			-	-	...	-	...	-
e_{12}			-	-	...	-	...	-
...		
e_{1i}			-	-	...	-	...	-
...		
e_{1m}			-	-	...	-	...	-

Given two event sets, $E_1 = \{e_{11}, e_{12}, \dots, e_{1m}\}$, $E_2 = \{e_{21}, e_{22}, \dots, e_{2n}\}$, $|E_1| = m$, $|E_2| = n$.

Definition 2. Matching-event

Given event set E_1 and E_2 , according to some a matching scheme, for each event e_{1i} in E_1 ($e_{1i} \in E_1$), identify the only event e_{2j} in E_2 that corresponding to it, then e_{1i} and e_{2j} are called matching-event or matching-event-pair, express as $ms(e_{1i}, e_{2j})$ or (e_{1i}, e_{2j}) .

Given event set similarity matrix $Sim(E_1, E_2)$, the algorithm for event matching based on maximum similarity priority is as followed:

Algorithm 1. MSPEM. Event matching based on maximum similarity priority

MSPEM(maximum similarity priority event matching):
 Input: $Sim(E_1, E_2)$ // event set similarity matrix, $|E_1| = m, |E_2| = n$
 Output:
 $Match_E(Em_1, Em_2, Sim(Em_1, Em_2)) =$
 $\{match_e(e_{1i}, e_{2j}, sim(e_{1i}, e_{2j})) \mid i \in [1, m], j \in [1, n]\}$
 // $\min(m, n)$ triples set of matching-event and the similarity // Em_1 is matched
 event set, $Em_1 \subseteq E_1$,
 // Em_2 is matched event set, $Em_2 \subseteq E_2$,
 $mc = \min(|E_1|, |E_2|) = \min(m, n)$ // the number of matching-event-pair
 begin
 Step 1. finding the current maximal $sim(e_{1i}, e_{2j})$ in $Sim(E_1, E_2)$,
 $match_e(e_{1i}, e_{2j}, sim(e_{1i}, e_{2j}))$ is a matching-event triple,
 adding the triple to $Match_E(Em_1, Em_2, Sim(Em_1, Em_2))$.
 Step 2. removing the row and column that $sim(e_{1i}, e_{2j})$ located in $Sim(E_1, E_2)$
 (i.e. the deleted-row is e_{1i} located, the deleted-column is e_{2j} located)
 Step 3. repeating step 1 and 2 until matrix $Sim(E_1, E_2)$ is empty, and there are
 mc matching-event-pairs.
 end

There is a one-to-one event matching in E_1 and E_2 only when the size of E_1 is equal to E_2 , otherwise there will be $\|E_1| - |E_2|\|$ unmatched events in the larger set, as follows:

If $|E_1| = |E_2|$, then $Em_1 = E_1, Em_2 = E_2$

If $|E_1| > |E_2|$, then $Em_1 \subset E_1, Em_2 = E_2$

If $|E_1| < |E_2|$, then $Em_1 = E_1, Em_2 \subset E_2$.

By algorithm 1, matching-event triples set $Match_E(Em_1, Em_2, Sim(Em_1, Em_2))$ is obtained, the formula for calculating event set matching M_E is defined as follows:

$$\begin{aligned}
 M(E_1 \rightarrow E_2) &= \frac{\sum_{i=1}^{mc} sim(e_{1i}, e_{2j})}{|E_1|} = \frac{\sum_{i=1}^{mc} sim(e_{1i}, e_{2j})}{m} \\
 M(E_2 \rightarrow E_1) &= \frac{\sum_{i=1}^{mc} sim(e_{1i}, e_{2j})}{|E_2|} = \frac{\sum_{i=1}^{mc} sim(e_{1i}, e_{2j})}{n} \\
 M_E &= \frac{M(E_1 \rightarrow E_2) + M(E_2 \rightarrow E_1)}{2}
 \end{aligned} \tag{1}$$

where $e_{1i} \in Em_1, e_{2j} \in Em_2, (e_{1i}, e_{2j})$ is matching-event-pair according to algorithm 1, $match_e(e_{1i}, e_{2j}, sim(e_{1i}, e_{2j})) \in Match_E(Em_1, Em_2, Sim(Em_1, Em_2))$, $\sum_{i=1}^{mc} sim(e_{1i}, e_{2j})$ is summation of similarities of matching-event. When the size of E_1 is equal to E_2 , formula 1 can be simplified as follow:

$$M_E = \frac{\sum_{i=1}^{mc} sim(e_{1i}, e_{2j})}{m} = \frac{\sum_{i=1}^{mc} sim(e_{1i}, e_{2j})}{n} \tag{2}$$

4 Relation Matching

4.1 Relation Similarity Based on Distance

There is neighboring relation; share event-element relation cause-effect relation and following relation are considered in the paper’s event network. These relations can be divided into two categories, as follows:

Semantic relationship: including cause-effect relation and following relation, the distance of these relationships is 1.

Syntactic relationship: including neighboring relation, share event-element relation, the distance of these relationships is 2.

The distance is smaller, the strength of relation is stronger, the distance and the strength of relation are reciprocals. The strength of semantic relationship is greater than syntactic relationship. If there were several relations between two events, then only the strongest relation is considered.

Based on the above, the relation is expressed as its distance. The distance-based relation similarity as shown in the table 2, where a is distance of r_1 , b is distance of r_2 . Semantic relationship is directed, the distance of the reversal of a directed relation is negative of it. “∞” denotes two events with no relation, “0” denotes diagonal elements of event network, i.e., the distance of an event and itself is zero.

Table 2. Relation similarity matrix

$a \quad sim(r_1, r_2) \quad b$	1	-1	2	∞	0
1	1	0.5	0.5	0	0.5
-1	0.5	1	0.5	0	0.5
2	0.5	0.5	1	0	0.5
∞	0	0	0	1	0
0	-	-	-	-	1

4.2 Relation Matching

This paper presents an algorithm for relation matching based on isotropic- -relational-distance matrix.

Definition 3. Matching relation:

Given event networks $EN_1 = (E_1, R_1)$ and $EN_2 = (E_2, R_2)$, event matching triples set $Match_E(Em_1, Em_2, Sim(Em_1, Em_2))$. If $match_r[(e_{1i}, e_{1k}), (e_{2j}, e_{2l})]$ is called matching relation, it should meet the following constraints:

- (1) $e_{1i} \in Em_1, e_{1k} \in Em_1, (e_{1i}, e_{1k}) \in R_1$
- (2) $e_{2j} \in Em_2, e_{2l} \in Em_2, (e_{2j}, e_{2l}) \in R_2$
- (3) $match_E(e_{1i}, e_{2j}, sim(e_{1i}, e_{2j})) \in Match_E(Em_1, Em_2, Sim(Em_1, Em_2))$

Where $r_1 \in R_1, r_2 \in R_2$, and $match_r[(e_{1i}, e_{1k}), (e_{2j}, e_{2l})]$ can also be expressed as $match_r[r_1, r_2]$, as shown in the Fig. 1:

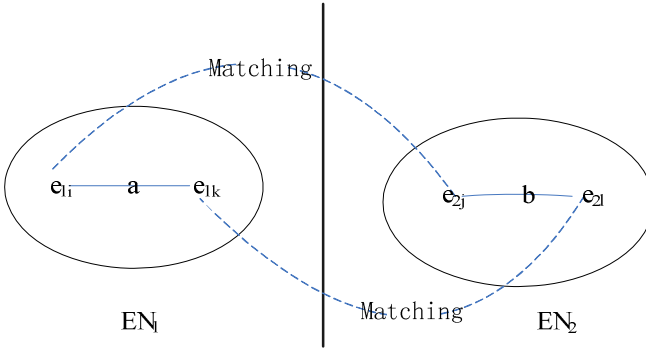


Fig. 1. Matching relation

The formula for calculating relation matching is defined as follows:

$$M_r = \frac{sim(e_{1i}, e_{2j}) + sim(e_{1k}, e_{2l})}{2} * [\alpha + \beta * sim(a, b)] \quad (3)$$

Where α and β are user-specified weighting adjustment factors, subject to $0 < \alpha, \beta < 1, \alpha + \beta = 1$.

Definition 4. Distance-based matrix of event network

Given event network $EN = (E, R)$, the distance-based matrix A is a $|E|$ order square matrix that saves relation distance of event relation, and has the following properties:

$$A(i, j) = \begin{cases} 1, (e_i, e_j) \in R_{semantic} \\ -1, (e_j, e_i) \in R_{semantic} \\ 2, (e_i, e_j) \in R_{syntax} \\ 0, (i = j) \\ \infty, (e_i, e_j) \notin R \end{cases},$$

Where $R_{semantic} \subseteq R$, is directed semantic relationship set, and $R_{syntax} \subseteq R$, is undirected syntactic relationship set.

Definition 5. Isotropic-element of matrix:

Element a_{rc} of matrix M_A and element $b_{r_1c_1}$ of matrix M_B are called isotropic-element, if and only if $r = r_1, c = c_1$, i.e. isotropic-element is a pair of two elements that are in two different matrices and has the same position in their respective matrices.

Definition 6. ($\xrightarrow{SetToVector}$): Set vectoring operation

According to some rule, convert each element of a set to a unique element of a vector.

For example:

Given set $s = \{a_1, a_4, a_8, a_6, a_3, a_2\}$, please convert s to vector according to set element subscript increasing:

$$\text{row vector } V_1 : s(\xrightarrow{SetToVector}) = (a_1, a_2, a_3, a_4, a_6, a_8)$$

$$\text{column vector } V_1^T : s(\xrightarrow{SetToVector}) = (a_1, a_2, a_3, a_4, a_6, a_8)^T$$

$$\text{and } |s| = |V_1| = |V_1^T|$$

Definition 7. \oplus : orderly union operation on vector

$$\text{row vector } V_1 = (a_1, a_2, \dots, a_m), \text{column vector } V_1^T = (a_1, a_2, \dots, a_m)^T,$$

$$\text{row vector } V_2 = (b_1, b_2, \dots, b_n), \text{column vector } V_2^T = (b_1, b_2, \dots, b_n)^T,$$

$$\text{null row vector } V_\phi = \phi, \text{null column vector } V_\phi^T = \phi,$$

\oplus : orderly union operation on vector with the definition is as follows:

$$V_1 \oplus V_2 = (a_1, a_2, \dots, a_m) \oplus (b_1, b_2, \dots, b_n) = (a_1, a_2, \dots, a_m, b_1, b_2, \dots, b_n)$$

$$V_1^T \oplus V_2^T = (a_1, a_2, \dots, a_m)^T \oplus (b_1, b_2, \dots, b_n)^T = (a_1, a_2, \dots, a_m, b_1, b_2, \dots, b_n)^T$$

$$(V_1 \oplus V_2)^T = V_1^T \oplus V_2^T$$

$$V_1 \oplus V_\phi = V_\phi \oplus V_1 = V_1$$

$$V_1^T \oplus V_\phi^T = V_\phi^T \oplus V_1^T = V_1^T$$

The algorithm for relation matching based on isotropic-distance-matrix, or isotropic- matrix, is as follows:

Algorithm 2. Isotropic-matrix-based relation matching

RMBIM(relation matching based on isotropic- matrix)

Variable: (r, c) // The position of an element in the matrix, r is row number and c is column number.

Function:

position (EN,element): return the position of “element” in matrix EN

elementIso(EN,r,c): return a element that is located at (r,c) in matrix EN

relation(EN,r,c): return a two-tuples, the former is row event and the latter is column event corresponds to position (r,c) in matrix EN, i.e. a relation.

Input: EN_1, EN_2 // distance-based matrix

$Match_E(Em_1, Em_2, Sim(Em_1, Em_2))$ // event-matching-triples set

Output:

$Match_R(Rm_1, Rm_2, M_R(Rm_1, Rm_2)) =$

$$\{match_r(rel_1, rel_2, M_r(rel_1, rel_2)) \mid rel_1 \in R_1, rel_2 \in R_2\}$$

/*triples set of matching relation and the matching degree, and $Rm_1 \subset R_1$, $Rm_2 \subset R_2$, $M_r(rel_1, rel_2)$ is matching degree between matching relation rel_1 and rel_2 , calculated by formula 3.*//

begin

Step 1. Set vectoring operation on event-matching-triples set:

letting $Match_E(Em_1, Em_2, Sim(Em_1, Em_2))$ orderly by subscript i increasing of e_{1i} ($e_{1i} \in Em_1$).

$$Vector(EN_1) = Em_1 \left(\xrightarrow{SetToVector} \right)$$

// event row vector of EN_1 , orderly by subscript i increasing of e_{1i} .

$$Vector(EN_1)^T = (Em_1 \left(\xrightarrow{SetToVector} \right))^T$$

// events column vector of EN_1 , orderly by subscript i increasing of e_{1i}

$$EN_1 = EN_1(Vector)$$

// Saving rows and columns of EN_1 that event vector elements are located, i.e. unmatched events are deleted from EN_1

Step 2. Calculating isotropic-matrix of EN_1 :

(Transformation on matrix EN_2 according to $Match_E(Em_1, Em_2)$, let isotropic-element of EN_1 and EN_2 is matching-event-pair by interchanging row-row, column-column of matrix EN_2 .)

$Vector(EN_{2*}) = \phi$ // Initialization of event vector of matrix EN_{2*} to null.

for each $match_e(e_{1i}, e_{2j}) \in Match_E(Em_1, Em_2)$

$$Vector(EN_{2*}) = Vector(EN_{2*}) \oplus (e_{2j})$$

transforming matrix EN_2 by elementary matrix transformation according to

$Vector(EN_{2*})$ and $Vector(EN_{2*})^T$ so that position of each element is consistent with the order of event in the event vector, after the transformation, EN_2 convert to EN_{2*} .

Step 3. matching based on isotropic-matrix:

(at off-diagonal, comparing element of EN_1 and its isotropic-element in EN_{2*} element-by-element, the isotropic-element pair denotes a matching relation.)

for each element in EN_1

if (element!=0)

$$(r, c) \leftarrow position(EN_1, element)$$

$$eleIso \leftarrow elementIso(EN_{2*}, r, c)$$

```

rel1(e1r, e1c) ← relation(EN1, r, c)
rel2(e2*r, e2*c) ← relation(EN2*, r, c)
// rel1 and rel2 are matching relation
// (e1r, e2*r) is matching-event-pair, (e1c, e2*c) is also.
Mr(rel1, rel2) =
    
$$\frac{sim(e_1^r, e_{2*}^r) + sim(e_1^c, e_{2*}^c)}{2} * [\alpha + \beta * sim(element, eleIso)]$$

// Calculating matching of relation according to formula 3
end

```

5 Event Network Matching

The last step is to integrate event matching and relation matching to calculate matching degree M_{EN} of event network EN_1 and EN_2 . The paper proposes two methods: averaging and weighted sum, as follows:

Averaging: All possible relation matching are calculated after algorithm 2 halts, matching degree M_{EN} is average of them.

Weighted sum:

Event set matching M_E is calculated according to formula 1. The formula for calculating relation set matching is defined as follows:

$$M_R(EN_1, EN_{2*}) = \frac{\sum_{i,j} sim(a,b)}{\sum_{i,j} \max[sim(a,a), sim(b,b)]} \quad (4)$$

Where a is off-diagonal element in distance-based matrix of EN_1 , b is off-diagonal element in distance-based matrix of EN_{2*} , a and b are isotropic-element each other, and EN_1 and EN_{2*} are isotropic-matrix each other.

The formula for calculating event network matching is defined as follows:

$$M_{EN}(EN_1, EN_2) = M_E * (\alpha + \beta * M_R) \quad (5)$$

Where α and β are user-specified weighting adjustment factors, subject to $0 < \alpha, \beta < 1, \alpha + \beta = 1$. If there is a very large difference in event network structure then relation set matching would approach to zero, and the general event network matching depends on event matching. In this case, the general event network matching is a fraction of the event matching, and the factor α indicates the value of the fraction. We can also specify factor α and β to decide emphasis on event or relation a little more, for example, if $\alpha > \beta$ then event matching is emphasized more than relation matching, and vice versa.

6 An Example

To illustrate the paper's method clearly, we show an example on a text titled "5 passengers were killed and more than 20 injured as buses collided in Jingdezhen, Jiangxi Province" and a text titled "10 passengers were killed and 2 serious injured as two buses collided in Liuzhou, Guangxi". Five unique events are extracted from the first text, ten events are extracted from the second text, put the same events together and six unique events are retained. Events lists are shown as Table 3 and Table 4.

Table 3. Events list of the first text

e_{i_i} : event ID	Event representation: verb(other event elements)
e_{11}	Collide(bus, bus, yesterday morning ,Jingdezhen)
e_{12}	Ran into(bush)
e_{13}	Die(5 persons , at the scene)
e_{14}	Seriously injured (2 bus drivers)
e_{15}	Injure(8 persons)

Table 4. Events list of the second text

e_{2_j} :event ID	Event representation: verb(other event elements)
e_{21}	Collide(car , coach)
e_{22}	Die(10 persons)
e_{23}	Seriously injured (2 men)
e_{24}	Traffic accident(Liuzhou City, Guangxi Province, early this morning)
e_{25}	Cross(car , coach, 323 National Highway, 6 a.m.)
e_{26}	Collide()

The distanced-based matrices of event network EN_1 and EN_2 of the two texts are shown as follows:

$$EN_1 = \begin{bmatrix} E_1 & e_{11} & e_{12} & e_{13} & e_{14} & e_{15} \\ e_{11} & 0 & 2 & 1 & 1 & 1 \\ e_{12} & 2 & 0 & 2 & \infty & \infty \\ e_{13} & -1 & 2 & 0 & 2 & \infty \\ e_{14} & -1 & \infty & 2 & 0 & 2 \\ e_{15} & -1 & \infty & \infty & 2 & 0 \end{bmatrix} \quad EN_2 = \begin{bmatrix} E_2 & e_{21} & e_{22} & e_{23} & e_{24} & e_{25} & e_{26} \\ e_{21} & 0 & 1 & 1 & \infty & 2 & 2 \\ e_{22} & -1 & 0 & 2 & -1 & \infty & -1 \\ e_{23} & -1 & 2 & 0 & -1 & 2 & -1 \\ e_{24} & \infty & 1 & 1 & 0 & 2 & \infty \\ e_{25} & 2 & \infty & 2 & 2 & 0 & 2 \\ e_{26} & 2 & 1 & 1 & \infty & 2 & 0 \end{bmatrix}$$

6.1 Event Matching

For each event in E_1 ($e_{i_i} \in E_1, |E_1|=5$), and each event E_2 ($e_{2_j} \in E_2, |E_2|=6$), event similarity $sim(e_{i_i}, e_{2_j})$ is calculated event-by-event according to the method introduced in section 3, triples set of matching-event and the event similarity is shown as follows:

$$Match_E(E_1, E_2, sim(E_1, E_2)) = \{match_e(e_{11}, e_{21}, 0.833), match_e(e_{14}, e_{23}, 0.8), \\ match_e(e_{13}, e_{22}, 0.75), match_e(e_{12}, e_{26}, 0.5), match_e(e_{15}, e_{24}, 0.07)\}$$

Event set matching is calculated according to formula 1:

$$M_E = \frac{\sum_{i=1}^5 sim(e_{1i}, e_{2j})}{5} = \frac{0.833 + 0.5 + 0.75 + 0.8 + 0.07}{5} = 0.596 \\ match_e(e_{1i}, e_{2j}, sim(e_{1i}, e_{2j})) \in Match_E(E_1, E_2, sim(E_1, E_2))$$

Matching-event set is vectored orderly by subscript i increasing of e_{1i} , the result as follows: $((e_{11}, e_{21}), (e_{12}, e_{26}), (e_{13}, e_{22}), (e_{14}, e_{23}), (e_{15}, e_{24}))$.

6.2 Relation Matching

Matrix EN_2 is transformed to EN_{2^*} according to vectoring-after matching-event list $((e_{11}, e_{21}), (e_{12}, e_{26}), (e_{13}, e_{22}), (e_{14}, e_{23}), (e_{15}, e_{24}))$ in section 6.1, and isotropic-matrix of EN_1 is worked out, and EN_{2^*} is shown as follows (e_{25} is unmatched event, it is omitted in EN_{2^*}):

$$EN_{2^*} = \begin{bmatrix} E_2 & e_{21} & e_{26} & e_{22} & e_{23} & e_{24} \\ e_{21} & 0 & 2 & 1 & 1 & \infty \\ e_{26} & 2 & 0 & 1 & 1 & \infty \\ e_{22} & -1 & -1 & 0 & 2 & -1 \\ e_{23} & -1 & -1 & 2 & 0 & -1 \\ e_{24} & \infty & \infty & 1 & 1 & 0 \end{bmatrix}$$

Comparing off-diagonal isotropic-element of matrix EN_1 and EN_{2^*} according to the step 2 of algorithm 2, the process can be diagramed as fitting EN_1 and EN_{2^*} , as shown in Fig. 2. Blue-shade-marked figures are the similarity of matching-event-pair, and the two events are above or left the figure. Each two-tuples represents a pair of isotropic-elements, the former is from EN_1 , and the latter is from EN_{2^*} , the former is a and the latter is b in formula 3, i.e. the former is *element* and the latter is *eleIso* in algorithm 2. Thus, each off-diagonal two-tuples denotes a matching relation, the data for calculating relation matching are the two elements of the two-tuples, and the event pair and similarity above or left the two-tuples. For example, the matching relation is denoted by two-tuples (2, 1) marked with underline, as shown in Fig.3, its matching is calculated as follows:

$$M_r = \frac{sim(e_{12}, e_{26}) + sim(e_{13}, e_{22})}{2} * [\alpha + \beta * sim(a, b)] \quad (\alpha = 0.5, \beta = 0.5) \\ = \frac{0.5 + 0.75}{2} * [0.5 + 0.5 * sim(2, 1)] = 0.46875$$

Its triple is $match_r((e_{12}, e_{13}), (e_{26}, e_{22}), 0.46875)$.

		EN ₁	e ₁₁	e ₁₂	e ₁₃	e ₁₄	e ₁₅
		EN _{2*}	e ₂₁	e ₂₆	e ₂₂	e ₂₃	e ₂₄
EN ₁	EN _{2*}		0.833	0.5	0.75	0.8	0.07
e ₁₁	e ₂₁	0.833	(0,0)	(2,2)	(1,1)	(1,1)	(1,∞)
e ₁₂	e ₂₆	0.5	(2,2)	(0,0)	(2,1)	(∞,1)	(∞,∞)
e ₁₃	e ₂₂	0.75	(-1,-1)	(2,-1)	(0,0)	(0,2)	(∞,-1)
e ₁₄	e ₂₃	0.8	(-1,-1)	(∞,-1)	(2,2)	(0,0)	(2,-1)
e ₁₅	e ₂₄	0.07	(-1,∞)	(∞,∞)	(∞,1)	(2,1)	(0,0)

Fig. 2. Fitting EN₁ and EN_{2*} matching event similarity

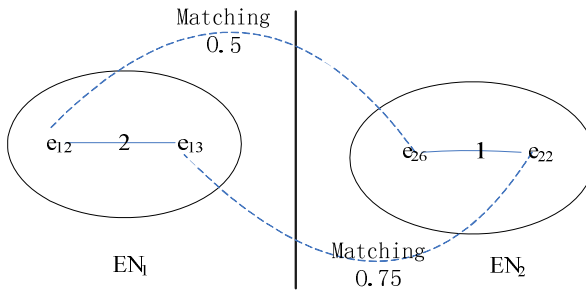


Fig. 3. An example of matching relation

6.3 Event Network Matching

The last step is to calculate matching M_{EN} between event network EN_1 and EN_2 by integrating event matching and relation matching, and $\alpha = \beta = 0.5$ in the example.

The first method: averaging

$M_{EN} = 0.47$, it is average of all relation matching calculated in section 6.2

The second method: Weighted sum

Event set matching is calculated in section 6.1: $M_E = 0.596$

Relation set matching is calculated according to formula 4: $M_R = 0.5$

Event network matching is calculated according to formula 5:

$$M_{EN} = M_E * (\alpha + \beta * M_R) = 0.596 * (0.5 + 0.5 * 0.5) = 0.447$$

In the example, there is little difference between the two methods, and both are within the acceptable limits.

7 Evaluating the Algorithm

The texts used in the experiments are collected from the website, can be divided into five categories: earthquake, fire, traffic accidents, terrorist attack and food poisoning.

7.1 Event Matching

Event matching is the first and most critical step throughout the matching process, it determines the relation matching strategy. Thus it is necessary to evaluate event matching specially.

Method for determining event matching: for any given two texts and their events list, matching-event-pairs are determined as much as possible by hand, and compare the automated results with manual results, the recall evaluation formula as follows:

$$\text{Recall} = \frac{|EM \cap AM|}{|AM|} \quad (6)$$

Where EM is the paper's matching-event-pairs set, AM is manual matching-event-pairs set. Table 5 shows evaluation results on our corpus.

Table 5. The evaluation for event matching

	Earthquake	Fire	Traffic-accidents	Terrorist-attack	Food-poisoning
Recall	94%	80%	90%	86%	85%

7.2 Event Network Matching

Matching results vary with relations and weights, the evaluation will be in different relation with different weights, and the final event network matching is an average of method 1 and 2 in section 5. Method for determining matching: for any given text, 5 and 10 mostly similar texts in the same category and 1 to 2 mostly similar texts in the different category are determined by hand.

The experiments are evaluated by recall and Kendal coefficient, the recall evaluation formula as follows:

$$\text{Recall} = \frac{|MS_{-|AS|} \cap AS| + |MC_{-|AC|} \cap AC|}{|AS| + |AC|} \quad (7)$$

Where AS is manual matching texts set in the same category, AC is manual matching texts set in the different category, $MS_{-|AS|}$ is the paper's $|AS|$ most similar texts set in the same category, $MC_{-|AC|}$ is the paper's $|AC|$ most similar texts set in the different category. Table 6 shows the paper's results considered different relation with three representative weights. From table 6 we can see, in the same weighting factor, performance of semantic relationship is better than share event-element

relation, and share event-element relation is better than neighboring relation. With same relation, the better result can be got only when event weighting factor is greater than relation weighting factor.

Table 6. The evaluation for event network matching

α	β	Semantic Kendall	Share Recall	event-element Kendall	Neighboring Recall	Recall Kendall
0.7	0.3	82%	0.414	81%	0.413	81% 0.411
0.5	0.5	83%	0.419	83%	0.407	77% 0.403
0.4	0.6	80%	0.409	80%	0.413	76% 0.400

8 Conclusions

Text match is one of the key steps in text information retrieval, and different text representation with different matching algorithm. In the paper, a text is represented as an event network, a novel algorithm for matching event network is proposed: event matching based on maximum similarity priority, relation matching based on isotropic-relational-distance matrix, the third and final step is event network matching by integrating event matching and relation matching. The experimental results demonstrate the feasibility and effectiveness of the paper's method.

Acknowledgements. This work is supported by four projects of National Science Foundation of China 60975033/F030503, Innovative Foundation for Graduates of Shanghai University(SHUCX091009).

References

1. Bengoetxea, E.: The graph matching problem, PhD Thesis, ch. 2 (2002)
2. Yang, G.C., Choi, Y.B., Oh, J.C.: CGMA: A Novel Conceptual Graph Matching Algorithm. In: Pfeiffer, H.D., Nagle, T.E. (eds.) *Conceptual Structures: Theory and Implementation*. LNCS (LNAI), vol. 754, pp. 252–261. Springer, Heidelberg (1993)
3. Montes-y-Gómez, M., Gelbukh, A., López-López, A., Baeza-Yates, R.: Flexible Comparison of Conceptual Graphs. In: Mayr, H.C., Lazanský, J., Quirchmayr, G., Vogel, P. (eds.) *DEXA 2001*. LNCS, vol. 2113, pp. 102–111. Springer, Heidelberg (2001)
4. Poole, J., Campbell, J.A.: A Novel Algorithm for Matching Conceptual and Related Graph. In: Ellis, G., Rich, W., Levinson, R., Sowa, J.F. (eds.) *ICCS 1995*. LNCS (LNAI), vol. 954, pp. 293–307. Springer, Heidelberg (1995)
5. Zhu, H.P.: *Semantic Ssearch By Matching Conceptual Graph*, Doctoral Dissertation, Shanghai Jiao Tong University (2006)
6. Chen, X.: Why did John Herschel fail to understand polarization? The differences between object and event concepts. *Studies in History and Philosophy of Science* 34, 491–513 (2003)
7. Zacks, J.M., Tversky, B.: Event structure in perception and conception. *Psychological Bulletin* 127, 3–21 (2001)

Static Patterns Matching for High Speed Networks

Kunpeng Jiang, Huifang Guo, Shengping Zhu, and Julong Lan

National Digital Switching System Engineering
& Technological Research Center
Zhengzhou, Henan, China
bjay371@163.com

Abstract. In response to the need of a large number of static pattern matching on high-speed network, this paper presents a FPGA-based hardware implementation of static pattern matching, which can process in parallel by using the matrix-and algorithm. This method can not only reduce the complexity of programming but also provide the basic of reconfigurable implementation. Experimental results show that the realization is able to reach the theoretical bandwidth multiplying clock frequency by input data width.

Keywords: network, FPGA, NFA, static Pattern, patterns Matching.

1 Introduction

On the one hand, current network bandwidth has been increased rapidly from 1000Mbps to 10Gbps, and now even to 40Gbps. Link bandwidth with 100Gbps has emerged. On the other hand, Network operating environment has become increasingly hostile. Internet worms, viruses, spam, DOS attacks, and malicious access etc. lead to serious problems. In order to ensure network security, a variety of network security devices have come forth. The firewall based on packet header processing, which is commonly used, is inadequate to meet the modern demand for network security. Intrusion detection system (IDS) and Intrusion Prevention System (IPS) which are based on packet content have been widely used. The traditional solution of IDS or IPS is to transfer network data packets on suspicion of malicious forward to some high-performance server, then the server performs deep packet identification (DPI) to identify the data, finally the input packets are processed based on the recognition results. In less than 1000Mbps network, the above method works well. But with the increase of link bandwidth, this method can't meet the requirements of link processing latency. It is inevitable that deep packet identification be performed by hardware directly. The main challenge of deep packet identification is to perform a large number of pattern matching at the same time to meet the requirement of high speed link bandwidth. For instance, SNORT IDS rule sets [1] contain more than 3000 static mode strings. The main demands of static patterns matching in the high-speed network are listed as follows:

- To meet the requirement of wire-speed processing of the link. With the development of network, Modern network bandwidth has been more than 1000Mbps and up to 40Gbps. In order to meet the link requirements, the bandwidth of static pattern matching needs up to 40Gbps.
- To meet the requirement of the number of pattern matching. Network IDS such as SNORT, Bro, Linux, and 7-layer filtering) contains thousands of pattern strings.
- To obtain all matching results. For thousands of pattern strings, the implementation required getting not only the result of matching or not, but also all of the matched pattern strings and matching number.

This paper mainly proposes a static patterns matching algorithm based on hardware in high-speed network, which can reach wire-speed processing with 64Gbps speed, and can process up to 10,000 pattern strings at the same time.

2 Related Works

The main task of patterns matching in modern network application is to identify pattern string from high speed dataflow. There have been many algorithms and hardware design on this aspect. These jobs can be classified three types: one is based on logic circuit and uses FPGA to validate. Another is based on Content Addressable Memory (CAM) or Ternary Content Addressable Memory (TCAM) which can be used to speed matching. The last one used Hash to locate the position of pattern string.

Floyd et al. first proposed implement of regular expression matching in hardware based on non-deterministic automaton (NFA) [2]. Then, Sidhu and Prasanna proposed a model of implement regular expressions based on NFA with combined circuit [3]. Hutchings et al. presented an optimized method of sharing prefix of patterns to reduce the die area [4]. Their methods process one character per one clock when the input text is identified. The shortage of the above methods is that the every character of input of text must be put into the comparator and all the static of matching process should be stored, which leads to high hardware cost. Moreover, the speed of matching can't meet the requirement of current link bandwidth since only one character can be processed per one clock. Clark et al. used pre-decoding to share the character comparators of their NFA implementations and thus hardware resources reduced [5]. The author also explored to process multi-byte per one clock cycle. The method can greatly improve the throughput of character process in high speed network. Hwang proposed a method to match text with string [13]. Because handling multiple characters as one character, their method can not be used to match the regular expression. More authors [14], [15], [16], and [17] did further research on the basis of the above, and some progresses were made. However, the maximum throughput is roughly 40Gbps (gigabit per second). The throughput can not meet the needs of the network nowadays.

If CAM is used to process input text, CAM can soon give the results of whether it matches with the patterns. So, CAM has been widely used in string

matching applications. Gokhale et al. used CAM to parallel find in high speed network [6]. Sourdis et al. made use of pre-coding based on pattern matching with CAM to reduce die area [7]. Yu et al. presented multi-pattern matching algorithm based on TCAM [8]. This method can process complex patterns such as correlative patterns and patterns with not logic.

Dharmapurikar et al. introduced a method based on Hash table which used to implement parallel Bloom filter [9]. Lockwood et al. implemented an intelligent gateway based on parallel Bloom filter which prevented local network from attack of internet worms, viruses [10]. The challenge of Hash algorithm is the requirement of large capacity Hash table. In addition, the confliction of Hash table is still puzzle which leads to the latency of matching variable and unpredictable.

Although, the method based CAM can perform rapid matching, the hardware cost is high and the limit capacity of CAM restricts the number of patterns in practical application. If rapid matching can be implemented by physical logic, the cost is low and the number of matching patterns can be increased with the development of FPGA. Based on character decoder Clark et al. presented in [5], this paper presents a new rapid static pattern matching algorithm implemented in FPGA.

3 Matrix-And Algorithm

Considering the characteristics of the network, without loss of generality, we assume that use 8 bits to represent a character. The maximum operating frequency of modern FPGAs is 300-500 MHz. In order to meet 10Gbps above link speed, multiple characters in a single cycle matching is becoming necessary. The paper proposed a matrix-And patterns matching algorithm to achieve the goal.

For convenience, the pattern of the paper is "uname" from the SNORT (snort rules-snapshot-2.8) rules. Since more than one characters required to be processed per clock cycle, and the status of the last cycle matching results must be considered, our approach not only produces the matching results of the current cycle, but also transfers current status to next cycle. For this example, a 4-bit partial matching result is produced per cycle. The each bit of this partial result represents whether all the previous characters are matched. Matching is denoted with '1', while not matching is denoted with '0'.

First, a 4x5 inverse lower triangular matrix, a 28x5 matrix, and a 4x4 inverse upper triangular matrix are produced based on decoded value of current 32-byte input text and pattern string which now is 5-byte string. The value of the first row of the first column from top left corner is from the result whether the final one character at last cycle match with 'u', and the result bit pass by a D Flip-Flop(DFF). The value of the first row of the second column is from the result whether the final two characters at last cycle match with 'un', and the result bit pass by a D Flip-Flop. The first rows of other columns are similar to these. Besides the first row, other rows of all columns are corresponding to every character of input text by order. Every column represents the matching results of the character in pattern string to corresponding input text character. That is to say, the first

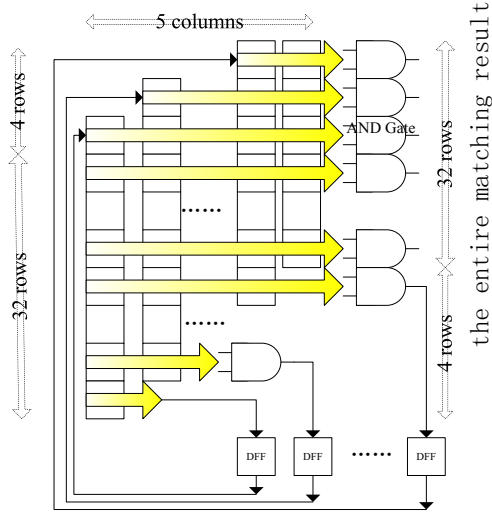


Fig. 1. the construction of the Matrix-And Algorithm

column is corresponding to 'u', the second column is corresponding to 'n', the third column is corresponding to 'a', and go on. The final column is corresponding to 'e'. The 4x4 matrix is used to generate partial matching results of current cycle. The 28x5 matrix is used to generate entire matching results of current cycle. The 4x5 matrix is used to generate the entire matching results across continue two cycles. If one of row in the 4x4 matrix is all equal to '1', we get a partial affirmative matching result. The entire matching result will be postponed until next cycle. If one of row in the 28x5 matrix is all equal to '1', we get a entire affirmative matching result. If one of row in the 4x5 matrix is all equal to '1', we get a entire affirmative matching result which is concluded by considering the previous partial matching result. Fig. 1 shows the construction of the Matrix-And Algorithm.

4 Experimental Results

This section will give experimental results of the Matrix-And algorithm. The algorithm has been validated in an EP4CE115 of Altera Cyclone IV E series. It can complete the task of pattern string matching with high throughput. The Clock frequency of it is 250MHz. The data Width of 128 bits, 256 bits and 512 bits have been respectively processed. Correspondingly, the throughput is 32Gbps, 64Gbps, and 128Gbps respectively. Fig. 2 shows the relation of them.

The experiments used input text widths of 16 to 64 characters and pattern string widths of 4 to 16 characters, which targeted a EP4CE115 of Cyclone IV E series. Table 1 gives corresponding pattern characters capacity of all case.

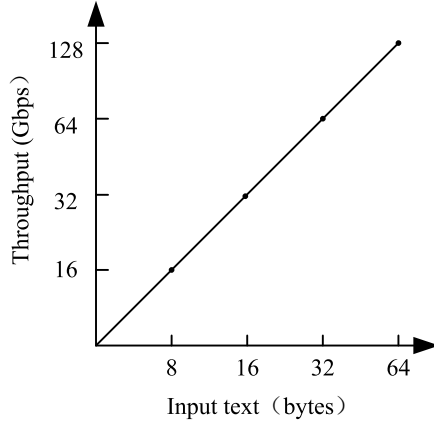


Fig. 2. The relation of throughput and input text

The experiment results shows:

- There is little decrease in character capacity as pattern width increases.
- There is a approximately linear decrease in character capacity as input text width increases.
- The experiments provide insight on the pattern character capacity and input text width trade-off, and input width is directly proportional to processing throughput as the network rate is invariant.

5 Comparison with Previous Work

The **LE** is the basic building block for the FPGA devices. One LE contains a 4-input look-up table(LUT), a 1-bit register, and additional carry and cascade logic, so the average number of LEs per symbol can be the criteria for evaluating the area cost.

Table 2 presents a comparison summary of related previous work with this work. It's clear that our approach is significant improvement in character density at most corresponding input widths. In addition, the throughput of our approach is better scalability than previous designs because LEs/character increases more slowly as input width is increased.

This section compares the results presented in section 3 with the results of previous work on pattern matchers for network applications [5]. The metrics used for comparison are throughput and LEs/character. The throughput of a design is calculated by multiplying the input text width (in bits) processed per cycle by the maximum clock frequency (in Megahertz). LEs/character is logic elements per character, which is a device neutral metric.

Table 1. Pattern characters capacity

pattern characters capacity (bytes)	Input text width (bits)		
pattern character (bytes)	16	32	64
4	33245	15678	7289
5	32774	15189	6814
6	32287	14723	6456
7	31946	14298	5989
8	31342	13856	5547
9	30887	13513	5038
10	30498	13098	4598
11	29864	12675	4115
12	29432	12148	3654
13	28960	11754	3198
14	28441	11219	2786
15	28021	10827	2224
16	27548	10365	1076

Table 2. Throughout and LEs per symbol

Method type	authors	device	Input width (Bytes)	Freq (MHz)	Throughput (Mbps)	LEs/ character
Brute Force	Cho [11]	EP20K	4	90.0	2880	10.6
	Sourdis [7]	Virtex-1000	4	171.0	5472	16.6
		Virtex2-1000	4	344.0	11008	16.6
		Virtex2-6000	4	252.0	8064	19.4
DFA	Moscola [12]	VirtexE-2000	1	37.0	296	5.5
		VirtexE-2000	4	37.0	1184	19.4
NFA	Hutchings [4]	Virtex-1000	1	30.9	247	2.6
		VirtexE-2000	1	52.5	420	2.6
		VirtexE-2000	1	49.5	396	2.5
Comp. NFA	Clark [5]	Virtex-1000	1	100.1	801	1.1
		Virtex2-8000	1	253.0	2024	1.7
		Virtex2-8000	4	218.9	7004	3.1
		Virtex2-8000	8	114.2	7310	5.3
		Virtex2P-125	16	129.0	16516	9.7
		Virtex2P-125	32	141.4	36194	31.5
Matrix-And	Kunpeng Jiang	EP4CE115	16	250.0	32000	3.4-4.1
		EP4CE115	32	250.0	64000	7.3-11.0
		EP4CE115	64	250.0	128000	15.7-106

6 Conclusion

This paper proposed a methodology for designing parallel pattern matching that enable a single FPGA to match ten thousands of static patterns at network rate up to 128Gbps and beyond. Our approach is scalable. The character capacity can

be adjusted by a specific input bandwidth. The lower is bandwidth; the larger can be the character capacity. Since our approach has lower LEs/character. The throughput of this design can go further with the development of network rate. If a larger throughput and character capacity matching system is required, we can use a larger capacity FPGA or ASIC chip, even more FPGAs.

References

1. SNORT Network Intrusion Detection System, www.snort.org
2. Floyd, R.W., Ullman, J.D.: The Compilation of Regular Expressions into Integrated Circuits. *Journal of ACM* 29(3), 603–622 (1982)
3. Sidhu, R., Prasanna, V.K.: Fast Regular Expression Matching Using FPGAs. In: *Field-Programmable Custom Computing Machines (FCCM 2001)*, pp. 227–238 (2001)
4. Hutchings, B.L., Franklin, R., Carver, D.: Assisting Network Intrusion Detection with Reconfigurable Hardware. In: *10th Annual IEEE Symposium on Field-Programmable Custom Computing Machines (FCCM 2002)*, p. 111 (2002)
5. Clark, C.R., Schimmel, D.E.: Scalable Pattern Matching for High Speed Networks. In: *Field-Programmable Custom Computing Machines (FCCM 2004)*, pp. 249–257 (2004)
6. Gokhale, M., Dubois, D., Dubois, A., Boorman, M., Poole, S., Hogsett, V.: Granidt: Towards Gigabit Rate Network Intrusion Detection Technology. In: Glesner, M., Zipf, P., Renovell, M. (eds.) *FPL 2002*. LNCS, vol. 2438, pp. 404–413. Springer, Heidelberg (2002)
7. Sourdis, I., Pnevmatikatos, D.: Fast, Large-Scale String Match for a 10 Gbps FPGA-based Network Intrusion Detection System. In: Y. K. Cheung, P., Constantinides, G.A. (eds.) *FPL 2003*. LNCS, vol. 2778, pp. 195–207. Springer, Heidelberg (2003)
8. Yu, F., Katz, R.H., Lakshman, T.V.: Gigabit rate packet patternmatching using TCAM. In: *Proc. 12th IEEE Int. Conf. Netw. Protocols (ICNP)*, pp. 174–183 (2004)
9. Dharmapurikar, S., Krishnamurthy, P., Sproull, T., Lockwood, J.: Deep packet inspection using parallel bloom filters. In: *Proc. 11th Symp. High Perform. Interconnects*, pp. 44–53 (August 2003)
10. Lockwood, J.W., Moscola, J., Kulig, M., Reddick, D., Brooks, T.: Internet worm and virus protection in dynamically reconfigurable hardware. Presented at the Military Aerosp. Program. Logic Device (MAPLD), Washington, DC, p. E10 (September 2003)
11. Cho, Y.H., Navab, S., Mangione-Smith, W.H.: Specialized Hardware for Deep Network Packet Filtering. In: Glesner, M., Zipf, P., Renovell, M. (eds.) *FPL 2002*. LNCS, vol. 2438, pp. 452–461. Springer, Heidelberg (2002)
12. Moscola, J., Lockwood, J., Loui, R.P., Pachos, M.: Implementation of a Content-Scanning Module for an Internet Firewall. In: *Proceedings of IEEE Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pp. 31–38 (2003)
13. Hwang, W.J., Ou, C.M., Shih, Y.N., Lo, C.T.D.: High throughput and low area cost FPGA-based signature match circuit for network Intrusion detection. *Journal of the Chinese Institute of Engineers* 32, 397–405 (2009)
14. Chang, Y.-K., Chang, C.-R., Su, C.-C.: The Cost Effective Pre-processing Based NFA Pattern Matching Architecture for NIDS. In: *Proceedings of the 2010* (2010)

15. Long, L.H., Hieu, T.T., Tai, V.T., Hung, N.H., Tinh, T.N., Vu, D.D.A.: Eceb: Enhanced constraint repetition block for regular expression matching on FPGA. *ectithailand.org*, pp. 65–74 (2010)
16. Chasaki, D., Wolf, T.: Fast regular expression matching in hardware using NFA-BDD combination. In: *Proceedings of the 6th ACM/IEEE Symposium on Architectures for Networking and Communications Systems, ANCS 2010*, pp. 12:1–12:2. ACM, New York (2010)
17. Nakahara, H., Sasao, T., Matsuura, M.: A regular expression matching circuit based on a modular non-deterministic finite automaton with multi-character transition. In: *The 16th Workshop on Synthesis and System Integration of Mixed Information Technologies, Ballroom*, pp. 359–364 (2010)

Information Propagation in Online Social Networks Based on User Behavior

Niu Li^{1,2} and Han Xiaoting³

¹ Key Laboratory of Ministry of Education for Data Engineering and Knowledge Engineering,
Renmin University of China, Beijing, China

² School of Information Resource Management, Renmin University of China, Beijing, China

³ School of Economics and Management, Beihang University, Beijing, China
libraniu@foxmail.com, hanxiaoting@buaa.edu.cn

Abstract. Along with the development of Internet and Web2.0, online social networks (OSNs) are becoming an important information propagation platform. Therefore, it is of great significance to study the information propagation rules in OSNs. An information propagation model named IP-OSN is proposed in this paper, and some simulation experiments are carried out to investigate the mechanism of information propagation. From the experimental results, we can see that along with the information propagation, the number of known nodes increases and reaches its maximum, then keep an unchanging status. Moreover, from the user behavior aspect, we find that different user behavior in OSNs causes different information propagation results, the more users who are willing to diffuse information, the more scope the information can propagate and the faster the information diffuses. Findings in this paper are meaningful for theory of information propagation and complex networks.

Keywords: Information Propagation/Diffusion, Online Social Networks, User Behavior.

1 Introduction

Along with the development of Internet and Web2.0, online social networks sites such as Facebook, MySpace, LinkedIn and Twitter have become a popular social media platform[1], while they have been developed massively for business and political purposes, such as viral marketing, targeted advertising, political campaigns, and even terrorist activities[2][3]. The users of these sites and the friendships among them constitute the so-called online social networks (OSNs).

Recently study on information propagation rules in online social networks have increased. These studies usually focus on the topology of these social networks [4] and the mechanism of information propagation [5][6]. However, limited work has been done from the user behavior aspect. Therefore, from the view of user behavior, to study the rules of information propagation, and then study how to control the information propagation process, has a very important theoretical value and practical significance.

2 Modeling the Information Propagation

2.1 Model Description

In order to describe the information propagation model clearly, notations of parameters used in the model is shown in Table 1.

In this paper, an information propagation model referred to communicable disease model SIR and SIS [7] named IP-OSN model is proposed in this section. The basic idea of this model is as follows:

Information is diffused in OSNs, and users in OSNs are divided into three types based on user behavior:

Type A users: Information in OSNs cannot affect their attitude of a product or service, and they don't diffuse the information.

Type B users: Information in OSNs can affect their attitude of a product or service in some extent, but they don't diffuse the information.

Type C users: Information in OSNs can affect their attitude of a product or service in some extent, and they are willing to diffuse the information.

At the initial time ($t=0$), there are few nodes in OSNs know the information, and most nodes haven't known the information.

When $t=t+1$, the information propagation process starts until the total running time arrives.

Nodes who know the information and is type C will diffuse it to all its neighbors.

If the neighbor node is type A, it will reject the information.

If the neighbor node is type B or C, it will accept the information.

Table 1. Notations of Parameters

Notation	Description	Notation	Description
N	number of nodes in OSNs	p_1	percentage of type A users
$\langle k \rangle$	averaged degree of nodes in OSNs	p_2	percentage of type B users
n_m	number of initial known nodes	p_3	percentage of type C users
k_m	degree of source node	n_1	number of type A users
k_i	degree of node i	n_2	number of type B users
p	probability of unknown node changed to known	n_3	number of type C users
T	total running time		

2.2 Model Algorithm

The process algorithm of IP-OSN model is shown in Fig. 1:

Step 1: When $t=0$, Initialize information. Set n_m nodes to known, and $(N-n_m)$ nodes to unknown. After the initialization, the number of type A, B, C users are:

$$n_1 = N \times p_1 \quad (1)$$

$$n_2 = N \times p_2 \tag{2}$$

$$n_3 = N \times p_3 \tag{3}$$

Step 2: When $t=t+1$, visit each node (suppose node i) and do Step 2.1 to Step 2.5, and the information propagation process starts until $t=T$.

Step 2.1: If the status of node i is unknown, the algorithm ends; else, go to Step 2.2.

Step 2.2: If node i is type A or B, the algorithm ends; else, go to Step 2.3.

Step 2.3: Visit all the neighbors of node i , then do Step 2.4 k_i times.

Step 2.4: Suppose the algorithm is visiting the neighbor node of node i , and we name the neighbor node j , then do Step 2.4.1 to Step 2.4.3.

Step 2.4.1: If node j is type A, the algorithm ends; else, go to Step 2.4.2.

Step 2.4.2: If status of node j is known, the algorithm ends; else, go to Step 2.4.3.

Step 2.4.3: Change the status of node j to known with the probability of p , and visit next neighbor of node i .

Step 2.5: Visit next node until all the nodes in OSN are visited.

Step 3: End.

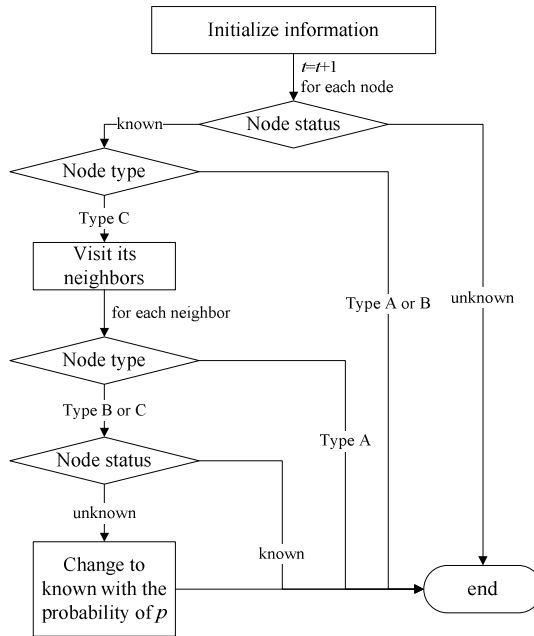


Fig. 1. Information Propagation Process in OSNs

Suppose the number of different nodes at time t is as shown in Table 2. According to mean-field theory [8], after an iteration of the above propagation process, when time = $t+1$, the number of each type of nodes is:

$$n'_{k1} = n_{k1} \quad (4)$$

$$n'_{n1} = n_{n1} \quad (5)$$

$$n'_{k2} = n_{k2} + n_{n2} - n_{n2} \times (1 - p_3 \times p)^{<k>} \quad (6)$$

$$n'_{n2} = n_{n2} \times (1 - p_3 \times p)^{<k>} \quad (7)$$

$$n'_{k3} = n_{k3} + n_{n3} - n_{n3} \times (1 - p_3 \times p)^{<k>} \quad (8)$$

$$n'_{n3} = n_{n3} \times (1 - p_3 \times p)^{<k>} \quad (9)$$

Table 2. Number of Different Nodes at Time t and $t+1$

Node Type	Status	Number at t	Number at $t+1$	Number increased
A	known	n_{k1}	n'_{k1}	
	unknown	n_{n1}	n'_{n1}	
B	known	n_{k2}	n'_{k2}	
	unknown	n_{n2}	n'_{n2}	
C	known	n_{k3}	n'_{k3}	
	unknown	n_{n3}	n'_{n3}	
all	known	n_k	n'_k	n_{kit}
	unknown	n_n	n'_n	n_{nit}

From Eq. (4)-(9), we can get that the number of known nodes increased from time t to $t+1$ is:

$$n_{kit} = n'_k - n_k = n'_{k2} - n_{k2} + n'_{k3} - n_{k3} = n_{n2} - n_{n2} \times (1 - p_3 \times p)^{<k>} + n_{n3} - n_{n3} \times (1 - p_3 \times p)^{<k>} \quad (10)$$

To be simple, that is:

$$n_{kit} = (n_{n2} + n_{n3}) \times (1 - (1 - p_3 \times p)^{<k>}) \quad (11)$$

Similarly, that the number of unknown nodes decreased from time t to $t+1$ is:

$$-n_{nit} = (n_{n2} + n_{n3}) \times (1 - (1 - p_3 \times p)^{<k>}) \quad (12)$$

From Eq. (11) and (12), we can see that the number of known nodes is increasing while the information is diffusing, but the number of increasing known nodes is gradually reduced. This phenomenon will be simulated in the next section.

3 Experiments and Results

3.1 Methodology

There are several methods to study the information propagation in OSNs, such as complex network analysis [9], cellular automata [10] and agent based modeling [11].

In these three methods, we choose agent based modeling as the method to simulate the information propagation process because of its flexibility. Agent based modeling method can adjust the various factors effecting information propagation, therefore how the different combination of factors causes different information propagation effect can be compared easily, which can provide strong evidence for controlling negative information and spreading the positive information. In the simulation process, specific data of each agent can be easily obtained to quantitatively analysis how different user behavior effects the information propagation in real-world OSNs.

In order to prove the efficiency of the above model, a network simulating OSNs is conducted in Netlogo [12], which is simulation software based on agent. Nodes in OSNs are modeled by agents, and interaction among agents is used to simulate the information propagation mechanism in the proposed IP-OSN model. In this way, parameters in IP-OSN model can easily be controlled, which can facilitate the observation of efficiency and effectiveness the model and obtain the data of simulation results to do quantitative analysis.

3.2 Experimental Setup

A randomly generated data set is used for the experiments. There are 2000 nodes in this data set, and the averaged degree is 6. Key features of this data set are summarized in Table 3.

Table 3. Parameters Setting of Data Set

parameter	value	parameter	value
N	2000	p_1	10%
$\langle k \rangle$	6	p_2	20%
n_{in}	4	p_3	70%
k_{in}	6	p	10%

The proposed IP-OSN model is implemented in Netlogo, on a Microsoft Windows 7 Professional platform with SP1 64bit edition. The experimental PC is with an Intel Core i7 2620M CPU, 4 GB DDRII 667 MHz RAM, and a Seagate Barracuda 7200.11 500GB hard disk.

3.3 Experimental Results

First, the running effect of the IP-OSN model is given in Fig. 2 and Fig. 3. In Fig. 2, red nodes represent nodes who don't know the information, while black nodes represent nodes who have already known the information. Fig. 3 shows the percentage change process of the two types of nodes while the information is diffused in OSNs.

As shown in the figures, at the initial time ($t=0$), there are only 4 nodes who know the information, while 99.8% of the initial nodes are unknown. As time passed,

number of known nodes rapidly increases. At $t=244$, the number of known nodes achieves maximum of 1743. After that, number of known nodes and unknown nodes remains unchanged.

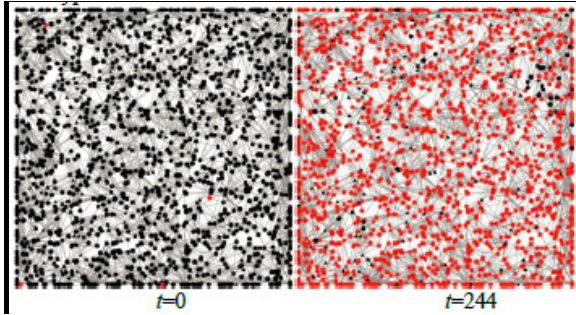


Fig. 2. Initial and Final Status of Information Propagation

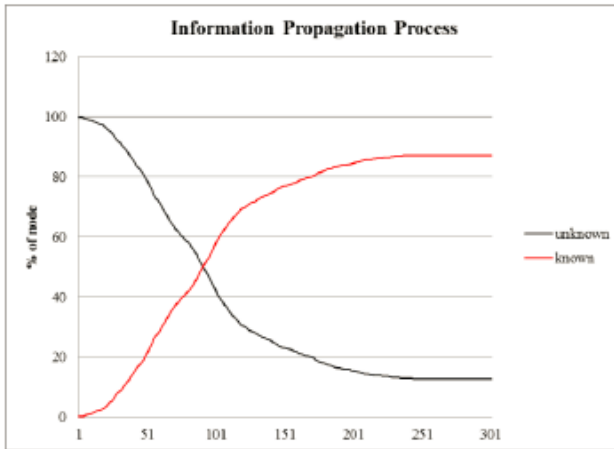


Fig. 3. Evolution Process of Information Propagation

3.4 Analysis of Different User Behavior

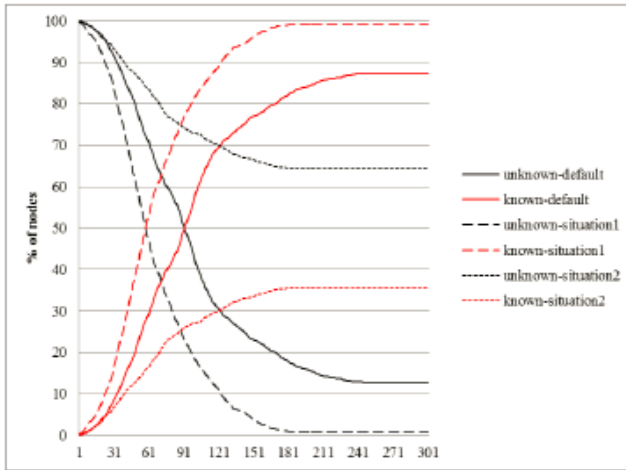
In this section, we change the percentage of different user types in order to analysis how different user behavior affects information propagation in OSNs. Two more situations are assumed. In situation 1, we decreases p_1 and p_2 (the percentage of Type A and Type B user), and in contrast, in situation 2, p_1 and p_2 are increased. The parameter setting is shown in Table 4.

Different propagation process based on the above three situations is shown in Fig. 4.

Table 4. Different Parameter Settings of User Behavior

	P_1	P_2	P_3
Default Situation	10%	20%	70%
Situation 1	0%	10%	90%
Situation 2	15%	30%	60%

As shown in Fig. 4, we can see that the final known nodes are more in situation 1 and less in situation 2 than that in default situation. It shows that when the number of users who want to diffuse information (Type C user) in OSNs increases, the final users who know the information will be more. Also, we can see in Fig. 4 that the propagation process changes in the three situations. Its changing trend goes steeper in situation 1 while gentler in situation 2. It shows that when the number of users who want to diffuse information in OSNs increases, the information speeding speed will be faster, and in contrast, when the number of users who want to diffuse information in OSNs decreases, the information speeding speed will be slower.

**Fig. 4.** Analysis of Different User Behavior

4 Summary

An information propagation model named IP-OSN is proposed in this paper. From the model, we can investigate some information propagation rules in OSNs. From the experimental results, we can see that along with the information propagation, the number of known nodes increases and reaches its maximum, then keep an unchanging status. Moreover, from the user behavior aspect, we find that different user behavior in OSNs causes different information propagation results, the more users who are willing to diffuse information, the more scope the information can propagate and the faster the information diffuses.

References

1. Ahn, Y.Y., Han, S., Kwak, H., Moon, S., Jeong, H.: Analysis of topological characteristics of huge online social networking services. In: Proceedings of the 16th International Conference on World Wide Web, New York, pp. 835–844 (2007)
2. Dunne, Lawlor, M.A., Rowley, J.: Young people's use of online social networking sites - a uses and gratifications perspective. *Journal of Research in Interactive Marketing* 4, 46–58 (2010)
3. Cheung, M.K., Chiu, P.Y., Lee, M.K.O.: Online social networks: Why do students use facebook? *Computers in Human Behavior* 27, 1337–1343 (2011)
4. Amaral, L.A.N., Uzzi, B.: Complex systems - a new paradigm for the integrative study of management, physical, and technological systems. *Management Science* 53, 1033–1035 (2007)
5. Bellini, V., Lorusso, G., Candini, A., Wernsdorfer, W., Faust, T.B., Timco, G.A., Winpenny, R.E.P., Affronte, M.: Propagation of Spin Information at the Supramolecular Scale through Heteroaromatic Linkers. *Physical Review Letters* 106, 227205 (2011)
6. Iribarrena, J.L., Moro, E.: Affinity Paths and information diffusion in social networks. *Social Networks* 33, 134–142 (2011)
7. Hethcote, H.W.: Qualitative Analyses of Communicable Disease Models. *Mathematical Biosciences* 28, 335–356 (1976)
8. Roudi, Y., Hertz, J.A.: Mean field theory for non-equilibrium network reconstruction. *Physical Review Letters* 106, 048702 (2011)
9. Kozinets, R.V., De Valck, K., Wojnicki, A.C., Wilner, S.J.S.: Networked Narratives: Understanding Word-of-Mouth Marketing in Online Communities. *Journal of Marketing* 74, 71–89 (2010)
10. Goldenberg, J., Libai, B.: Talk of the Network: A Complex Systems Look at the Underlying Process of Word-of-Mouth. *Marketing Letters* 12, 211–223 (2001)
11. Smith, T., Coyle, J.R., Lightfoot, E., Scott, A.: Reconsidering Models of Influence: The Relationship between Consumer Social Networks and Word-of-Mouth Effectiveness. *Journal of Advertising Research* 47, 387–397 (2007)
12. Wilensky, U.: <http://ccl.northwestern.edu/netlogo/>

Research on Operation and Management of Railway Transport of Dangerous Goods in Third-Party Logistics Enterprises

Xin Li and Yue-fang Yang

Beijing Jiaotong University, 100044 Beijing, China
11120961@bjtu.edu.cn, yyf583@126.com

Abstract. With China's rapid economic development, the demand for railway transportation of dangerous chemicals is getting stronger and stronger. Consequently, the construction of chemical logistics parks has become hotter than ever. This paper is aimed at describing the necessity of developing the third-party logistics enterprises and studying the operation and management strategies of the third-party logistics enterprises in railway transportation of dangerous goods. Based on the existing laws, regulations and technical standards concerned with dangerous goods transportation, the paper proposes the methods of vehicle management, qualification management and security management of the third-party logistics enterprises in railway transportation of dangerous goods. The legitimate operation of the third-party logistics enterprises in railway transportation of dangerous goods will greatly ease the capacity tension situation, improve the transportation efficiency and help promote the development of China's railway transportation of dangerous goods, which will also make a great significance to China's economy development.

Keywords: transport of dangerous goods, the third-party logistics, qualification management, vehicles management, safety management.

1 Introduction

The definition of third-party logistics in China's national standard <logistics terminology>(GB/T18354-2006) is: A logistics service mode which is independent of the supply and demand sides, provide special or comprehensive logistics system design or system operation for customers. It is a mode that the production enterprises which want to concentrate their efforts on primary industry, give their own logistics activities to the professional logistics service enterprises by way of contract entrusted, at the same time, keeping in close contact with the logistics service enterprises through the information systems, so as to reach the goal of the entire logistics management and control.

Thus we can make the definition of the third-party logistics of railway transport of dangerous goods: it is independent of supply and demand sides, providing specific or comprehensive logistics system design, and transport mainly by railway. It is generally not directly engaged in the business of production, operation or use, but just

providing logistics services for enterprises which related to dangerous goods production, management and use.

The logistics industry developed rapidly in recent years, while there are still several problems: low vehicle utilization, complex business procedures, lacking of transport capacity etc. In order to integrate resources of dangerous goods logistics sectors, enterprises related to dangerous goods should be encouraged to try to use third-party railway transport to organize the implementation of the dangerous goods logistics activities[1].

2 The Necessities to Develop the Third-Party Logistics of Railway Transport of Dangerous Goods

2.1 The Needs of Developing Railway Transport of Dangerous Goods

Third-party logistics enterprises emphasis on the customer's entire full-service, in order to develop the modern logistics, railway transport of dangerous goods enterprise should fully adapt to the multi-functional, systematic demand, build an alliances between the logistics enterprises in warehousing, distribution, information, and packaging, expand extension of logistics service functions and attract logistics and high-tech enterprises. Establishing specialized, integrated third-party logistics enterprises can make a new profit growth point in the railway transport of dangerous goods.

Railway transport of dangerous goods can form logistics enterprises alliances with air, road, waterway and pipeline transport enterprises. The establishment of specialized third-party logistics enterprises can provide a variety of modes of transport, which dominate the transport industry. The third-party logistics will bring new business ideas and management concepts, and promote the service quality and management level of the railway transport. Therefore, the development of railway transport of dangerous goods in the third-party logistics is the needs of their own.

2.2 Meet the Specialized and Agglomerated Needs of the Chemical Enterprises

In recent years, China's chemical logistics has been held fairly rapid growth, accompanied by the development of the industrial park. Chemical logistics has a certain degree of particularity. Compared with the transportation and warehousing of other materials, chemical logistics operation is comprehensive, with high security requirements, in particular in the drug, dangerous goods logistics, which need more stringent technical requirements and high degree professional equipments.

As the division of labor is refined, in order to focus on core competitiveness, most of the international chemical enterprises are trend to outsourcing logistics operations to logistics service providers. The chemical itself has the characteristics of the complex, professional and dangerous, so the healthy development of the chemical logistics is also the reflection of scale economies effect of the chemical industry

cluster which is fully embodied in the guarantee. Therefore, the development of third-party logistics is the needs of the specialized and agglomerated of the chemical enterprises.

2.3 The Needs of Integrating the Cluster Chemical Logistics Resources

Due to its own natural ports and other transport in eastern China, southern coastal areas, the chemical enterprises are gathering here, which become arrival or sending areas of dangerous goods raw materials and other intermediate products. Domestic large-scale chemical enterprises exploit crude oil, liquefied petroleum gas, liquefied natural gas and other resources in the foreign, and these resources need to be transported to domestic processing and sales, which also led to many chemical enterprises gather in the development of the eastern coastal areas.

These areas will have a greater amount chemical logistics needs. Relatively large radius of most of the transport of dangerous goods, taking the cost and security issues into account, long-distance transport of dangerous chemicals should be mainly depending on railway, and the short-distance transport mainly depending on road or port. Therefore, the development of third-party logistics is the needs of integrating cluster chemical logistics resources[2].

3 Analysis of Operations Management

3.1 Analysis of Qualification Management

In China the transport of dangerous goods prohibition agent, thus the third-party logistics should realize the logistics functions by the identity of the shipper or consignee. Besides, according to the relevant provisions of “ Regulations on the Control over Safety of Dangerous Chemicals (2011 edition)”, the third-party logistics enterprises cannot obtain a production license, all in all, third-party logistics enterprises had better in the name of the hazardous chemicals enterprises for handling dangerous goods shipper qualification[3].

According to <The Railway Dangerous Goods Transport Management Rules>, the third-party logistics enterprises should get the qualification certificates of the shipper, logistics enterprises should meet the qualification provisions of “the Railway Dangerous Goods Shipper Qualification License Governing Permission”.

In China, the right to use of the line facilities can be rented, but the command line is still in the hands of the Ministry of Railways, which makes it more complex to get the qualification certificates of the shipper. The involvement of third-party logistics enterprises will make many chemical enterprises get rid of this problem, all transport operations can be done by third-party logistics enterprises, which not only reduce the burden on businesses, but also bring convenience to the management of dangerous goods transport enterprises.

3.2 Analysis of Transporting Vehicle Management

Due to the railway property rights tanker are limited to load crude oil, gasoline, kerosene, diesel, naphtha and non-dangerous goods, heavy oil, lubricating oil. Combined with the transport of dangerous goods regulations, third-party transportation enterprises must have their own dedicated vehicles, and the vehicle performance determines frequency of the risk and the accident, extent in the process of the transport loading and unloading of dangerous goods. It is necessary to set strict vehicle purchase management in order to control accidents caused by vehicle condition during transport.

Third-party railway transport of dangerous goods logistics enterprises should strictly meet the relevant requirements and technical conditions. Third-party railway transport of dangerous goods logistics enterprises should be responsible to the vehicle's operating status and maintenance, should test the performance and status of the periodic of vehicles, the things loaded in the process should meet the requirements, besides, it should do the cleaning when loaded different substances, at the same time the enterprises should also be equipped with fire control equipment, in order to protect the security, reduce costs and improve efficiency.

Transport operations is borne by the third-party logistics and vehicles unified deployment, the same substance between different enterprises can be loaded by the vehicles of the third-party logistics enterprises, all of the logistics needs are united in the third-party logistics enterprises, which not only improve the vehicle utilization, but also expand the transport capacity and bring conveniences to the entire column loading and unloading, which increase the efficiency of transport of goods, and ultimately, enhancing economic efficiency.

3.3 Analysis of Safety Management

3.3.1 Apply to Share of the Special Railway Line

The special railway line is owned to the third-party logistics enterprises, so if the chemical enterprises want to transport, they should apply to share of the special railway line according to the regulations in China. They can acquire the application after they have the "dangerous chemicals business license" or "dangerous chemicals production license".

3.3.2 The Training of Employees

According to the "dangerous goods railway transportation management rules", the related employees in the third-party railway transport of dangerous goods logistics enterprises should be received technical business training, familiar with the dangerous goods knowledge of the position, and grasp the rail transport of dangerous goods regulations, then improve the technical and professional qualities of the practitioners.

According to the regulations, the owner should have at least 15 members who are trained and the shared enterprises should have at least 10, with the involvement of third-party rail logistics enterprises, there will be a reduction in the overall number of trainers, saving human and financial resources, what's more, the employees in the

third-party logistics enterprises have richer experience and more specialized skills, which also increases the safety[4].

3.3.3 The System of Contingency Plans

With the involvement of the third-party logistics enterprises, it achieves the goal of co-ordinated management of the transport of dangerous goods, compared with the usual individual enterprises contingency plans, it is more stringent convergence and careful, when the accident happened, the emergency team equipped with third-party logistics enterprises is more easily unified disposal.

The shared enterprises' contingency plans should have a good convergence with the third party logistics enterprises, with regular exercise and well records, and manner to exclude security risks exposed in the drill timely.

In addition, the third-party logistics enterprises can achieve unified management of the operation area of the special railway line, not only simplifies the process of transport of dangerous goods, but also improve transport efficiency. The unified arrangements for its security and management, can greatly improve the security of working areas.

4 The Problems and Suggestions for the Development of Third-Party Logistics Enterprises

According to the regulations in China, the railway transport of dangerous goods prohibit agent[5]. In fact, the third-party logistics enterprises are different from the freight forwarders, because the freight forwarder is just an enterprise between the shipper (consignee) and the carrier. It is only responsible for the arrangement for the transport of goods, generally do not take the actual transport operations; but third party logistics enterprise refers to an external supplier to provide all or part of the logistics services for the enterprise. it is, non-producers, and non-sales side, but provide third party service in the entire logistics process from production to sales, which generally does not own the goods, but only provide customers with storage distribution and other logistics services.

In China, we can know that roads, waterways, air transport allows the specialized transport of dangerous goods enterprises involved in rail transport, but in railway, the laws and administrative regulations have on relevant provisions of the transport of dangerous goods enterprises, and the third-party logistics enterprises can only be operated by handling dangerous goods business, this is not conducive to third-party logistics enterprises in-depth development of railway transport of dangerous goods business and its value-added services in a long run, it is recommended that the relevant departments to establish the transport of dangerous goods license for railway transport. By the outsourcing of logistics operations to third-party logistics enterprises which are engaged in railway dangerous goods transport, it is not only in favor of the chemical enterprises focusing on their core business, but also conducive to the rapid expansion of third-party logistics enterprises, to obtain economies of scale, and in

addition, the third-party logistics enterprises can also according to the needs of itself, deciding which category to deal with in order to obtain greater profit margins.

However, due to the particularity of the railway transport of dangerous goods, we need to further establish and improve the transport of dangerous goods management system, the enterprise which does not meet the requirements is strictly banned, the license issued by a rigorous examination to ensure that the transport of dangerous goods rank ordered, safe and efficient operation.

5 Conclusion

In recent years, economic develop rapidly in China, accompanied by a strong demand for transport of dangerous goods, as a result, the third-party railway transport of dangerous goods logistics enterprises have a huge space for development, particularly the development of chemical industry park railway transport of dangerous goods in recent years. Do a good job on the third-party rail dangerous goods transport and logistics operations management, based on the related regulations, do our best to simplify the management, it can not only save resources, improve efficiency, but also enhance job security, which will be the new power of the development of the transport of dangerous goods.

References

1. Zhang, X.: Considerations to the development of the logistic market of the railway dangerous goods. *Railway Freight Transport*, 7–9 (2007)
2. Wang, H., Li, Z.: Research on Developing Strategy for Railway Dangerous Goods Logistics. *Logistics Technology*, 32–34, 80 (2009)
3. Wang, D.: The management of dangerous goods of third-party logistics. *China Logistics & Purchasing*, 76–77 (2009)
4. Jia, C.: Study on Operation Management and Investment Mode of the Third Party Railway Dangerous Goods Transport Logistics Enterprise. Beijing Jiaotong University, Beijing (2010)
5. The Ministry of Railways of the People's Republic Of China. The regulations of railways transport of dangerous goods. *Railways transport* [2008]174. China railway publishing house, Beijing (2008)
6. Hai, T.: Problems and protection measures of China's railway transport of dangerous goods. *Railway Transport and Economy*, 61–65 (2011)
7. Hai, T.: Research on Onboard Dynamic Monitoring System of Railway Dangerous Goods Tank Cars. *Railway Transport and Economy*, 36–40 (2011)
8. Qu, Y.: The Problem of Safety Management of Road Dangerous Cargo Transportation and Its Solution. *Logistics Technology*, 26–34 (2011)
9. Bao, M.: Analysis on the safety influence factors during railway transportation of dangerous goods. *Logistics Engineering and Management*, 129–130 (2010)
10. The Ministry of Railways of the People's Republic of China, <http://www.china-mor.gov.cn/>

A Password Authentication Scheme against Smart Card Security Breach^{*}

Jing Shen¹ and Yusong Du^{2,3}

¹ South China Institute of Software Engineering
Guangzhou University, Guangzhou 510990, P.R. China

² School of Information Management
Sun Yat-sen University, Guangzhou 510006, P.R. China

³ Key Lab of Network Security and Cryptology
Fujian Normal University, Fuzhou 350007, P.R. China
szsj_ren@163.com, yusongdu@hotmail.com

Abstract. Remote user authentication is very important for identifying whether communicating parties are genuine and trustworthy. Using a password and a smart card between a login user and a remote server is necessary. Recently, C.T. Li *et al.*'s noted that smart card security breach was not considered in the password authentication scheme given by S.K. Kim *et al.*'s in 2009, then they proposed a remote user authentication scheme against smart card security breach, which was presented in *Data and Applications Security and Privacy* 2011. However, we note that Li *et al.*'s scheme needs a verification table in the remote server. It is well-known that a verification table should not be involved in a good password authentication scheme with smart cards. In this paper, we propose a password authentication scheme against smart card security breach and without maintaining verification tables.

Keywords: password authentication, hash function, smart card.

1 Introduction

Remote user authentication is a procedure that allows a server to authenticate a remote user through an insecure channel. Remote user authentication is very important for identifying whether communicating parties are genuine and trustworthy [1]. Using passwords and smart cards at the same time is common in authentication schemes to check the validity of the login message and authenticate the user [2-6]. It is interesting to design a password authentication scheme with smart cards for remote user authentication.

It has been concluded that a good password authentication scheme with smart cards should satisfy the following requirements [2, 7]. (1) users can freely choose

* This work is supported by the Research Fund of South China Institute of Software Engineering from 2011 to 2012 and the Open Funds of Key Lab of Fujian Province University Network Security and Cryptology (2011008).

and update passwords; (2) low computation complexity; (3) session key agreement; (4) mutual authentication between login users and remote servers; (5) prevention of all the possible attacks such as impersonation attacks, off-line password guessing attacks, replay attacks and parallel-session attacks; (5) resistance to password disclosure to the server, i.e., privacy of passwords; (6) forward security of session keys; (7) without maintaining verification tables; (8) prevention of smart card security breach attacks.

For the smart card security breach, it is important to note that secret information stored in a smart card may be extracted by some physical methods [8, 9]. If a legal user's smart card is lost and it is picked up by a malicious attacker, or an attacker steals the user's smart card, the user's sensitive data may be derived out by the attacker.

Smart card security breach sometimes was neglected when designing a password authentication scheme. K.Shim showed in [10] that three password authentication scheme with smart cards, which were proposed by S.K. Kim *et al.* [3], H.C. Hsiang *et al.* [4] and Y.Y. Wang *et al.* [5] in 2009 respectively, can not resist smart card security breach attacks. In 2011 C.T. Li *et al.* also noted smart card security breach was not considered in S.K. Kim *et al.*'s scheme and they proposed an authentication scheme against smart card security breach instead of S.K. Kim *et al.*'s scheme [6].

However, we note that C.T. Li *et al.*'s scheme needs a verification table in the remote server. A verification table should not be involved in a good password authentication scheme with smart cards, i.e., requirement (7) mentioned above, since verification tables will cause potential security threatens such as stolen-verifier attacks and insider attacks. C.T. Li *et al.*'s scheme solves the problem on smart card security breach attacks by maintaining a verification table in the remote server. Their scheme can not be considered as a good password authentication scheme using smart cards.

In this paper, we propose a password authentication scheme that satisfies all the requirements mentioned above. The proposed scheme can resist smart card security breach attacks and does not need verification tables.

The rest of the paper is organized as follows. Section 2 gives some notations used in this paper. Section 3 recalls S.K. Kim *et al.*'s scheme and C.T. Li *et al.*'s scheme. Section 4 describes our scheme. The security analysis of our scheme is presented in Section 5 and Section 6 concludes the paper.

2 Some Notations

For the convenience of description some notations used in the paper are summarized as follows:

- U_i : The login user.
- ID_i, PW_i, SC_i : The identity, password and the smart card of U_i .
- Server: The remote server.
- x : The master secret key, which is kept secretly and only known by Server.
- $H(\cdot)$: The secure hash function used in authentication schemes.

- \implies : A secure channel.
- \longrightarrow : A public (insecure) channel.
- \oplus : The bitwise XOR operation.

3 Kim *et al.*'s Scheme and Li *et al.*'s Scheme

In this section, we recall S.K. Kim *et al.*'s and C.T. Li *et al.*'s scheme. We point out that C.T. Li *et al.*'s scheme needs to maintain a verification table in the remote server.

Since security flaws appear in the registration phase of Kim *et al.*'s scheme, we only describe its registration phase here. For the full description of the scheme the reader is referred to [3].

(R.I) $U_i \implies$ Server : ID_i, PW_i

U_i chooses his/her identity ID_i and password PW_i , submits $\{ID_i, PW_i\}$ to the authentication server over an secure channel, and remembers PW_i .

(R.II) Server : $\implies SC_i : K_1, K_2, R$

Upon receiving U_i 's login request, the server finishes the following steps.

1. Generate a unique b as the private key of U_i .
2. Compute $K_1 = H(ID_i \oplus x) \oplus b$, $K_2 = H(ID_i \oplus x \oplus b) \oplus H(PW_i \oplus H(PW_i))$, and $R = K_1 \oplus H(PW_i)$.
3. Store $\{K_1, K_2, R\}$ into SC_i and release SC_i to U_i over a secure channel.

Kim *et al.*'s scheme can not resist smart card security breach attacks. If a user's smart card is lost and it is picked up by an attacker or an attacker steals user's smart card. The secrets stored in the smart card may be extracted by some physical methods, then the attacker can off-line guess user's password and can impersonate a legitimate user.

In this scheme, the attacker can breach the secrets $K_1 = H(ID_i \oplus x) \oplus b$, $R = K_1 \oplus H(PW_i)$ and the secure hash function $H(\cdot)$ used in the scheme, which are stored in the smart card. Then, with K_1 and R the attacker off-line guess user's password by performing the following three steps. (1) Select a guessed password PW_i^* ; (2) Compute $K_1' = R \oplus H(PW_i^*)$; (3) Compare K_1 and K_1' .

A match in Step 3 above indicates the correct guess of user's password PW_i . Since password PW_i is usually kept in user's mind it is possible in the off-line case for the attacker to correctly guess the password because of its low-entropy. Thus Kim *et al.*'s scheme is vulnerable to off-line password guessing attacks.

After noting the security flaw of smart card security breach in Kim *et al.*'s scheme, C.T. Li *et al.* described an authentication scheme which resolves the security flaw of smart card security breach. We will not give the description of Li *et al.*'s scheme. The reader is referred to [6]. Here we only point out that Li *et al.*'s scheme needs a verification table in the server.

In the registration phase of Li *et al.*'s scheme, U_i chooses his/her identity ID_i and password PW_i , and generates a random number RN_1 . Then U_i computes $H(H(PW_i \oplus RN_1))$ and submits ID_i and $H(H(PW_i \oplus RN_1))$ to the server over an secure channel, and remembers PW_i . Upon receiving U_i 's login request, the

server maintains a account table (AT) for the registration service and the format of the AT is shown as follows:

User's identity	Registration times	Verification parameter
ID_i	$N = 0$	$H(H(PW_i \oplus RN_1))$

The 2nd field of AT records $N = 0$ if it is U_i 's initial registration, otherwise, the server sets $N = N + 1$ in the existing field for U_i . The verification parameter will be used in verification phase and will be changed into $H(H(PW_i \oplus RN_2))$ after a legal login of U_i , where RN_2 is a new random number generated by U_i for the login request. In the password update phase the verification table is also involved.

It is well-known that a verification table should not be involved in a good password authentication scheme with smart cards. Li *et al.*'s scheme can not work without the verification table. Their scheme can not be considered as a good password authentication scheme using smart cards.

4 The Proposed Scheme

In this section, we describe a password authentication scheme which resolves security flaws related to smart card security breach and verification tables mentioned above. There are four phases in our scheme: registration, login-verification, session key agreement and password update.

4.1 Registration Phase

(R.I) $U_i \implies$ Server : ID_i

U_i chooses his/her identity ID_i and submits the registration request to the server over an secure channel.

(R.II) Server $\implies SC_i$: ID_i, r, E, \bar{D}

Upon receiving U_i 's registration request, the server finishes the following steps.

1. Generate a random number r and a random number b .
2. Compute $E = H(ID_i \oplus x \oplus r) \oplus b$ and $\bar{D} = H(ID_i \oplus x \oplus b)$.
3. Store $\{ID_i, r, E, \bar{D}\}$ into smart card SC_i .
4. Release SC_i to U_i over an secure channel.

After receiving the smart card U_i chooses and remembers his/her password PW_i , then enters the password into the smart card. The smart card chooses a random number s , then computes $D = \bar{D} \oplus H(H(PW_i) \oplus s)$ instead of \bar{D} . Finally the smart card keeps $\{ID_i, r, E, D, s\}$, and disposes \bar{D} and PW_i securely.

4.2 Login and Verification Phase

(LV.I) $U_i \implies SC_i$: PW_i

The user U_i enters PW_i into smart card SC_i .

(LV.II) $SC_i \rightarrow \text{Server} : ID_i, r, E, M_1, T_{U_i}$

With PW_i smart card SC_i takes current time-stamp T_{U_i} and computes

$$M_1 = H(H(D \oplus H(H(PW_i) \oplus s) \oplus E) \oplus T_{U_i}).$$

Then SC_i submits $\{ID_i, r, E, M_1, T_{U_i}\}$ to the server.

(LV.III) Server : $\rightarrow SC_i : r^{new}, E^{new}, \bar{D}^{new}, M_2, M_3, T_S$

Upon receiving the login request, the server finishes the following steps.

1. Verify the validity of T_{U_i} . If it is invalid, reject U_i 's login request.
2. Derive $b = E \oplus H(ID_i \oplus x \oplus r)$ and check whether $H(H(H(ID_i \oplus x \oplus b) \oplus E) \oplus T_{U_i})$ is equal to the received M_1 . If it does not hold, terminate the communication. If it holds, authenticate U_i .
3. Choose a new random number r^{new} and a new random number b^{new} .
4. Take current time-stamp T_S and compute

$$E^{new} = H(ID_i \oplus x \oplus r^{new}) \oplus b^{new},$$

$$\bar{D}^{new} = H(H(ID_i \oplus x \oplus b) \oplus r^{new}) \oplus H(ID_i \oplus x \oplus b^{new}),$$

$$M_2 = H(H(H(ID_i \oplus x \oplus b) \oplus E^{new}) \oplus T_S),$$

and

$$M_3 = H(H(H(ID_i \oplus x \oplus b) \oplus \bar{D}^{new}) \oplus T_S).$$

5. Send $\{r^{new}, E^{new}, \bar{D}^{new}, M_2, M_3, T_S\}$ to SC_i .

Upon receiving the message from the server, SC_i finishes the following steps.

1. Verify the validity of T_S . If it is invalid, terminate the communication.
2. Check whether $H(H(D \oplus H(H(PW_i) \oplus s) \oplus E^{new}) \oplus T_S)$ is equal to received M_2 . If it does not hold, terminate the communication.
3. Check whether $H(H(D \oplus H(H(PW_i) \oplus s) \oplus \bar{D}^{new}) \oplus T_S)$ is equal to received M_3 . If it does not hold, terminate the communication.
4. If both of them hold, authenticate the server.
5. Choose a new random number s^{new} instead of s and compute

$$D^{new} = H(D \oplus H(H(PW_i) \oplus s) \oplus r^{new}) \oplus \bar{D}^{new} \oplus H(H(PW_i) \oplus s^{new}).$$

6. Change E, D, r and s into $E^{new}, D^{new}, r^{new}$ and s^{new} respectively.

4.3 Session Key Agreement Phase

The session key establishment phase is based on a legal login and Diffie-Hellman key exchange mechanism.

Suppose that the server and smart cards share a cyclic group \mathcal{G} of large enough order l with generator $g \in \mathcal{G}$. After authenticating U_i the server choose random number α with $1 \leq \alpha \leq l - 1$ and sends

$$K = H(H(ID_i \oplus x \oplus b^{new}) \oplus r^{new}) \oplus g^\alpha$$

to the smart card.

Upon receiving the message from the server, the smart card derive g^α by

$$g^\alpha = H(D^{new} \oplus H(H(PW_i) \oplus s^{new}) \oplus r^{new}) \oplus K.$$

Then the smart card choose random number β with $1 \leq \beta \leq l - 1$ and sends

$$H(D^{new} \oplus H(H(PW_i) \oplus s^{new}) \oplus r^{new}) \oplus g^\beta$$

to the server. The smart card uses $g^{\alpha\beta}$ as the session key.

Finally, upon receiving the message from the smart card, the server derive g^β similarly and uses $g^{\alpha\beta}$ as the session key.

4.4 Password Update Phase

The password update phase includes a legal login and verification phase.

(PU.1) $U_i \implies SC_i : PW_i, PW_i^{new}$

U_i enters the old password PW_i and the new password PW_i^{new} to SC_i .

(PU.2) $SC_i \longrightarrow \text{Server} : ID_i, r, E, C, M_0, M_1, T_{U_i}$

According to password PW_i , smart card SC_i chooses s^{new} and computes

$$C = H(D \oplus H(H(PW_i) \oplus s) \oplus r) \oplus H(H(PW_i^{new}) \oplus s^{new}),$$

$$M_0 = H(H(D \oplus H(H(PW_i) \oplus s) \oplus C) \oplus T_{U_i}),$$

and

$$M_1 = H(H(D \oplus H(H(PW_i) \oplus s) \oplus E) \oplus T_{U_i}),$$

where T_{U_i} is current time-stamp. Then SC_i sends $\{ID_i, r, E, C, M_0, M_1, T_{U_i}\}$ to the server.

(PU.3) Server : $\longrightarrow SC_i : r^{new}, E^{new}, \bar{D}^{new}, M_2, M_3, T_S$

Upon receiving the update request, the server finishes the following steps.

1. Derive $b = E \oplus H(ID_i \oplus x \oplus r)$.
2. Check whether $H(H(H(ID_i \oplus x \oplus b) \oplus C) \oplus T_{U_i})$ is equal to the received M_0 and check whether $H(H(H(ID_i \oplus x \oplus b) \oplus E) \oplus T_{U_i})$ is equal to the received M_1 . If one of them does not hold, terminate the communication.
3. Derive $H(H(PW_i^{new}) \oplus s^{new})$ from $H(H(PW_i^{new}) \oplus s^{new}) = C \oplus H(H(ID_i \oplus x \oplus b) \oplus r)$.
4. Choose r^{new} and b^{new} . Compute $\bar{D}^{new} = H(H(ID_i \oplus x \oplus b) \oplus r^{new}) \oplus H(H(PW_i^{new}) \oplus s^{new}) \oplus H(ID_i \oplus x \oplus b^{new})$.
5. Compute E^{new} , M_2 and M_3 as in the verification phase.
6. Send $\{r^{new}, \bar{D}^{new}, E^{new}, M_2, M_3, T_S\}$ to SC_i .

Upon receiving the message from the server, SC_i finishes the following steps.

1. Authenticate the server via M_2 and M_3 .
2. Compute $D^{new} = H(D \oplus H(H(PW_i) \oplus s) \oplus r^{new}) \oplus \bar{D}^{new}$.
3. Change E , D , r and s into E^{new} , D^{new} , r^{new} and s^{new} respectively.

5 Security Analysis of the Proposed Scheme

In this section, we analyze the security of the proposed scheme. We point out that the proposed scheme can satisfy 8 requirements mentioned in Section 1 (Introduction).

The password PW_i is freely chosen by U_i in the registration phase. U_i can update his/her password by a password update phase. The proposed scheme includes a session key agreement phase that can generate a secure session key. The proposed scheme has low computation complexity since it is mainly based on operations of hash functions except the session key agreement phase.

Impersonation attacks cannot work in the scheme without U_i 's smart card. Even if the attacker steals the U_i 's smart card he/she can not impersonate U_i to login the server since he/she does not know U_i 's password PW_i .

Off-line password guessing is also impossible. Login and verification phase will be completed by the server after smart cards submit the login request and the verification information. Without the remote server smart cards are not able to determine if a guessed password is valid.

The attacker cannot deceive the server by replay attacks. The verification information including ID_i, r, E, M_1, T_{U_i} or the password update information including $ID_i, r, E, C, M_0, M_1, T_{U_i}$ used in a previous run of the protocol is useless in a coming run of protocol. Without the smart card SC_i and password PW_i the attacker cannot derive the new information for a coming run of protocol from the old information used in a previous run of the protocol.

The proposed scheme achieves mutual authentication between login users and remote servers. The server authenticates the user by checking M_1 while the user authenticates the server by checking M_2 . Mutual authentication between login users and remote servers is the countermeasure for man-in-the-middle attacks.

Parallel-session attacks cannot work in the scheme. A set of verification information $\{ID_i, E, M_1, r, T_{U_i}\}$ in the current run of protocol uniquely corresponds to the user ID_i and is useless for other users. The attacks cannot derive any useful information that can be used in another run of protocol.

The user does not submit directly the password PW_i to the server. The server cannot derive U_i 's password PW_i according to the information from U_i 's smart card. The verification information $\{ID_i, E, M_1, r, T_{U_i}\}$ from the smart card do not include any password information. In the password update phase the server can derive $H(H(PW_i^{new}) \oplus s^{new})$ from C . But it is still very hard for the server to guess PW_i^{new} since the server does not know s^{new} .

We called session keys satisfy forward security in the scheme if the attacker cannot get any session key established earlier than the time point at which the master key of the server x was lost. The Diffie-Hellman key exchange mechanism guarantees the forward security of session keys. The attacker cannot compute $g^{\alpha\beta}$ because of the Diffie-Hellman intractable problem even if the master key of the server x is lost and g^α and g^β are derived.

It clear that there is no a verification table in the scheme. The server only need keep a master key x . Except the password PW_i all the verification information of U_i 's is kept by the smart card.

The scheme is able to prevent smart card security breach attacks. Suppose that U_i 's smart card is lost and it is picked up by an attacker, or an attacker steals U_i 's smart card. The secrets stored in the smart card may be extracted, i.e., $\{E, D, r, s\}$ are obtained by the attacker. However, the attacker is still unable to derive U_i 's password PW_i from $\{E, D, r, s\}$. The attacker has to guess x, b and PW_i at the same time. The attacker faces the problem of finding an original image of a secure hash function, which is intractable.

6 Conclusion

In this paper, we note that C.T. Li *et al.*'s scheme, that can prevent smart card security breach attacks, needs a verification table in the remote server. The existence of verification tables is not good for a secure password authentication scheme. We proposed a password authentication scheme against smart card security breach and without maintaining verification tables.

References

1. Lamport, L.: Password authentication with insecure communication. *Communications of the ACM* 24(11), 770–772 (1981)
2. Liao, I.E., Lee, C.C., Hwang, M.S.: A password authentication scheme over insecure networks. *Journal of Computer and System Sciences* 72(4), 727–740 (2006)
3. Kim, S.K., Chung, M.G.: More secure remote user authentication scheme. *Computer Communications* 32(6), 1018–1021 (2009)
4. Hsiang, H.C., Shih, W.K.: Weaknesses and Improvements of the Yoon-Ryu-Yoo Remote User Authentication Scheme using Smart Cards. *Computer Communications* 32(6), 649–652 (2009)
5. Wang, Y.Y., Liu, J.Y., Xiao, F.X., Dan, J.: A More Efficient and Secure Dynamic ID-Based Remote User Authentication Scheme. *Computer Communications* 32(6), 583–585 (2009)
6. Li, C.-T., Lee, C.-C., Liu, C.-J., Lee, C.-W.: A Robust Remote User Authentication Scheme against Smart Card Security Breach. In: Li, Y. (ed.) *DBSec 2011*. LNCS, vol. 6818, pp. 231–238. Springer, Heidelberg (2011)
7. Mao, W.: *Modern Cryptography: Theory and Practice*. Prentice Hall, New Jersey (2003)
8. Kocher, P., Jaffe, J., Jun, B.: Differential Power Analysis. In: Wiener, M. (ed.) *CRYPTO 1999*. LNCS, vol. 1666, pp. 388–397. Springer, Heidelberg (1999)
9. Messerges, T.S., Dabbish, E.A., Sloan, R.H.: Examining smart-card security under the threat of power analysis attacks. *IEEE Transactions on Computers* 51(5), 541–552 (2002)
10. Shim, K.A.: Security Flaws in Three Password-Based Remote User Authentication Schemes with Smart Cards. *Cryptologia* 36(1), 62–69 (2012)

A Vulnerability Attack Graph Generation Method Based on Scripts

Bo Han^{1,*}, Qing Wang¹, Fajiang Yu², and Xianda Zhang³

¹ International School of Software, Wuhan University

² School of Computer, Wuhan University

³ School of Geodesy and Geomatics, Wuhan University

bhan@whu.edu.cn

Abstract. The vulnerability attack graph is an important method for prevention of network attacks. However, the huge amount of vulnerability has caused great difficulties for attack graph generation. By using the general search methods, we often retrieve many unrelated vulnerabilities from database, difficult to locate the true exploits of points. In this paper, we proposed an attack graph generation method based on scripts. It applied text mining technology to analyze in-depth vulnerability information. We have got a relatively small range of vulnerability. By combinations of these related vulnerability, we generated the attack graphs. The approach helps attack graph play better defense functionality.

Keywords: information security, vulnerability database, attack graph, attack scripts, TF-IDF.

1 Introduction

With the continuous development of computer and Internet technology, the network has become an integral part of people's daily work and life. The Internet technology facilitate the sharing of network resources, but it also brings a variety of security risks at the same time. Information security has become a focus of great concern to national safeguard. As today's security technology continues to evolve, the use of a single vulnerability has basically been unable to achieve the successful invasion of the information network. A hacker now exploits multiple vulnerability exists in the multi-hosts and multi-class operating systems [1-5]. Each vulnerability achieves a particular goal, and the purpose of invasion can be reached step by step. Compared with the previous attacks, this approach has a clear target and its destructive power is very large. Once it is outbreak, it is very likely to have serious consequences. For example, the Stuxnet virus [6], known as the first network super weapon , swept through the world of industry. If we can detect the vulnerability combination in an early stage, then we can prevent such virus effectively, reducing the losses caused by virus attack or even in advance to prevent the outbreak of the virus.

* Corresponding author.

2 Vulnerability Attack Graph

From the perspective of an attacker, attack graph enumerates all possible attack paths from the attack starting point to the target, which intuitively provides a visual representation model of attack process scene. It also helps defenders understand the relationship between various vulnerabilities and potential threats [7-9].

Building attack graph is a challenging task [10-11]. In early studies, an attack graph is only hand-built by network specialists. However, with the increase of vulnerabilities, it is hard to hand-construct attack graphs. For example, China has also established a National Vulnerability Database (CNNVD). It has accumulated nearly 40,000 vulnerability records. In general, vulnerability database lists detailed description information for each vulnerability: date of publication, a brief text description of vulnerability, vulnerability type, degree of risk, correlated operating system, vulnerability testing and certification file name. Faced with such a large amount of information, the use of purely manual way to build a large-scale attack graph has become an impossible mission.

3 The Method of Building Attack Graph Based on Scripts

Although the huge amounts of data in the vulnerability database poses challenges to us. It contains important security information and provides favorable conditions for text mining. In this paper, we proposed the method of building attack graphs based on scripts. Its main idea is to extract keywords from the attack scripts and then perform retrieval in the vulnerability database by using text mining technology. Next, by combination of the most matched vulnerabilities, an attack graph is built.

3.1 Script Generation

Attack scripts are the texts which describe the main characteristics and behavior of an attack on the attack path. In general, information security experts or hackers hand-code an attack script by using security knowledge.

3.2 Keywords Extraction

The keyword is the basic unit representing the attack scripts. It describes the scripts' behavioral characteristics. Keywords have certain properties:

- 1) Keywords able to really identify the contents of the attack script;
- 2) Keywords have the ability to distinguish between attack scripts and other text;
- 3) The number of keywords is not large. But the number of keywords can be much larger than the traditional query with just a few keywords.

After extracting keywords, we can then make use of text mining technology onto vulnerability database, and to find the best fit vulnerabilities, and finally generate an attack graph with a controllable scale.

3.3 The Process of Building an Attack Graph

For each attack node in the attack path, we perform the following steps,

- 1) Extraction of keywords from attack scripts;
- 2) Using TF-IDF algorithm to retrieve relevant vulnerability;
- 3) Sort vulnerabilities according to the calculated degree of relevance;
- 4) Put the vulnerability with high degree of relevance into the candidate set.

3.4 Retrieval of Relevant Vulnerability

The main technical means for Script-based vulnerability attack graph generation method is TF-IDF text mining algorithm.

Vulnerability database contains tens of thousands of vulnerabilities, and it is still in expanding. The description for each vulnerability includes many of the terms. General search treated equally between important terms and unimportant terms, which results in retrieval of some irrelevant vulnerability. By using TF-IDF text mining algorithm, we can measure the relevant extent between a vulnerability description and a few keywords, and only returns highly relevant vulnerability as results, thus narrowing the scope of building attack graphs.

The vulnerability description consists of a large amount of text, which contains explanation about how vulnerability affects the software, the ways and means for affecting, the exploiting scene and so on. To deal with the text, we apply the concept of vector space[12]. The vector space means that the representation of a vulnerability description with a vector, where each element corresponds to the weight about a keywords. In another word, if a vulnerability description contains a keyword, the corresponding element in a vector is assigned to a weight; otherwise it is assigned with 0. After a text has been transformed into a vector, the computation to text can be applied on the vectors.

For comparing two text vectors, we need to apply TF-IDF (Term Frequency-Inverse Document Frequency) techniques.

TF (Term Frequency) is the occurrences number of a word t in a text. It can be used to measure the degree of association between the term t and a given text d . Generally, if the text does not contain the term, TF is defined as zero, otherwise it is defined as a non-zero frequency value. However, only using term frequency is not a good measure since the term frequency is related to text length. Next, the relevance degree of a text with a keyword occurring 10 times is not 10 times than a text with a keyword only occurring once. By considering the above factors, we define the following equation for computing TF value:

$$TF(d,t) = \log(1 + \frac{n(d,t)}{n(d)})$$

Where, $n(d)$ represents the number of terms in the text d , $n(d,t)$ represents the number of occurrences of term t in text d .

The IDF (Inverse Document Frequency) represents the importance of the word t . If the word t appears in many texts, it reduces its function for text distinction and also

reduces its importance. If a word has a very small IDF value, it can be ignored when using text mining techniques. We use the following equation for computing IDF,

$$IDF(t) = \log \frac{1 + |d|}{|d_t|}$$

Where $|d|$ is the number of texts in collection, $|d_t|$ is the number of texts containing term t .

In vector space model, we combine TF and IDF value together to form TF-IDF weight,

$$W_i = TF(d, t) \times IDF(t)$$

So, the relevance degree between a text d to a term collection Q is defined as the following,

$$r(d, t) = \sum_{t \in Q} W_i$$

The extracted keywords from the attack scripts can be regarded as term set Q . The vulnerability description can be regarded as a text d . We can apply TF-IDF algorithm to calculate the relevance degree between vulnerability and keywords.

4 Experiments

4.1 The Experiment Settings

Script-based attack graph generation experiments are based on the CVE vulnerability database, with 53967 vulnerability instances. The attack scripts are provided by information security experts, describing the attack path for invading the system. Attack script is shown in Fig. 1.

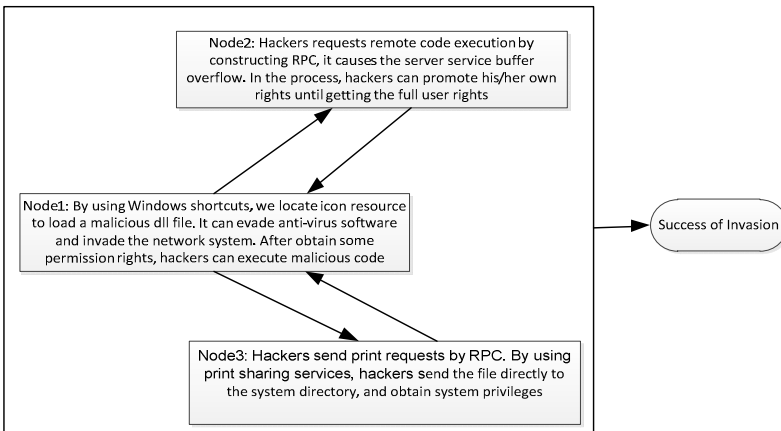


Fig. 1. Attack Script Example

We extract the keywords from attack scripts for each attack node as follows:

The keywords in attack node 1:lnk,shortcut,icon,malware,dll,windows explorer, remote attack.

The keywords in attack node 2:server service, rpc request, remote attack, overflow, permission, microsoft.

The keywords in attack node 3:print spooler, system directory, permission, remote attack, rpc request, microsoft.

4.2 Experimental Results and Analysis

We compute TF-IDF values for the extracted keywords and text descriptions in vulnerability database, and then sort them by their relevance degrees. Only vulnerabilities with relevance degree larger than a threshold (it is set to 0.39) are put into vulnerability candidate set. The computation results are shown in Table 1.

These vulnerabilities are combined to form a complete attack graph. The number of candidate vulnerability in an attack step is controlled within 10, effectively reducing the workload of manual selection.

The common keyword search method will result in up to dozens or even hundreds of records. For example, if we use rpc request as keywords to perform a general search in the vulnerability database, 87 records will be popped up. Using printer as a keyword, 40 records will be resulted in. And the search result has not be processed by standardization and quantitative computing, some relevant results have not put on the top of the result list. In comparison, the script-based attack graph generation method has resulted in a small size of vulnerability. Especially, we notice that the attack path from Stuxnet virus is listed in the results (combination of vulnerability: cve-2010-2568, cve-2010-2729, cve-2008-4250). This experiment illustrates the practical significance of the attack graph generation algorithm, and proves that it has a positive effect to prevent unknown attacks.

Table 1. Relevant Vulnerability for Each Attack Node

	Attack node1	Attack node2	Attack node3
	CVE-1999-0280		
Vulnerability	CVE-2004-0537		CVE-2000-0184
Candidate	CVE-2001-1386	CVE-2002-0642	CVE-2009-0228
Set	CVE-2005-0519	CVE-2008-4250	CVE-1999-0353
	CVE-2005-0520	CVE-2004-1560	CVE-1999-0564
	CVE-2007-6535		CVE-2002-2201
	CVE-2009-4965		CVE-2004-1856
	CVE-2010-2568		CVE-2010-2729
	CVE-2010-1640		

5 Conclusion

Based on the text mining techniques, we proposed a script-based method for vulnerability attack graph generation. We apply TF-IDF algorithm to quantitatively

measure the relevance degree between script keywords and vulnerability description. Our experimental results show that the approach can effectively select the relevant vulnerability for building an attack path. Compared to the general search method, our approach can narrow the size of vulnerability in 7-10 times, greatly reducing the complexity of building attack graphs. We observe the attack path from Stuxnet virus is shown in our result list. It confirms our method can effectively prevent exploits and has wide application in information security area.

Acknowledgment. This work was supported by grant 2009CB723905 from National Basic Research Program of China (973 Program), by grant 2011CDB447 from natural science foundation of Hubei province, by grant 216271207 for talent development at Wuhan University. It was also partially supported by grant 61103220 from National Natural Science Foundation of China and by grant 3101044 from the Fundamental Research Funds for the Central Universities in China.

References

1. Dawkins, J., Hale, J.: A systematic approach to multistage network attack analysis (2004)
2. Murphy, C.T., Yang, S.J.: Clustering of multistage cyber attacks using significant services. In: 13th Conference on Information Fusion (2010)
3. Mathew, S., Upadhyaya, S., Sudit, M., et al.: Situation awareness of multistage cyber attacks by semantic event fusion. In: IEEE Military Communications Conference, MILCOM, pp. 1286–1291 (2010)
4. Du, H., Liu, D., Holsopple, J., et al.: Toward ensemble characterization and projection of multistage cyber attacks. In: ICCCN (2010)
5. Yang, S.J., Stotz, A., Holsopple, J., et al.: High level information fusion for tracking and projection of multistage cyber attacks. *Information Fusion* 10(1), 107–121 (2009)
6. Hudson, J.: Weaponised malware: how criminals use digital certificates to cripple your organization. *Network Security* 6, 12–14 (2011)
7. Barik, M.S., Mazumdar, C.: A novel approach to collaborative security using attack graph. In: IMSAA 2011 (2011)
8. Somesh, J., Oleg, S., Jeannette, M.W.: Minimization and reliability analyses of attack graphs. Technical Report CMU-CS-02-109, Carnegie Mellon University (2002)
9. Wang, L., Noel, S., Jajodia, S.: Minimum-cost network hardening using attack graphs. *Computer Communications* 29(18), 3812–3824 (2006)
10. Zhong, S., Xu, G., Yang, Y., Yao, W., Yang, Y.: Algorithm of generating host-based attack graph for overall network. *Advances in Information Sciences and Service Sciences* 3(8), 104–110 (2011)
11. Zhang, B., Lu, K., Pan, X., Wu, Z.: Reverse search based network attack graph generation. In: CiSE (2009)
12. Han, J., Kambe, M.: *Data mining: concepts and techniques*, 2nd edn. Morgan Kaufmann Publishers (2006)

DFP-Growth: An Efficient Algorithm for Mining Frequent Patterns in Dynamic Database

Zailani Abdullah¹, Tutut Herawan², A. Noraziah², and Mustafa Mat Deris³

¹Department of Computer Science University Malaysia Terengganu 21030 Kuala Terengganu, Terengganu, Malaysia

²Faculty of Computer Systems and Software Engineering University Malaysia Pahang Lebuhraya Tun Razak, 26300 Kuantan Pahang, Malaysia

³Faculty of Computer Science and Information Technology University Tun Hussein Onn Malaysia Parit Raja, Batu Pahat 86400, Johor, Malaysia
zailania@umt.edu.my, mmustafa@uthm.edu.my

Abstract. Mining frequent patterns in a large database is still an important and relevant topic in data mining. Nowadays, FP-Growth is one of the famous and benchmarked algorithms to mine the frequent patterns from FP-Tree data structure. However, the major drawback in FP-Growth is, the FP-Tree must be rebuilt all over again once the original database is changed. Therefore, in this paper we introduce an efficient algorithm called Dynamic Frequent Pattern Growth (DFP-Growth) to mine the frequent patterns from dynamic database. Experiments with three UCI datasets show that the DFP-Growth is up to 1.4 times faster than benchmarked FP-Growth, thus verify its efficiencies.

Keywords: Efficient algorithm, Frequent patterns, Dynamic database.

1 Introduction

In the several decades, mining frequent patterns [1-5] or association rules [6-22] have been received much research attentions and developments. The first effort was done by Agrawal et al. [6] in 1993 and it is still became an evergreen topic in data mining. In mining frequent patterns, the main objective is to discover the interesting association rules from data repositories. Based on data mining interpretation, a set of item in pattern is defined as an itemset. The itemset is said to be frequent, if it appears equal or exceed the predefined minimum support thresholds. The item (or itemset) support is defined as a probability of item (or itemset) occurs in the transaction. Besides that, confidence is another alternative measurement used in pair with support. The confidence is defined as the probability of the rule's consequent that also contain the antecedent in the transaction. Association rules are said to be strong if it meets the minimum confidence.

At the moment, FP-Growth [5] is one of the famous and benchmarked algorithms to mine the frequent patterns. It is based on the utilization of a compact trie data structure called FP-Tree. However, the major drawback occurs in FP-Growth is, the FP-Tree must be entirely rebuilt again if the original database is changed. For simplicity, let assume that the original database consist of 50,000 transactions. Thus,

FP-Growth will build the FP-Tree and finally mine the frequent patterns. Subsequently, the numbers of the transactions in database are increased into 80,000. In order to mine back the recent frequent patterns, the FP-Tree must be rebuilt again. This is because the previous FP-Tree is unable to accept and adapt new transactions. As a result, the processing time to build and mine the new FP-Tree will be proportionally increased due to the latest size of database.

Therefore, in this paper we propose a new algorithm called DFP-Growth in an attempt to mine the frequent patterns from updatable database. In this algorithm, only new transactions will be involved in building a new FP-Tree. The previous patterns that have been generated from the past FP-Tree will be reused again in producing the latest frequent patterns. Indeed, both set of patterns will be merged together to produce the final frequent patterns. Prior to this, they will be sorted in canonical order and not in standard support descending order.

In summary, the contribution of this paper is as follows. First, we propose a scalable DFP-Growth algorithm that can focus on building a new FP-Tree purely based on the new transactions. Second, we sort the itemset from previous and latest patterns in canonical order before producing the final canonical frequent patterns. Third, we do experiment with three benchmarked datasets from UCI Data Repository to evaluate the efficiency between DFP-Growth and benchmarked FP-Growth algorithm. Fourth, our proposed algorithm is outperformed the benchmarked algorithm in term of processing time at 1.4 times.

The rest of the paper is organized as follows. Section 2 describes the related work. Section 3 explains the basic concepts and terminology of association rule mining. Section 4 discusses the proposed method. This is followed by performance analysis in section 5. Finally, conclusion and future direction are reported in section 6.

2 Related Work

Mining frequent pattern or association rules from data repository have been a focused theme in data mining research for more than a decade. Apriori [6] was the first algorithm to generate the frequent pattern based on generate-and-test strategy. It employs a level-wise searching, where k -itemsets (an itemset that contains k items) are used to produce $(k+1)$ -itemsets. These k -itemsets are also known as candidate itemsets. As an attempt to optimize and increase the Apriori efficiencies, several variations based on Apriori have been proposed such as AprioriTid and Apriori-Hybrid [7], Dynamic Itemset Counting [8], Direct Hashing & Pruning [9], Partition Algorithm [10], High-Dimension Oriented Apriori [11], Variable Support-based Association Rule Mining [12], etc.

Due to the limitation in Apriori-based algorithms, FP-Growth [5] algorithm without candidate itemsets has been proposed. This method constructs a compact data structure known as FP-Tree from the original transaction database. The main focus is to avoid cost generation of candidate itemsets, resulting in greater efficiency. However, among the main challenge faces by FP-Growth algorithm is, the new FP-Tree must be rebuilt all over again once the original database is updated. Thus, it will directly increase the computational cost due to latest database size. Indeed, this algorithm is obviously unsuitable for incremental pattern mining or dynamic database. It is only fitting for static database. Thus, several variations of FP-Growth algorithm

have been proposed to mitigate this problem. H-mine algorithm [2] uses array-based and trie-based data structures to deal with sparse and dense datasets respectively. PatriciaMine [3] employs a compressed Patricia trie to store the datasets. FPgrowth* [4] uses an array technique to reduce the FP-tree traversal time. However, for any FP-Growth based algorithms, recursive construction of the new FP-Tree is usually will give a great impact in downgrading the performance of proposed algorithms.

3 Proposed Method

Throughout this section the set $I = \{i_1, i_2, \dots, i_{|A|}\}$, for $|A| > 0$ refers to the set of literals called set of items and the set $D = \{t_1, t_2, \dots, t_{|U|}\}$, for $|U| > 0$ refers to the data set of transactions, where each transaction $t \in D$ is a list of distinct items $t = \{i_1, i_2, \dots, i_{|M|}\}$, $1 \leq |M| \leq |A|$ and each transaction can be identified by a distinct identifier TID.

3.1 Definition

Definition 1. (Frequent Items). An itemset X is called frequent item if $\text{supp}(X) > \alpha$, where α is the minimum support.

The set of frequent item will be denoted as Frequent Items and

$$\text{FI} = \{X \subset I \mid \text{supp}(X) > \alpha\}$$

Definition 2. (Canonical Frequent Items). An itemset X is called canonical frequent items if $\text{supp}(X_i) > \alpha$ and sorted in canonical order and

$$\text{CFI} = \{X_i \subset \text{FI}, |1 \leq i \leq k, k = |\text{Items}|, \text{supp}(X_i) > \alpha\}$$

Definition 3. (Canonical Frequent Patternset). An itemset X from FP-Tree is called a canonical frequent patternset if $\text{supp}(X) > \alpha$, and sorted in canonical order and

$$\text{CFP} = \{X_i \subset \text{CFI}, |1 \leq i \leq k, k = |\text{Items}|, \text{supp}(X_i) > \alpha\}$$

Definition 4. (Final Canonical Frequent Pattern). An itemset X is called final canonical frequent pattern if $\text{supp}(X) > \sum \alpha$ and

$$\text{FCFP} = \{X_i \subset \text{CFP}, |1 \leq i \leq k, k = |\text{Items}|, \text{supp}(X_i) > \sum \alpha\}$$

3.2 Algorithm Development

Determine Minimum Supports. Let I is a non-empty set such that $I = \{i_1, i_2, \dots, i_n\}$, and D is a database of transactions where each T is a set of items such that $T \subset I$. An

itemset is a set of item. A k-itemset is an itemset that contains k items. From Definition 1, an itemset is said to be frequent if it has a support count more than α .

Construct FP-Tree. A Frequent Pattern Tree (FP-Tree) is a compressed representation of the frequent itemset. It is constructed by scanning the dataset of single transaction at a time and then mapping onto a new or existing path in the FP-Tree. Items that satisfy the minimum support are only captured and used in constructing the FP-Tree. Fig. 1 shows a complete procedure to construct the FP-Tree.

Canonical Frequent Patterns Generator Algorithm (CFP-Gen)	
1:	Input: Dataset D_n , α_d
2:	Output: Canonical Frequent Patterns CFP_n ,
3:	for ($t_i \in D_n$) do
4:	if ($t_i > \alpha_n$)
5:	$FI \leftarrow t_i$
6:	endif
7:	endfor
8:	for ($t_i \in D_n$) do
9:	if ($t_i > \alpha_n$)
10:	$PP \leftarrow FI \cap t_i$
11:	if ($PP_{\text{items}} \neq FPTree_{pp}$)
12:	$FPTree_{pp,supp} \leftarrow \cup PP_{\text{items},supp}$
13:	else
14:	$FPTree_{supp} \leftarrow \cap PP_{supp}$
15:	endif
16:	endif
17:	endfor
18:	for ($PP_{\text{items}} \in FPTree_{pp}$)
19:	for ($P_i \in FI$)
20:	if ($PP_{\text{items}} \cap P_i \neq 0$)
21:	$FP^{tmp} \leftarrow \text{perm}(PP_{\text{items} \neq P_i}) \cup P_i$
22:	if ($FP \neq FP^{tmp}$)
23:	$FP_{\text{items},supp} \leftarrow \cup FP_{\text{items},supp}^{tmp}$
24:	else
25:	$FP_{supp} \leftarrow \cap FP_{supp}^{tmp}$
22:	endif
23:	endif
24:	endfor
25:	endfor
26:	$CFP_j \leftarrow \cup \text{cano}(FP)$

Fig. 1. CFP-Gen Algorithm

Mining FP-Tree. Once the FP-Tree is fully constructed, the mining process will be started. Hybrid ‘divide and conquer’ method is employed to decompose the tasks of mining desired pattern. The process of mining the canonical frequent itemset and frequent patterns is mentioned in Definition 2 and Definition 3, respectively. The final output is a Final Canonical Frequent Pattern (FCFP) which is stated in Definition 4. Fig. 2 depicts a complete procedure in producing the FCFP from FP-Tree.

Dynamic Frequent Patterns Growth Algorithm (DFP-Growth)	
1:	Input : Dataset D_n, α_a, α_b
2:	Output: Final Canonical Frequent Patterns FCFP
3:	Capture $(D^{\text{newTrans}} \in D_n)$
4:	Execute CFP-Gen
5:	$CFP^{\text{tmp}} \leftarrow CFP^{\text{newTrans}} \cup CFP^{\text{oldTrans}}$
6:	for $(cfp_i \in CFP^{\text{tmp}})$
7:	if $(cfp_i > (\alpha_a + \alpha_b))$
8:	$CFP \leftarrow \cup cfp_i$
9:	endif
10:	endfor
11:	FCFP \leftarrow CFP

Fig. 2. DFP-Growth Algorithm

4 Comparison Tests

In this section, we do comparison tests between DFP-Growth and FP-Growth algorithms. The performance analysis is made by comparing the processing time and number of iteration required. We used three benchmarked and famous datasets in mining frequent patterns. These experiments have been conducted on Intel® Core™ 2 Quad CPU at 2.33GHz speed with 4GB main memory, running on Microsoft Windows Vista. All algorithms have been developed using C# as a programming language and running in .NET Framework 4.

4.1 UCI Dataset from [23]

The three benchmarked datasets were used in the experiments were T10I4D100K, Retail and Mushroom. For DFP-Growth, each dataset was divided into two portions with equal size of transaction. However, for FP-Growth, it was used the original size of transaction in database. Here, the first dataset is T10I4D100K, and was developed by the IBM Almaden Quest research group. It is categorized as a sparse dataset. It consists of 100,000 transactions and 1000 unique items. The second dataset is Retail and based on retail market-basket data from an anonymous Belgian retail store. It has 88,136 transactions and 16,471 items. The third dataset is Mushroom. The dataset

contains 23 species of gilled mushroom in the Agaricus and Lepiota Family. The dataset comprises of 8,124 transactions and 120 items. In the experiments, variety of minimum supports thresholds (α) were employed for each dataset.

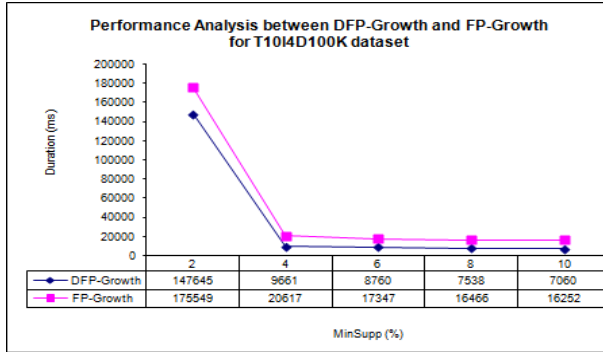


Fig. 3. Computational Performance for Mining T10I4D100K dataset

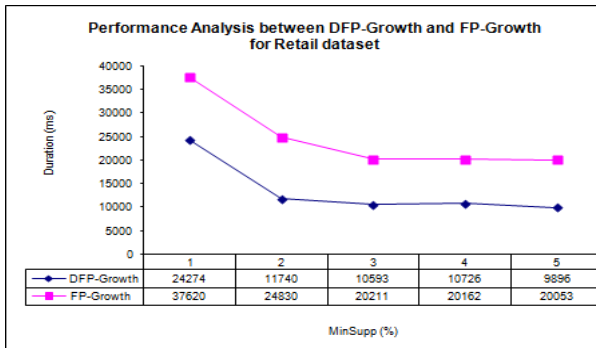


Fig. 4. Computational Performance for Mining Retail dataset

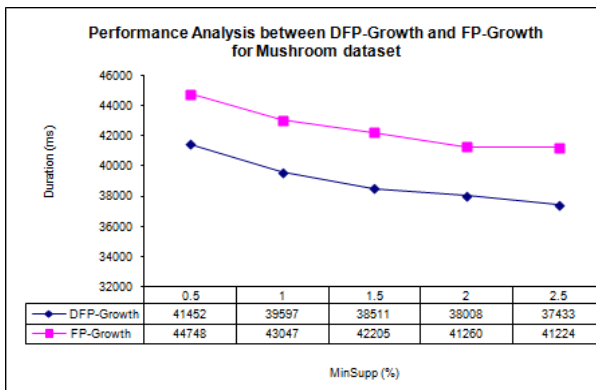


Fig. 5. Computational Performance for Mining Mushroom dataset

Fig. 3 shows the computational performance of both algorithms in mining T10I4D100K dataset. The average time taken to mine the frequent pattern sets using DFP-Growth was 1.83 times faster than FP-Growth. For Retail dataset, the average time taken to mining the frequent patterns using DFP-Growth was 1.36 times faster than FP-Growth as shown in Fig. 4. For Mushroom dataset, the average time taken to mine frequent pattern using DFP-Growth was 1.09 times faster than FP-Growth as presented in Fig. 5. In average, DFP-Growth is outperformed at 1.4 times better than FP-Growth in term of processing to mine the frequent patterns.

5 Conclusion

Mining frequent patterns is a very important study and one of the research themes in data mining. Due to the limitation in Apriori-based algorithms, FP-Growth has been introduced and apparently becomes the benchmarked algorithm in mining frequent patterns. However, the typical problem in FP-Growth is, its FP-Tree data structure must be rebuilt all over again if the number of transactions in the original database is increased. In other words, it is only suitable for static database. Therefore, in this paper we introduce an efficient algorithm called Dynamic Frequent Pattern Growth (DFP-Growth) to mine the frequent patterns from dynamic database. Experiments with UCI datasets show that the DFP-Growth is outperformed FP-Growth up to 1.4 times faster, thus verify its efficiencies.

Acknowledgement. This research is supported by FRGS from Ministry of Higher Education of Malaysia No. Vote RDU 100109.

References

1. Abdullah, Z., Herawan, T., Deris, M.M.: Mining Significant Least Association Rules Using Fast SLP-Growth Algorithm. In: Kim, T.-H., Adeli, H. (eds.) AST/UCMA/ISA/ACN 2010. LNCS, vol. 6059, pp. 324–336. Springer, Heidelberg (2010)
2. Pei, J., Han, J., Lu, H., Nishio, S., Tang, S., Yang, D.: Hmine: Hyper-Structure Mining of Frequent Patterns in Large Databases. In: Proceedings of IEEE International Conference on Data Mining, pp. 441–448 (2001)
3. Pietracaprina, A., Zandolin, D.: Mining Frequent Item sets Using Patricia Tries. In: IEEE ICDM 2003, Workshop on Frequent Itemset Mining Implementations, pp. 3–14 (2003)
4. Grahne, G., Zhu, J.: Efficiently using Prefix-Trees in Mining Frequent Itemsets. In: Proceeding of Workshop Frequent Itemset Mining Implementations, pp. 123–132 (2003)
5. Han, J., Pei, H., Yin, Y.: Mining Frequent Patterns without Candidate Generation. In: Proceeding of the 2000 ACM SIGMOD, pp. 1–12 (2000)
6. Agrawal, R., Imielinski, T., Swami, A.: Database Mining: A Performance Perspective. *IEEE Transactions on Knowledge and Data Engineering* 5(6), 914–925 (1993)
7. Agrawal, R., Srikant, R.: Fast Algorithms for Mining Association Rules. In: Proceeding of 20th VLDB Conference, pp. 487–499. Morgan Kaufmann, Santiago (1994)
8. Brin, S., Motwani, R., Ullman, J.D., Tsur, S.: Dynamic Itemset Counting and Implication Rules for Market Basket Data. In: Proc. ACM SIGMOD, International Conference on Management of Data, pp. 255–264. ACM Press, New York (1997)

9. Park, J.S., Chen, M., Yu, P.S.: An Effective Hash Based Algorithm for Mining Association Rules. In: International Conference Management of Data. ACM SIGMOD, vol. 24(2), pp. 175–186. ACM, San Jose (1995)
10. Hipp, J., Guntzer, U., Nakhaeizadeh, G.: Algorithms for Association Rule Mining – A General Survey and Comparison. In: The Proceedings of SIGKDD Explorations. ACM SIGKDD, vol. 2(1), pp. 58–64. ACM, New York (2000)
11. Ji, L., Zhang, B., Li, J.: A New Improvement of Apriori Algorithm. In: Proceeding of International Conference on Computer Intelligence and Security 2006, pp. 840–844. Springer, Guangzhou (2006)
12. Anad, R., Agrawal, R., Dhar, J.: Variable Support Based Association Rules Mining. In: Proceeding of the 33rd Annual IEEE International Computer Software and Application Conference, pp. 25–30. IEEE Computer Society, Washington (2009)
13. Herawan, T., Vitasari, P., Abdullah, Z.: Mining Interesting Association Rules on Student Suffering Study Anxieties using SLP-Growth Algorithm. International Journal of Knowledge and Systems Science 3(2), 24–41 (2012)
14. Abdullah, Z., Herawan, T., Deris, M.M.: Scalable Model for Mining Critical Least Association Rules. In: Zhu, R., Zhang, Y., Liu, B., Liu, C. (eds.) ICICA 2010. LNCS, vol. 6377, pp. 509–516. Springer, Heidelberg (2010)
15. Abdullah, Z., Herawan, T., Noraziah, A., Deris, M.M.: Extracting Highly Positive Association Rules from Students' Enrollment Data. Procedia Social and Behavioral Sciences 28, 107–111 (2011)
16. Abdullah, Z., Herawan, T., Noraziah, A., Deris, M.M.: Mining Significant Association Rules from Educational Data using Critical Relative Support Approach. Procedia Social and Behavioral Sciences 28, 97–101 (2011)
17. Abdullah, Z., Herawan, T., Deris, M.M.: An Alternative Measure for Mining Weighted Least Association Rule and Its Framework. In: Zain, J.M., Wan Mohd, W.M.B., El-Qawasmeh, E. (eds.) ICSECS 2011, Part II. CCIS, vol. 180, pp. 480–494. Springer, Heidelberg (2011)
18. Herawan, T., Yanto, I.T.R., Deris, M.M.: Soft Set Approach for Maximal Association Rules Mining. In: Ślęzak, D., Kim, T.-H., Zhang, Y., Ma, J., Chung, K.-I. (eds.) DTA 2009. CCIS, vol. 64, pp. 163–170. Springer, Heidelberg (2009)
19. Herawan, T., Yanto, I.T.R., Deris, M.M.: SMARViz: Soft Maximal Association Rules Visualization. In: Badioze Zaman, H., Robinson, P., Petrou, M., Olivier, P., Schröder, H., Shih, T.K. (eds.) IVIC 2009. LNCS, vol. 5857, pp. 664–674. Springer, Heidelberg (2009)
20. Herawan, T., Deris, M.M.: A Soft Set Approach for Association Rules Mining. Knowledge Based Systems 24(1), 186–195 (2011)
21. Herawan, T., Vitasari, P., Abdullah, Z.: Mining Interesting Association Rules of Student Suffering Mathematics Anxiety. In: Zain, J.M., Wan Mohd, W.M.B., El-Qawasmeh, E. (eds.) ICSECS 2011, Part II. CCIS, vol. 180, pp. 495–508. Springer, Heidelberg (2011)
22. Abdullah, Z., Herawan, T., Deris, M.M.: Visualizing the Construction of Incremental Disorder Trie Itemset Data Structure (DOSTrieIT) for Frequent Pattern Tree (FP-Tree). In: Badioze Zaman, H., Robinson, P., Petrou, M., Olivier, P., Shih, T.K., Velastin, S., Nyström, I. (eds.) IVIC 2011, Part I. LNCS, vol. 7066, pp. 183–195. Springer, Heidelberg (2011)
23. Frequent Itemset Mining Dataset Repository, <http://fimi.ua.ac.be/data/>

Analysis on Key Nodes Behavior for Complex Software Network

Xizhe Zhang¹, Guolong Zhao¹, Tianyang Lv^{2,3}, Ying Yin¹, and Bin Zhang¹

¹ College of Information Science and Engineering, Northeastern University, Shenyang, 110819, China

² College of Computer Science and Technology, Harbin Engineering University, Harbin, 150001, China

³ College of Computer Science and Technology, Tsinghua University, Beijing, 100084, China
zhangxizhe@ise.neu.edu.cn

Abstract. It is important to understand software structural complexity and execution behavior in controlling the software development and maintenance process. Difference from previous work which based on structure network built on method association of software, we explore the topological characteristics of software execution behavior based on complex network and model the software execution network based on method invocation sequences. Taking typical open-source software under Linux for example, we build execution network based on the method call relationships, and then explore typical topology measurements of the key node and its adjacent network in software execution network. The result shows that the key nodes with high degree or high betweenness plays an important role in execution process of software system and the execution network can be divided into several levels, which has an important significance for maintenance and quality assurance for software.

Keywords: Complex network, Open-source software, Key node, Dynamic execution, Structural characteristic.

1 Introduction

With the rapid development of the Internet computing technologies such as SOA, cloud computing, Software products are showing a more dynamic, network characteristic [1]. The complexity of software systems gradually change from the complex structure in creation to a complex behavior in execution. The research and understanding of behavioral characteristics in execution network has become a hot issue [2-4].

In recent years, the research and applications of complex network [5-7] theory have penetrated into from physics to biology, from the social sciences to computer science and other disciplines, and it provides a strong analytical methods and tools to analyze the evolution characteristics and behavior patterns of the individuals in complex system. Software as a series of artificial complex systems [8-10] have complex structure and the feature of difficult to control, which gives the software design,

development, maintenance and management many difficulties. Considering the software system as a complex network, from the overall and global perspective to explore and discover structure characteristics, dynamic behavior and evolution of the software, contributes to a comprehensive understanding of the core characteristics of software systems and quantify the complexity of software and performs well in improving the efficiency and quality of software [11].

This paper focuses on the influence of the key nodes to the software execution process, thus providing quality assurance for software. In this paper, taking typical open-source software for example, we take sequence of function calls in the real running environment as a class of software execution network, given the corresponding network model and define the structure statistical indicators such as betweenness, clustering coefficient and degree distribution of network structure. This paper analyzes the behavior of the key nodes with special characteristics, given the structure characteristics of the adjacent subnet of the nodes which have high degree, high betweenness or high clustering coefficient. We make an analysis of the causes and effects of dynamic behavior, so as to find the influence factors for the software reliability.

2 Software Behavior Network Model

Software execution network is a network based on the function calls of a software system execution process, it considers function entity as a network node and the function call $i \rightarrow j$ as a directed edge from node i to node j . The execution network G can be described as triple $\langle V, E, A \rangle$, in which $V = \{v_1, v_2, \dots, v_m\}$ is the set of nodes in the network, $E = \{e_{ij}; i, j = 1, 2, \dots, m\}$ is the set of edges in the network, A is a collection of network attributes, divided into node attributes and edge attributes.

In the execution network, the node strength is the expansion of the node degree in the topology, taking into account both the number of neighbor nodes and the weights between the nodes and neighbors. It is noteworthy that the node execution cost is not the cost of a single function call by itself, but the total consumption time of the function call in the entire implementation process, which has a certain relationship with the intensity of the node.

In order to track the function call relationships of open source software, we use Gnu tool-chain to debug and compile the source code to generate the necessary debug information file `gmon.out` in the run-time, then use the Gnu environment tools `Gprof` to analyze `gmon.out`. Gnu/Gprof is a profile analysis tool for `c/c++` open source projects in Unix-like platform, It can record code-level information such as the calling relationships between functions while the program is running, the number of each function is called, the time consumed by each function. The gcc compiler adds each function a function called "`mcount`" (or "`_mcount`", depends on the compiler or operating system) in the program, which holds a function call graph in memory can be used to find the address of the parent function and the sub function using the form of call stack and to get the function calling relationships, as well as the number of each function call, running time and other information. Gprof is a source code analysis

tool, it can print out the number of each function call, the calling relationships between functions and the time consumption of function call. In order to obtain the desired function call relations from the run-time information, we use specific script file to filter the analysis results of Gprof, extract function call relations to form a sequence file, and finally do network structure calculations for the sequence of function calls based on network model.

Fig. 1 shows the complete experimental procedure. We use Gcc tool chain to debug and compile the source code to generate the necessary debug information file “gmon.out” in running time, then use the Gnu environment tools Gprof to analyze the file in order to extract function call relations to form a sequence file and finally make a function calls based on network model. In addition, in order to simulate the software usage in the real environment, we generate a use case library for each open source software in the experiment, each use case library contains 100 use cases, each use case is a real case used the software which describes the behavior of the user interaction with the software, including a series of continuous button moves.

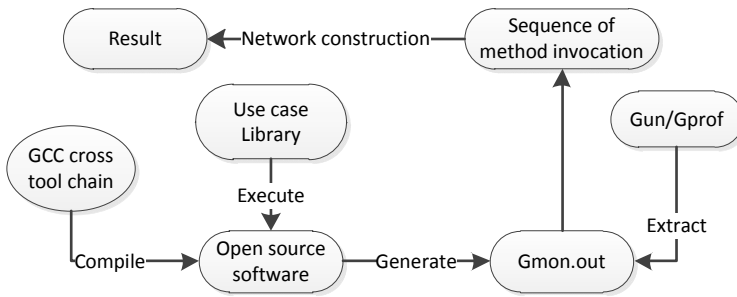


Fig. 1. The experiment process

In this paper, we use the typical open-source software GUI as an object in Linux environment, studying structural features and its dynamics evolution of the execution software which generated in a real environment. Following is the selected open-source GUI software, Dia is a GTK+ based graphical drawing tools, it has a lot of graphical objects in order to achieve flow, UML diagrams, network diagrams, schematics and other graphics rendering work. In addition, users can also write an XML file to add new graphic shape to meet the special requirements of the drawing.

3 Structural Characteristics of Key Node

In software execution network, the key-node can be defined as the node which has a high value in a particular structural indicator. Every node is labeled by the corresponding function name. In this section by extracting these node in execution network and combined with the structure statistical measurement of its neighbor nodes, given an analysis of their behavior in their respective roles in the subnet.

Fig. 2 is an execution network formed by a single Dia use case.

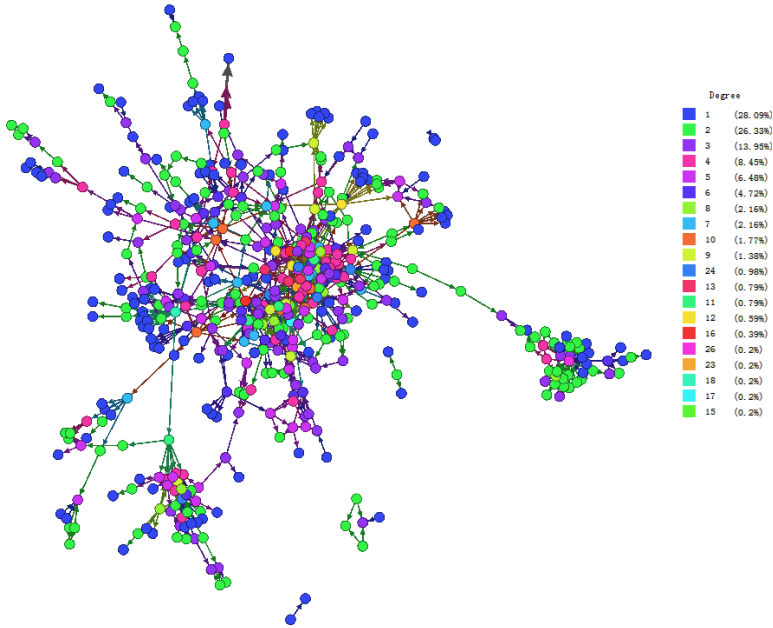


Fig. 2. An example of execution network of Dia

Betweenness is a measure of the centrality or the importance of a node in software execution process, and is normally calculated as the fraction of the shortest paths between any two nodes that pass through the node of interest. Nodes with high betweenness may be regarded as key players that have a highly active role. As a matter of fact high betweenness nodes have an important position in software network. The second-order subnet of the node with the highest betweenness is shown in Fig.3. The id of the node with the highest betweenness is 96, and the name of the function called “object_add_updates”. The function can update the chart when new graphic objects are added to workspace.

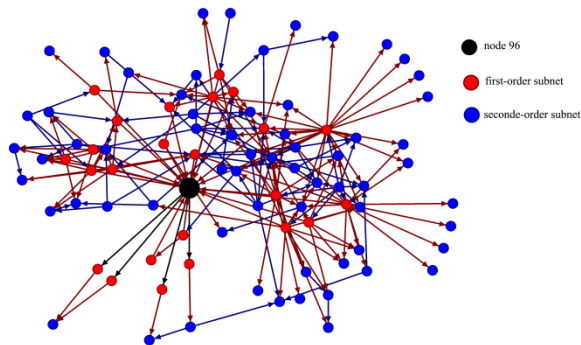


Fig. 3. The first-order and second-order subnet of node 96

The betweenness of node 96 is 460.0. The value is significantly higher than the nodes in its first-order and second-order subnet, which indicates that the node 96 has absolute control to other nodes in the subnet. At same time node 96 has relatively small closeness (1.9), which shows that the node is the heart of the subnet in the topological structure. The correlations between the adjacent nodes are sparse which has led to very small clustering coefficient. Second-order subnet is more likely to a further expansion of the first-order subnet. By observing the values we have found that the nodes in the first-order subnet have a higher value of betweenness than the nodes in respective first-order subnet, which means that the network is divided into multi-layer by node 96. The node 96 plays the role of connecting multiple local subnets. Absolute control in the center and small clustering coefficients has made the node become critical. If you remove the node, the local subnets may become a mess due to loss of contact. Based on these findings, if we control and protection these key nodes in the process of software development and maintenance can effectively improve the software reliability and performance.

The clustering coefficient quantifies how well connected are the neighbors of a vertex in a graph, which has been taken as a signature of the network hierarchical structure. The highest value of clustering coefficient is 0.5 in our software execution network. But the node with the highest clustering coefficient is not unique, which means that many local close call groups are exist in the network. With the node 94 as an example, the statistical indicators of each node in its first-order subnet are shown in table 1.

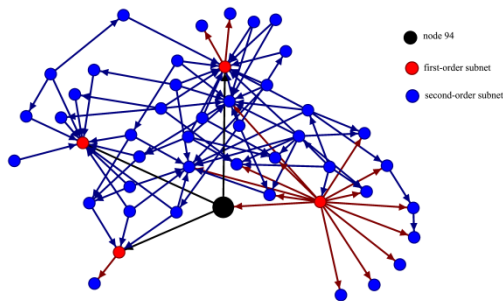


Fig. 4. The first-order and second-order subnet of node 94

In Table 1, node 94 has high clustering coefficient, which means that the first-order subnet of the node 94 is closely related. Its neighbor nodes can be divided into two categories according to the value of betweenness, one category is the node with high betweenness (such as the node 92,101 and 229, they play an critical role in their respective subnets), the other is just the opposite with low closeness (such as node 203, it plays a “marginal” role in their respective subnets), which shows that the node 94 has obvious opposite properties on the relationship of neighbors node. Either connects key node or connects marginal nodes in second-order subnet. In addition, the clustering coefficient of each node in the first-order subnet of node 94 is equal to zero, which indicates that the capability of high clustering no longer exists in their neighbors.

Table 1. Index statistics for 1 step child-network of node 94

nid	name	Betweenness	Closeness	CC	k_{in}	k_{out}
92	create_object_button_press	15.92	2.456	0.0	1	14
94	diagram_add_objects	20.0	2.357	0.5	1	3
101	diagram_modified	190.0	1.714	0.0	14	2
203	diagram_tree	0.0	0.0	0.0	12	0
229	diagram_tree_add_object	68.0	1.75	0.0	4	1

The in-degree and out-degree of node respectively represent the number of edges pointing to the node and the number of ties that the node directs to others. Intuitive point of view, the degree of node denotes its importance in the network. For the software network, in-degree of a function node represents the number of times called by other function. Out-degree is the number of times the function calls other function. The function nodes with high in-degree are at the bottom foundation module in the software, which can provides support for other functions as a public method. The higher the in-degree of the function node is, the more important the function node is in software execution process, such as data process method in the Management Information System. The nodes with high out-degree are like the entrance node in software execution process and also highly dependent on other nodes. The following figure shows the statistical indicators of first-order subnet of nodes with the maximum in-degree and out-degree in Dia software. The node with the highest in-degree is 35 and the node with highest out-degree is 153. The statistical indicators are shown in Table 2.

Table 2. Statistical indicators for node with maximum in-degree and out-degree

nid	name	Betweenness	Closeness	CC	k_{in}	k_{out}
35	dia_ps_renderer_get_type	0.0	0.0	0.5	26	0
153	ddisplay_canvas_events	0.0	2.589	0.0	0	24

As can be seen from the above table, the betweenness, closeness and out-degree of node 35 is 0, which means that they are mostly called by other functions and can only play a role of basic function. At the same time the largest clustering coefficient of the nodes is 0.5, which shows that the contacts between the functions who call them are closely related. Its neighbors are likely to implement a specific function in the correlation function. The largest out-degree of the nodes is 153. The betweenness, clustering coefficient and in-degree is equal to zero, and closeness is only 2.589(relative to the larger sets), these indicators show that they are upper abstraction and need to call a number of other functions to help achieve a specific function. The role and behavior of the node 153 is stands in stark contrast to the node 35. It is more or less related to that they are the largest out-degree and in-degree in the network. Because of low betweenness the two nodes have a common characteristic which is marginal effect. A high value of out-degree means high dependence on the other nodes. Therefore, to encourage reuse of loosely coupled modular design principle, the critical point of the node with big value of out-degree is mostly in the same module. The first-order subnet closely linked and forms a small group.

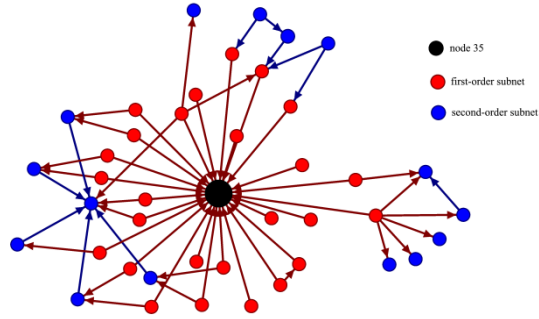


Fig. 5. The first-order and second-order subnet of node 35

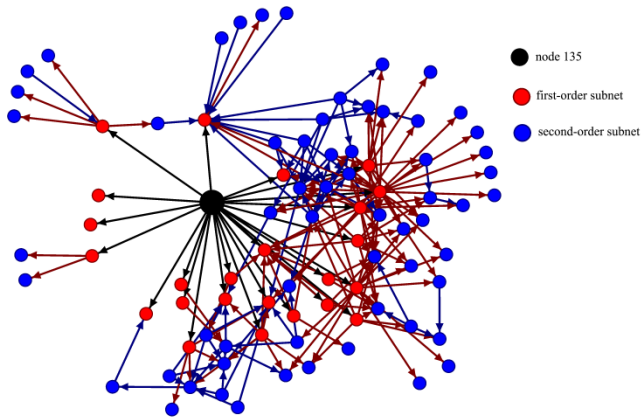


Fig. 6. The first-order and second-order subnet of node 135

4 Conclusions

The study of the structural characteristics of software system in dynamic execution can more essential, in-depth understand the complexity of such artificial systems, which has an important practical significance in measuring software, improving the study can improve architecture and securing the quality. In this paper, based on open-source software, we use the complex network theory to do network modeling and structural characteristics definition of function call relations during the software execution. By simulating the use of software in a real environment, we study the structure characteristics of the type of network, relevant conclusions reveal the special nodes with high degree or high betweenness play an important role in normal execution of software system. The execution network can be divided into several levels, which has an important significance for operation, maintenance and quality assurance of software.

Acknowledgements. This work is sponsored by the Natural Science Foundation of China under grant number 60903009, 61073062, 61100028, 61100090, 61100027 and the Fundamental Research Funds for the Central Universities under grant number N090104001.

References

1. Gonzalez, M.C., Barabasi, A.L.: Complex networks – From data to models. *Nature Physics* 3, 224–225 (2007)
2. Zhang, H.H., et al.: Using the k-core decomposition to analyze the static structure of large-scale software systems. *J. Supercomput.* 53, 352–369 (2010)
3. Moyano, L.G., Mouronte, M.L., Vargas, M.L.: Communities and dynamical processes in a complex software network. *Physica A* 390, 741–748 (2011)
4. Myers, C.R.: Software systems as complex networks: Structure, function, and evolvability of software collaboration graphs. *Phys. Rev. E* 68, 046116 (2003)
5. Fortunato, S.: Community detection in graphs. *Physics Reports* 486, 75–174 (2010)
6. Liu, Y.Y., Slotine, J.J., Barabasi, A.L.: Controllability of complex networks. *Nature* 473, 167–173 (2011)
7. Jenkins, S., Kirk, S.R.: Software architecture graphs as complex networks: A novel partitioning scheme to measure stability and evolution. *Information Sciences* 177, 2587–2601 (2007)
8. Li, H.A., Li, B.: A Pair of Coupling Metrics for Software Networks. *Journal of Systems Science & Complexity* 24(1), 51–60 (2011)
9. Zhang, H.H., et al.: Using the k-core decomposition to analyze the static structure of large-scale software systems. *Journal of Supercomputing* 53(2), 352–369 (2010)
10. Savic, M., Ivanovic, M., Radovanovic, M.: Characteristics of Class Collaboration Networks in Large Java Software Projects. *Information Technology and Control* 40(1), 48–58 (2011)
11. Cai, K.Y., Yin, B.B.: Software execution processes as an evolving complex network. *Information Sciences* 179(12), 1903–1928 (2009)

Webpage Information Hiding Algorithm Based on Integration of Tags and Data

Junling Ren¹ and Li Zhang²

¹ School of Information Management, Beijing Information Science & Technology University, 100192, Beijing, China,

² School of Automation, Beijing Information Science & Technology University, 100192, Beijing, China

renjunling@bistu.edu.cn, zzzrl@163.com

Abstract. To investigate the webpage information hiding technology and especially improve its hiding capacity, the webpage structure and its elements were analyzed. Then according to the structural characteristics that the web was composed of the tags and data, the integration strategy using both the tags and data to implement information hiding was established, and accordingly a hiding information algorithm was proposed. This algorithm improves the hiding capacity of the webpage information hiding, and also an attempt is made for the method of taking the webpage as the channels to transmit the secret information.

Keywords: information hiding, webpage, tags, data, integration.

1 Introduction

Webpage information hiding technology takes the webpage as the carrier and transmission method to conduct the information hiding [1-3] and it is the combination of webpage and information hiding technology. By means of changing the important information transmitting in the webpage into the secret information, the security performance of webpage delivery can be improved. Simultaneously, because the information in the network is multitudinous and updates very fast, the information hiding is more covert and secure.

From the structure, the webpage mainly consists of data and HTML tags. So there are two methods to hide the information. One is hiding the information in the tags, and the other is hiding the information in the data elements. Now the researches are mainly made on the webpage tags [4-9]. This method is easy to realize. But since the tag file is rather small, the hidden capacity is very limited. Compared with that of the information hiding method based on the tags, the hiding capacity of the method based on data is much more. In view of the characteristics above, the integration information hiding strategy is established, which combines the strategy based on the tags with the strategy based on the data together, and the webpage information hiding algorithm is proposed.

2 Webpage Information Hiding Model Based on the Tags

The webpage is formed by analyzing the HTML language through the web browser. While the webpage is transmitted in the network, the source code can not be seen directly, what is been viewed is the contents parsed by the browser. On basis of this feature, information hiding for webpage is mainly by means of modifying the syntax and tag of the HTML files according to the secret information to make the parsed contents same before and after modification. The specific model is shown in Fig. 1.

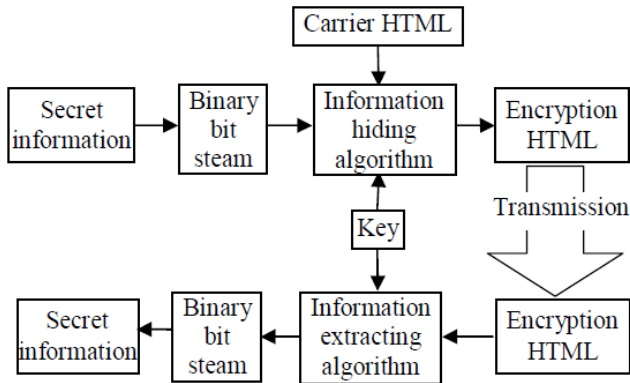


Fig. 1. Webpage information hiding model based on tags

Known from Fig. 1, the basic flow of webpage information hiding is described as below: transform the secret information into a binary bit stream; select a certain syntax or tag of HTML which has at least two different forms and is exactly the same content after being parsed, and set 0 and 1 respectively on its two different forms, give out the hiding algorithm; modify the HTML file style according to 0 and 1 in the binary bit stream of the secret information, and information hiding is accomplished, and the modified HTML file is just the file hiding the secret information, that is, the encryption HTML. At the same time, taking the robustness of the algorithm into account, start of secret information hiding location can be randomly selected by keys.

The process of extracting the secret information is to convert received HTML file to 0 and 1 according to the pre-set different forms of syntax, and the binary bit stream of the secret information is obtained, and then it is translated into the secret information.

3 Webpage Information Hiding Model Based on the Data

The webpage contains not only tags but also the multimedia data of the text, image, audio, animation and video which are also called as webpage data elements. Some multimedia data of the webpage is selected as the hiding carrier in the information hiding method based on the webpage data. There is no change to vision after and before hiding secret information. The specific model is shown as Fig. 2.

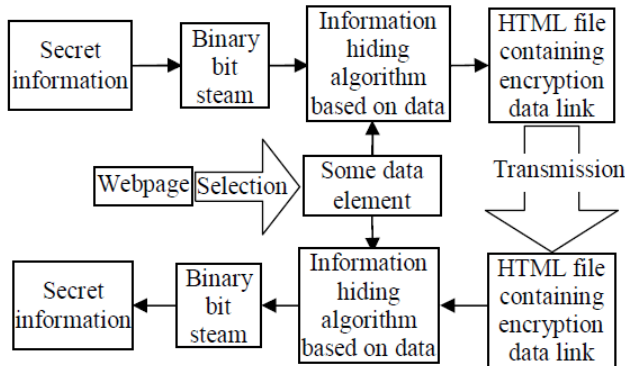


Fig. 2. Webpage information hiding model based on data

The procedure of the information hiding based on the data is described as below: transform the secret information into binary bit stream, and select some data element of the webpage as the carrier. Then select the corresponding algorithm for the type of the carrier data to realize the information hiding.

Since only the linking locations of the multimedia data elements are stored in the webpage source files, the webpage source files don't have to be modified in the information hiding based on the data. And just the data which is appointed to hide information by sender and receiver is needed. The link of the data element carrying secret is found in the webpage received. And thus fix a position on the data element carrying secret. Then use the extraction algorithm based on the data to accomplish the extraction of the secret information.

4 Webpage Information Hiding Algorithm Based on the Integration of the Tags and the Data

4.1 Integration Strategy of Information Hiding Combining Tags and Data

As we all known, the webpage information hiding algorithm based on tags is easy to realize. But the hiding capacity is very limited and the secret information of a large quantity can't be accomplished. This is also the bottleneck to practical application for the algorithm. By contrast, a large quantity of the secret information can be hidden by means of the information hiding algorithm based on the data. But because the sender and receiver should appoint the carrier data elements, the algorithm is lack of flexibility. Simultaneously, the hiding algorithm based on the data elements is more complex than that based on the tags. When the quantity of secret information is not large, the achieved efficiency will be reduced if the information hiding method based on the data is used.

On the basis of the analysis above, an integration strategy of the webpage information hiding based on both the tags and data is proposed. It is specifically described as bellow: when hiding the secret information, for the selected carrier webpage, it is automatic to judge if the hiding capacity is enough to hide the secret

information after the secret key and secret information are input. If it is enough, then the secret information is directly hidden into the carrier webpage. If the hiding capacity of the carrier webpage is not enough, then jump to the information hiding part based on the data. The user will select the carrier data and hide the secret information. After hiding the secret information, the link of the data carrying the secret information is saved in the webpage. When extracting the secret information, secret key is input into the selected carrier webpage. If the extraction information is the link of some data element, then the corresponding data element is found, and gain the secret information by means of the extraction algorithm of this type of data.

4.2 Basic Process of the Webpage Information Hiding Algorithm Based on the Integration of the Tags and Data

4.2.1 Basic Hiding Algorithm

The webpage information hiding algorithm based on integration of tags and data is illustrated in Fig. 3, and the specific procedure is described as below:

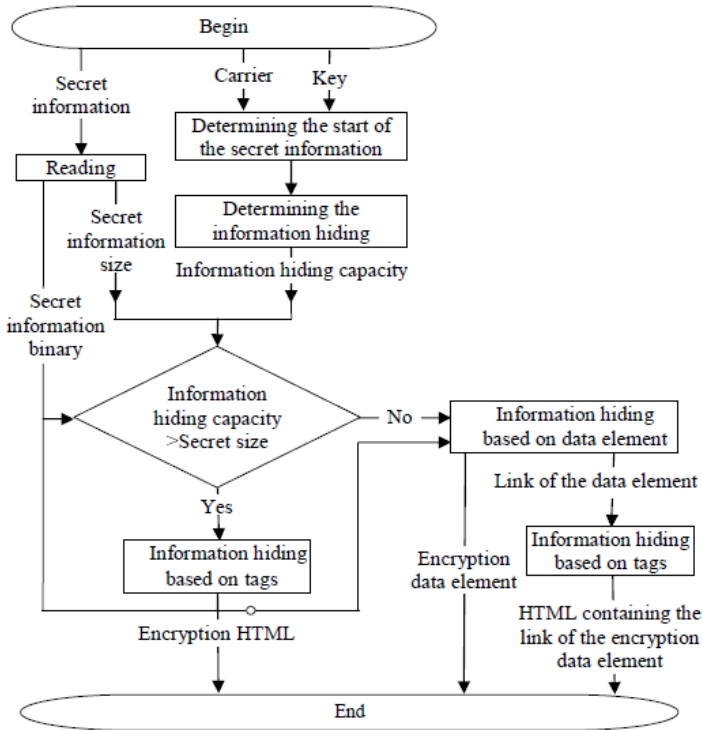


Fig. 3. Webpage information hiding algorithm based on integration of tags and data

Step1: According to the secret key, determine the initial hiding position of the secret information in the carrier webpage files.

Step2: Compute the hiding capacity of the webpage HTML files.

Step3: Compare the hiding capacity with the secret information.

If the hiding capacity is more than the secret information, the secret information is directly hidden into the webpage by mean of the algorithm based on the tags. Otherwise go to the next step to make the secret information hiding based on data.

Step4: Hide webpage information based on the data element.

Take the multimedia data element in the carrier webpage files as the information hiding carrier.

According to the type of the data element, select the corresponding hiding algorithm, and hide the secret information into the data carrier by means of the information hiding algorithm based on data.

Hide the link of the data carrier in which the data elements are hidden, that is, according to the information hiding algorithm based on the webpage tags, the link of the data carrier taking the secret information is hidden into the carrier webpage HTML files.

4.2.2 Basic Extraction Algorithm

Corresponding to the hiding process, the extraction algorithm of the webpage information hiding which is based on the integration of tags and data is illustrated in Fig. 4. The specific follow is as bellow:

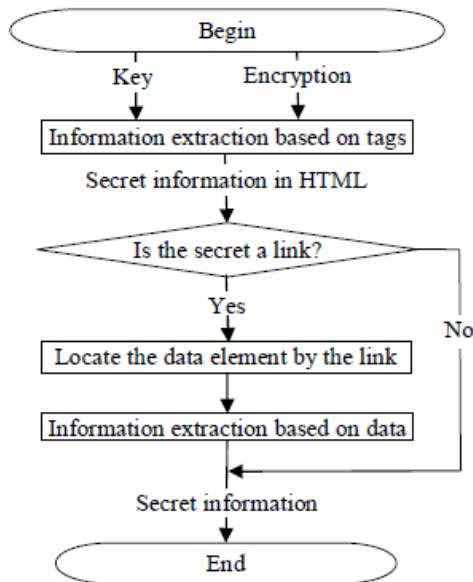


Fig. 4. Webpage information extraction algorithm based on integration of tags and data

Step1: Extract the secret information based on the webpage tags.

According to the extraction algorithm of the information hiding based on the tags, the secret information is gained from the webpage HTML files.

Step2: Analyze the secret information in the webpage HTML files.

If the secret information extracted from the webpage HTML files is the link of the data element carrier, then go to Step 3 and take the extraction of the secret information from the data carrier. Otherwise, the extraction information is just the secret information, and the process of the extraction ends.

Step 3: Extract the secret information from the data carrier

Locate the data element carrier

According to the link of the data element carrier, find the data element.

Extract the secret information

According to the type of the data element, the corresponding extraction algorithm is used to obtain the secret information from the data element files.

5 Experiments

In order to verify the effectiveness of the algorithm and analyze related performance, we have studied its invisibility, robustness and hiding capacity. The method based on the tags dictionary in the literature [10] is applied for the webpage information hiding algorithm based on the tags. Taking the universality of the data elements and the performance of their information hiding algorithm into account, here JPEG image widely used is taken as the webpage data element to have a test.

5.1 Invisibility

In this system, the secret information is usually hidden in the picture of the webpage, and it is more invisible than hiding the secret information in the webpage. Because the carrier picture is randomly selected in the webpage, the invisibility of the system is increased.

The link of the picture carrier is covered in the webpage after information hiding. Normally the link of the picture is very small in size, so when the webpage is transmitted in the network, the visitor can't find any difference when visiting the webpage by the browser. Simultaneously since the link of the picture carrier is very small, the hiding density change higher by means of adjusting the contents of the dictionary file. Even if the source code is viewed, it is also less likely to arouse the visitor suspected.

5.2 Robustness

On the one hand, the robust of the algorithm which selects the image as the data carrier in this paper is much better. Even the some part of the image is modified or it is under mosaic attack, most of the secret information is able to be recovered. The hiding algorithm based on the uppercase and lowercase of the tag dictionary is used in

webpage information hiding, and even if the contents are modified, the secret information can be also recovered.

On the other hand, under the usual circumstances, the contents of the webpage can't be modified when transmitting in the network. Even the attacker intentionally tampers with the contents of the webpage, since the data carrier is selected at random, the probability of its being attacked will be greatly reduced, and thereby the anti-aggressive performance of the algorithm is increased. So the robustness of the system is rather better.

5.3 Information Hiding Capacity

The information capacity comparison between different webpage information hiding algorithms is shown in Table 1. In Table 1, Meth.1 is the information hiding algorithm based on the general tags, and Meth.2 stands for the hiding algorithm based on the integration of the tags and data elements.

From Table 1, the hiding capacity of the hiding information algorithm based on the integration of the tags and data elements increases greatly than the webpage information hiding based on the tags. Thereby the webpage information hiding technique can be widely practical applied.

Table 1. Information hiding capacity based on different webpage information hiding algorithms

Meth.	Hiding ratio in HTML file	Hiding ratio in picture	Final webpage hiding ratio
Meth.1	0.6-0.8%	/	0.6%-0.8%
Meth.2	0.6-0.8%	15%	15.6%-15.8%

6 Conclusions

From the structure of the webpage, the integration strategy of the hiding information which is based on the tags and the data is proposed in this paper, and webpage information hiding algorithm is realized. On basis of the tags characteristic, the secret information of little quantity is hidden in the HTML files of the webpage, and the large quantity secret information is hidden in the data element of the webpage, and then the link of the carrier data elements is covered in the HTML files of the webpage. By means of the experiments, compared with the algorithm based on the tags, this algorithm is more invisible and robust. The hiding capacity of the algorithm increases greatly, and the bottleneck of the webpage hiding information algorithm based on the tags is overcome. An infeasible scheme which takes the webpage as the covert channel is provided.

Acknowledgements. The research is sponsored by Beijing Municipal Commission of Education Science & Technology Development Program under grant No.KM201010772017, No. KM201110772012, supported by Beijing Municipal Administration of college and university middle-aged backbone teacher training

program under grant No.PHR201108252, and supported by Beijing Municipal Administration of college and university academic innovation team program under grant No.PHR201107133.

References

1. Petitcolas, F.A.P., Anderson, R.J., Kuhn, M.G.: Information hiding—A survey. *Proceedings of the IEEE, Special Issue on Protection of Multimedia Content, USA* 87(7), 1062–1078 (1999)
2. Moulin, P., O’Sullivan, J.A.: Information-theoretic analysis of information hiding. *IEEE Transactions on Information Theory* 49(3), 563–593 (2003)
3. Provos, N., Honeyman, P.: Hide and seek: an introduction to steganography. *IEEE Security & Privacy* 1(3), 32–44 (2003)
4. Yong, S.: A scheme of information hiding based on HTML document. *Journal of Wuhan University (Nature & Science Ed.)* 50(S1), 217–220 (2004)
5. Long, Y.-X.: Model of information hiding based on HTML tags. *Application Research of Computers* 24(5), 137–140 (2007)
6. Huang, H.-J., Tan, J.-S., Sun, X.-M.: On Steganalysis of Information in Tags of a Webpage Based on Higher-order Statistics. *Journal of Electronics & Information Technology* 32(5), 1136–1140 (2010)
7. Mi, Z., Yong, F.: Information Hiding Model Based on Attributes Scheduling of XML/HTML Label. *Communications Technology* 43(5), 106–108 (2010)
8. Wu, D.-S.: Research and Implementation of the Information Hiding Technology Based on Hypertext Markup Language. *Journal of Software Guide* 10, 66–67 (2011)
9. Sun, L., Zhang, D.-S., Chen, P.: Algorithm Study of Information Conceal Based on Many Web Pages. *Journal of Value Engineering* 30(23), 129–130 (2011)
10. Ren, J.-L., Wang, C.-Q.: A Webpage information hiding algorithm based on tag dictionary. In: 2012 International Conference on Computer Science and Electronic Engineering (CSEE 2012), pp. 546–550 (2012)

Location Method of Underground Pipeline Monitoring Point Based on Cavity Detection

Wei Zhu¹, Ping Sun^{1,2}, Ying nan Ma¹, Rui Song², Shi wei He², and Ke Hui Liu¹

¹ Beijing Research Center of Urban System Engineering, Beijing 100089, China

² School of Traffic and Transportation, Beijing Jiaotong University, Beijing 100044, China
zhuweianquan@126.com

Abstract. Aiming at the situation of major underground pipeline collapse accident to happen occasionally, we analyzed the cause of underground cavity, establishing underground pipeline risk index system based on cavity, and then gave the risk classification method of the surrounding environment of underground pipeline. Based on the radar detection data, we used the expert scoring method to calculate the weight of the index at all levels, and assigned its basic indicators. For the first time, we had a risk assessment on the weak section of the underground pipeline, to identify the underground pipeline monitoring points preliminarily, and established the maximum coverage location model, and had a preliminary discussion of cavity monitoring points location around the underground pipeline, to make full use of the present resources, and to realize dynamic monitoring on the cavity evolution around the underground pipeline, and to predict and prevent major underground pipeline accidents effectively and scientifically.

Keywords: cavity, underground pipeline, detection, monitoring, maximum coverage location model.

1 Introduction

Urban water supply underground pipeline is the important urban infrastructure, which is a necessary condition of the material base of the survival and development of the city, the main artery maintaining the normal operation of city, and with the name of "the urban lifeline". Urban underground pipeline is bearing the supply of urban life, public welfare undertakings and the security of the industrial production, which is inseparable with People's daily life and one part of the urban public utilities, and is also one of the most important infrastructures for city constructions.

Along with the quick city development, urban construction changes daily, the urban underground pipeline construction accelerates constantly, with urban underground pipeline management problems more and more complex, especially having complex effects on the ground and underground construction, and the smooth operation and safety of underground pipeline [1-4]. Underground cavity problems caused by all reasons can't be ignored, which may cause the ground collapse, leading to the water supply, gas, water drainage, heating, communication and other pipeline damages, and even makes the ground building damaged and extremely easily causes losses of lives

and property [5-10]. The surrounding environment especially the underground cavity, and the urban underground pipeline operation influences each other, major pipeline accidents occur frequently. At present, researches on the underground pipeline leakage, water pollution, and other aspects of the monitoring method are more comparatively [11-14]. There is a great important practical significance to research on monitoring method of cavity around underground pipelines, cavity detection and to take some measures to control the weak links of the pipeline.

2 The Formation and the Classification of the Underground Cavity

The laying methods of underground pipe can be roughly divided into three kinds of ditch buried, on buried, tunnel. The stress state and the vertical soil pressure on the tube top of the underground pipeline directly depends on the laying method of the buried pipes. When cavity exists, the extent of damage is closely related to pipe laying method.

The formation of the underground cavity is very complicated, according to relevant data analysis and some practice summary, the formation of the urban underground cavity can be summarized by the following several aspects:

- (1) Geology, groundwater function. Such as the 4th soil (including clay, sandy soil, etc) thickness are mostly for dozens of meters or more than hundreds of meters, due to long geological movement to form the fissure, with the erosion of groundwater, fracture becomes more and more big hole or large area of loose soil.
- (2) Long-term vibration. With city's large-scale machinery, large vehicles passing on the road so that the strata under the road is in the long-term vibration state, together with the decline of underground water level, and the pressure drops of underground space fissures and voids, and with the crack expansion and the "floating" phenomenon, which leads to the formation of cavity under the road.
- (3) The existence of the civil air defense projects underground. Due to the underground civil air defense projects constructions already for a long time, there is some section lack of management yearly to cause failure and landslides, with the action of water, making a lot of mud into the hole to cause a large area of hole above the air defense projects.
- (4) Infrastructure construction effect. Such as subway tunnels, deep foundation pit construction, etc. With the stratigraphic disturbance of construction plenty of groundwater seeps. A large number of sediment in its upper or surrounding level of the loose soil around them is brought away, and which gradually forms the hollow.
- (5) The impact of pipeline construction. Pipeline excavation and construction backfill is in the loose state, the trenchless block grouting is lax, so that groundwater flow along the pipeline. It brings away the upper loose sediment strata and then the cavity forms.
- (6) The rolling of road construction is loose, forming a loose layer. For the long-term vibration and rolling, the soil under the roadbed with rigidity gradually sinks to the formation of voids; and then flexible roadbed of the local area forms collapse.

- (7) The underground pipeline damage. Leaking caused by water supply, drainage pipeline, which flushes the soil around the pipe and then hollow holes form nearby.

At the same time, according to the formation reasons of the cavity, cavity can be divided into three types: cavity, watering cavity and loose soil layer, whatever the form is, they all have a common characteristic, which is the different scale and irregular shape, having no obvious trend and extension.

3 The Construction of Risk Evaluation Method Based on Cavity Detection

Underground pipeline system is intricate, and the gas, water supply, power, transportation, roads and bridges, often centers in some the same area, the interaction and impact among the structure of all kinds of lines should not be ignored, and artificial excavation, neighbor structures such as the construction of all kinds of engineering activities have become an important factor influencing the safety of the lifeline system.

Table 1. Underground pipeline risk index system

Disaster level	First-level Index	Second-level Index	Third-level Index	
Running Environment Risk G	Pregnancy disaster environment V ₁	Natural environment V ₁₁	Soil category V ₁₁₁ Road level V ₁₂₁ Load level V ₁₂₂ Traffic volume V ₁₂₃ Pipeline position V ₁₂₄	
		Social environment V ₁₂	Road surface level V ₁₂₅ Road surface category V ₁₂₆ Drainage infrastructure V ₁₂₇	
		Vulnerability of environmental disasters V ₂₁	Pipeline category V ₂₁₁ Use function V ₂₂₁ Pipeline material V ₂₂₂ Laying methods V ₂₂₃ Construction methods V ₂₂₄	
		Disaster bearing the bearing ability of the hazard-affected body V ₂₂	Interface form V ₂₂₅ Pipeline diameter V ₂₂₆ Buried depth V ₂₂₇ service period V ₂₂₈ Pipeline in good condition V ₂₂₉	
		Disaster causing factor V ₃	Geology V ₃₁	Area V ₃₁₁ Depth V ₃₁₂ Soil properties V ₃₁₃

Table 1. (continued)

Disaster level	First-level Index	Second-level Index	Third-level Index
Running Environment Risk G	Disaster causing factor V ₃	Hydrogeology V ₃₂	Water violations occurrence probability V ₃₂₁
		Biology V ₃₃	Construction effects V ₃₃₁
		Disaster types V ₄₁	Accident type V ₄₁₁
	Disaster situation V ₄	geographical position V ₄₂	The number of pipelines affected V ₄₂₁ The distribution formation of pipeline affected V ₄₂₂ Position relationship V ₄₂₃ Distance relation V ₄₂₄

The complex environment risks of underground pipeline mainly come from four aspects: pregnancy disaster environment, disaster bearing body, disaster causing factor, and the disaster situation factors. The influencing factors of the complex environment risk index of the underground pipeline are more, combined with the data analysis of pipeline accidents, building a risk index system of underground pipeline, shown in table 1. Using the method of expert scoring, the weight index at various levels is scored, with the total score, based on table 2 we can have an evaluation on pipeline risk status, and then we can have a key monitoring on higher-risk pipeline section(risk levels reached A and B class).

Table 2. The assessment standards of the risk source impact on pipeline

G	$G \geq 3$	$3 > G \geq 2$	$2 > G \geq 1$	$G < 1$
Evaluation grade	A	B	C	D

Explanation:

A—Such risk source has seriously affected the safe operation of the pipeline, we must take the technical measures necessary to handle in time.

B—Such risk source has a major influence on pipeline, we suggest having accurate drilling in imperfect area if conditions allowed, taking corresponding measures according to detailed survey results. The tracking observation intervals of the unconditional drilling should be 1 to 3 months.

C—Such risk source has certain effect on pipeline, the imperfect area won't seriously influence the pipeline safety. We should have a 3-6-month-interval tracking observation.

D—Such risk source to pipeline has a little influence and can be temporarily ignored.

4 Monitoring Site Selection Principle

(1) Spreadability

For fixed target area, coverage was measured by the valid range monitored by activation nodes of the sensor network. In the target area there is always some

unactivated sensor and corresponding unapped area, namely cover holes. When the energy of the sensor node is run out, or the sensor nodes randomly arranged is less than the amount required, this kind of circumstance appears. The spreadability directly decides the data integrity and accuracy or the integrity degree provided by the sensor network.

(2) Connectivity

There is different description according to different needs and applications . For example, the node connectivity is usually referred to the adjacent number connecting to a network node; All the connectivity is referred to a connected path between any

two nodes in the network; Network connectivity $c = \frac{\sum_{i=1}^N \sum_{j=1}^N y_{ij}}{N^2}$, N is the network

nodes in total, $y_{ij}=0$ means nodes can't be connected, $y_{ij} = 1$ shows nodes can be connected. On certain occasions the meaning of connectivity can be weakened. As in some wireless sensor network, there is no business requirement among common nodes, this time the connectivity is referred to whether there is a routing between any activated network nodes.

(3) The relationship between coverage and connectivity

In consideration of the coverage, there is a certain contact between the coverage and connectivity. When the node communications coverage is two times larger than the range of node sensor, we think the node coverage problem can contain cover connection completely, so when the communication distance and perception distance satisfy certain relations, coverage can ensure connectivity, and vice versa.

5 The Maximum Coverage Location Model

Coverage location problem is to use the least amount of sensors covering (monitoring) all the cavities around underground pipeline (hereinafter referred to as cavities), the relative importance of different cavities are same. In practice, because of some economic and technological and administrative reasons, we can't guarantee that all the cavities can be monitored, so this sensor location problem becomes a problem when given a number of the budget, we can monitor the largest number of cavities around the underground pipeline.

Assuming that the midpoint of the weak underground pipeline section is cavity point, the possible leakage events set of the pipe section around the cavity are H , $l_{ij}^h = l_i(t_j^h)$ shows when the leakage $h \in H$ happens, and the sensor of cavity j firstly monitoring the average soil moisture surrounding the underground pipeline abnormal, the leakage volume of the section around the i th cavity. Assuming that the allowable leakage volume of the section around cavity i is ordered as l_i^{\max} , if $l_{ij}^h \leq l_i^{\max}$, we think as for leakage event h , we can have a monitoring in cavity j with a sensor in cavity i . Sensor optimization problem can be summed as the maximum coverage problem:

$$\max w = \sum_{i \in I} \sum_{h \in H} y_i^h \quad (1)$$

$$s.t. \sum_{j \in N_i^h} x_j - y_i^h \geq 0, \forall h \in H, N_i^h \neq \emptyset, i \in I \quad (2)$$

$$\sum_{j \in J} x_j = q \quad (3)$$

$$x_j \in \{0, 1\} \forall j \in J \quad (4)$$

$$y_i^h \in \{0, 1\} \forall i \in I, h \in H \quad (5)$$

$$N_i^h = \{j \in J \mid l_{ij}^h \leq l_i^{\max}\} \forall h \in H, i \in I \quad (6)$$

$$N_i^h = \emptyset \text{ if } l_{ij}^h > l_i^{\max} \forall j \in J, i \in I \quad (7)$$

Among them, W shows, as for the leakage events H , the times of the cavity can be monitored; I shows cavity set; J shows the cavity set around the candidate sensors, q shows the sensor budget number.

The objective function (1) maximize the number of the cavity that q pieces of sensors monitoring on the leakage events. Constraints (2) show when the leakage events l occurs, $y_i^h = 1$ when and only when at least one sensor can monitor the cavity. Constraints (3) limits the number of the sensor equal to q . Constraints (4) shows x_j is the decision variable, if a sensor is set in cavity j , then $x_j = 1$, otherwise $x_j = 0$, constraints (5) limits y_i^h as a binary variable, according to the constraint (2) we decides its value. Constraints (6) show N_i^h as the candidate sensor set that can monitoring joint i in the current leakage event. Constraints (7) shows on the current leakage incident l , if every candidate sensor node j satisfies with $l_{ij}^h > l_i^{\max}$, the N_i^h is the empty set.

6 Conclusion

Based on the study of the cause of cavities around the underground pipeline, with the cavity detection data, we have a recognition of the weak section of underground pipeline, selecting the midpoint of the weak pipeline section as the candidate installation point, considering the maximum location model of the monitoring points covering the cavity when economic conditions allowed (only q pieces of sensor), and we have a preliminary discussion on the monitoring location method of the cavity around the underground pipeline, providing a certain basis for monitoring the cavities around the underground pipeline scientifically.

References

1. Du, X., Liu, R., Liu, Y., et al.: Study of Geological Disaster Prevention in Beijing Natural Gas Network. *Chengshi Ranqi* 411(5), 11–17 (2009)
2. Zhou, J.: Gas transmission pipeline corrosion situation testing and evaluation. *Oil-Gasfield Surface Engineering* 27(8), 49–50 (2008)
3. Yan, J., Wu, B., Wang, G., et al.: Estimation of Fault with Classified Evaluation Method of for Ground Weighting Pipeline. *Municipal Engineering Technology* 27(2), 125–128 (2009)
4. Zhang, Z.: The discussion of the pressure measuring point optimum arrangement of water supply network. *Water Technology* 1(1), 59–61 (2007)
5. Zhou, S., Xu, S., Liu, J.: Study on Optimal Location of Flow Measurement Station in Urban Water System. *Journal of Shaoyang University (Natural Science)* 2(2), 86–88 (2005)
6. Huang, Y.D.: Optimal sensor placement in water distribution systems, Hang zhou (2007)
7. Zhou, S.K., Xu, S.R., Liu, J.X.: Optimal Location of Flow Measurement Station of Urban Water Scada System. *Journal of Zhuzhou Institute of Technology* 19(4), 118–121 (2005)
8. Wang, S.W., Li, S.P., Liu, X.P.: The internal flow monitoring technology and its application of the water supply network. *Water & Wastewater Engineering* 35(10), 107–111 (2009)
9. Guo, J., Xin, K.L.: The optimum arrangement of water quality monitoring points of the water supply network. *Public Utilities* 21(4), 21–24 (2007)
10. Shen, H.: Research on Optimal Configuring of Power Quality Monitors in Distribution Network, Qing Dao (May 2010)
11. Sun, P., Zhu, W., Xing, T.: The application of “3S” technology in underground pipeline emergency management, pp. 175–179. IEEE Computer Society, Wuhan (2010)
12. Sun, P., Zhu, W., Xing, T.: Survival Environment Risk and Accident Mechanism Analysis of Underground Pipeline, pp. 153–160. Beijing Institute of Technology Press, Beijing (2010)
13. Sun, P., Zhu, W., Zheng, J.: Failure Risk Analysis and Control of Urban Flow Pipe. In: *Advances and Experiences with Pipelines and Trenchless Technology for Water, Sewer, Gas, and Oil Applications*, Shang Hai, October 2009, pp. 1589–1595 (2009)
14. Sun, P., Zhu, W., Zheng, J.: Research of Risk Evaluation on Urban Water Supply Underground Pipeline. *Journal of Beijing Institute of Technology (English Edition)* 19(suppl. 1), 31–36 (2010)

Research on Purified Internet Environment for College Students

Qichun Zhong¹ and Jinghong Hu²

¹Department of the humanities, Shandong Jiaotong University, Jinan, 250023, China

²Shandong University of TCM, Jinan, 250355, China
{qczhong2006,hujinghong97}@126.com

Abstract. Together with the information and knowledge age, here comes the widely-extend internet and swarming information. More and more college students have good command of internet application and overall utilization of internet resources. On the same time the internet society has produced negative impacts on the college students. This article analyses the advantages and disadvantages of internet affecting college students in China, especially the moral impact and puts forward some feasible measures to regulate the internet behavior of students and provide a healthy and safe internet environment for students to fully utilize favorable opportunity and convenient condition brought about by network.

Keywords: Internet Environment, College Students, Internet Behavior, Regulation, Moral Education, Virtuality.

1 Introduction

With the development of computer technology, the internet begins to spread to every corner of the world, and opens up a new field of knowledge and practice. The students are the most active internet users in the group. On the one hand, they use the internet to obtain information, acquire knowledge; the other hand, as technology development and the spread of personal computers, the negative effects of internet become increasingly prominent.

At present, amongst the users in China, 11% are over 10 years old and half of them have a higher degree. Thus, the population using internet in China is young students, especially college students. The network quickly won the students generally welcomed as a new thing; the reason is that its inherent features fit the characteristics of the students themselves.

Firstly, students are the generation requiring the high-speed and timeliness information and internet fits this perfectly. Secondly, the richness and openness of internet content meet college students on the pursuit of knowledge and information. Thirdly, the freedom of the network meet the students a strong pursuit of personality psychology. Fourth, the hidden nature of the network interaction meets the students' psychological characteristics of the desire to the truth and doubt of it at the same time. Network provides a perfect golden-distance of a right: to communicate directly, and

retain their privacy at the same time. Fifth, the timeliness of the network meets the requirement of the students' pursuit of the fashion.

However, the network is like a double-edged sword. On one hand, the rich and timely recourses have provided a great convenience for the students' learning and life, on the other hand, the vitality, invisibility and non-binding characteristics has encouraged the students' fluke and indulge, resulting in many students' network moral anomie, triggering a series of social problems, including cyber crime. Many students are addicted to the internet. The dramatic contrast between virtual and reality worlds, plus the lots of pressure of the fierce competition in real world, many students are infatuated with the internet to escape from reality and even suffering from internet addiction syndrome. The survey showed that they most couldn't be able to fit in with their surroundings. They often don't have strong interests in the study, frustration of social relationship, block-minded and distress, etc.

2 Students' Network Moral Enemies

2.1 Extravagance Is Prevalent in Network

The virtual state has protected the behavior of the security barrier, but also immoral to put on a virtual coat, resulting in the spread of false information of the network society and the occurrence of non-ethical behavior. The network moral system is still in the construction and specification, the network society itself is difficult to let student netizens "excepted and independent", coupled in a growing lack of students' self-discipline awareness. It's easily produce free without limit whatever they want and then to make some of the real world rarely do immoral things, such as malicious insult, personal attacks, online "polygonal" love, making up and spreading computer viruses, peeping others' emails, browsing the yellow information, infringement of intellectual property rights, contrary to the Code of Ethics, or even illegal behavior.

2.2 The Fuzzy Concept of Network Value

In a survey by a university of "what kinds of basic moral traits do you think that it should have in the internet?", some students didn't choose the "honest and trustworthy", for they believed that they can fool each other, but only limit to protect their privacy and should not hurt others; some even believed that "no one can see each other in online chat and it's normal that cheating each other". To the "uncivilized language", some believed that "those can be used online but couldn't in reality", while some believe that "it could be used in both online and in reality". Some students neither think that online copy and copying articles is immoral, nor is it immoral for unauthorized use of someone else's online account.

2.3 Personality Conflicts

The virtuality of the network is likely to cause the students self-lost in reality. It's the interaction platform, and also an interaction barrier. It provides the open, free and no

constraints space for them to show themselves, at the same time, it also conceal the true identity. Therefore, it has three kinds of self: true self, real self and network self. The meaning of three “self” sometimes are intertwined sometime conflicted. On one hand, the network has provided a platform to show them, on the other hand, it makes them become more extroverted or introverted offline. The dislocation of the characters, may result in the generation of multiple personality, and thus cause a serious self-loss and impact the development.

The over-lance on the network has caused some students low personality. Some students have used the advantage of the richness and timeliness of network resources. They abandon the assiduously style of study and like to use the short cuts to gain the information that they need. They become falsifier, eager for quick success and instant benefit and blundering. It also causes the negative impact on their studies with less time spending on reading books and thinking over the questions that have been occupied by the network.

3 Basic Measures to Regulate the Students' Network Behavior

3.1 To Adjust the School Moral Education Strategy

To Update the Notion of Moral Education and Infiltrate the Network Moral Education to All Aspects of School Education

Today, moral education in schools and network are at the same time and space. It always faces the impact and challenges of the network media. The school walls can no longer and nor any more necessary penetration of the barrier network. When the students sit before the computers, click the mouse, fly in the network spaces, it's their personality to determine their behavior. Due to their limitation and immature of cognitive level and analyzing ability, it's often difficult for them to make the right decision, especially when it has the multi-dimensional values. To establish a concept of “learn how to choose” is the primary measures for the education adapts the network society. The purpose of moral education no longer require students to accept a few code of ethics, but to help establish a correct outlook on life and values, especially in this complex information environment, and constantly improve their moral judgment and the ability to choose.

To Fully Use Network Resources and Expand the Ways of Education of Moral Education in Schools

The openness, freedom of interaction and democratic, these advantages could help the education of moral education in schools. We should encourage students to directly reflect and express frankly via network. This could help educators understand the real thinking of the students and grasp the problem, thus to develop targeted educational measures and perform the education close to the students, close to reality, and ultimately reach into the brain, heart, fair and reasonable.

To Emphasize and Strengthen Campus Network Construction, Actively Occupied In Online Ideological Position

Internet culture is an open and pluralistic culture. The cultural diffusion and collision blend together. These value orientations increase the conflicts and make it more difficult to make right value choice. In this case, the school must construct the internet culture into the overall planning of the campus culture. It should keep updating the party's line, principles, policies and civilized and healthy cultural information online; and guide students to absorb nutrients in the network, strive to enhance the conscious ability to resist unhealthy information. To strengthen the campus network, it could open the virtual area of the school moral education online.

3.2 To Improve the Quality of the Main Body of Network

To Improve Their Own Abilities and Qualities to Positively Create the Moral Education of New Areas of the University Network for Moral Educators

At present, the students' ideas and value orientation increasingly diversified. The closure of the traditionally ideological education is a great challenge and impact of the openness of network society. The speed and methods for students to receive information have even gone ahead of the educators. To face these new topics of network moral education, moral educators must change their minds and face reality. It's necessary to accept online education, openness and equality, and understand the diversity and complicity of network, but also continue to strengthen the network technology learning and research and innovative ideas and methods.

To Improve the Ideological and Moral Quality of the Educated and Cultured Students A High Degree of Moral Responsibility

As the subject of network moral education, students are a widest range of groups of network users. They are the owner of the high-tech network technology. If they violate the Code of Ethics, the harmful consequences are often greater. On one hand, we should continuously enrich the content of moral education; on the other hand, we should also focus on the moral responsibility to educate students and improve their level of moral character.

3.3 To Strengthen the Network Communication Management and Optimize Network Environment

To Strengthen the Network Legislation

Moral principles and public opinion could constraint and standardize the people's behavior, but can not punish bad behavior beyond the moral boundaries. For the behavior of people in the virtual space of the network, it should be regulated by law. But so far, in terms of punishment for the crimes of the network, it's mainly scattered in the laws and regulations of the Penal Code, Civil Law. The law lags behind the process of the network. It's difficult to adapt to the norms of online information. We

should actively learn from the practices of the United States, Germany, Japan, and create relevant laws and regulations to prevent young people's bad behavior in network. We should amend, add and improve the relevant legal content, in order to protect the interests of young people network and use of management.

To Strengthen the Network Technology Research and Control of Harmful Information Dissemination

Relative to the network in terms of the legislation and moral self-discipline, technical control is the most objective means. We should take advantage of modern high-tech means to purify the internet environment to control the spread of bad information from the source, to fill gaps in network vulnerability, to enhance firewall functions and closely monitor the internet portal, and eliminate restrictions on pornography, crime and other unhealthy content and information, and strive to the creation of the positive, healthy and orderly network environment.

To Strengthen the Management of Internet Bars In Accordance with the Law

It should stringent requirements of telecom operators to provide internet access services in accordance with state regulations for operating without a license, overtime, and daily limited time for all internet cafes to stop network access services. Install network security monitoring software in all cafes to effectively prevent young people in the cafe, or even harmful use of the network on a home computer. In addition, to guide the cafes policy and strengthen the scientific and rational way to establish the scale of development the number of around cafes, to prevent haphazard development and adverse competitive, and to guide the development of internet cafes to the direction of information services.

3.4 To Found a Group of Outstanding Chinese Sites and Make They Become the Main Body of the Network Civilization as Soon as Possible

The problems of the students in the process of using the internet are largely focused on the lack of high-quality professional web site. Education, culture, industry and commerce, public security departments should support a group of healthy, civilized, scientific, especially the Chinese website that could provide young people with a correct role in guiding and attractive information. It's the most active and effective way to enable students to learn effective use of network information, network communication, and use the internet to conduct scientific research, develop and improve their abilities.

References

1. Nancy, O.: Net knowledge: Performance of new college students on an Internet skills proficiency test. *The Internet and Higher Education* 5, 55–66 (2002)
2. Rolf, H.W.: Internet of Things-New security and privacy challenges. *Computer Law & Security Review* 26, 23–30 (2010)

3. Steven, F., Valleria, T., Andy, P.: Security beliefs and barriers for novice Internet users. *Computers & Security* 27, 235–240 (2008)
4. Eric, Y.L., Hongyan, M., Sandra, T., Wayne, H.: Wireless Internet and student-centered learning: A Partial Least-Squares model. *Computers & Education* 49, 530–544 (2007)
5. Steve, J., Camille, J.Y., Sarah, M., Francisco, S.P.: Academic work, the Internet and U.S. college students. *The Internet and Higher Education* 11, 165–177 (2008)
6. Byung, C.K., Yong, W.P.: Security versus convenience? An experimental study of user misperceptions of wireless internet service quality. *Decision Support Systems* 53, 1–11 (2012)
7. Melike, K., Hafize, K., Necmettin, T.: Reviewing unethical behaviors of primary education students' internet usage. *Procedia-Social and Behavioral Sciences* 28, 1043–1052 (2011)
8. Nicola, D., Elizabeth, S.: It won't happen to me: Promoting secure behaviour among internet users. *Computers in Human Behavior* 26, 1739–1747 (2010)
9. Aashish, S.: Is internet security a major issue with respect to the slow acceptance rate of digital signatures? *Computer Law & Security Review* 21, 392–404 (2005)
10. Chien, C., Hsinyi, P.: Promoting awareness of Internet safety in Taiwan in-service teacher education: A ten-year experience. *The Internet and Higher Education* 14, 44–53 (2011)

A Mobile-Certificate Security Method of Satellite-Earth Integration Networks

Qianmu Li^{1,*}, Qiugan Shi¹, Jun Hou³, Yong Qi², and Hong Zhang¹

¹ School of Computer Science, Nanjing University of Science and Technology, Nanjing 210094, China

² Wuxi Institute, Nanjing University of Science and Technology, Wuxi 214192, China

³ Zi Jin College, Nanjing University of Science and Technology, Nanjing 210094, China

liqianmu@126.com,

{848853385, 727545277, 790815561, 475019943}@qq.com

Abstract. With high-speed move and frequent cross-domain, the problem of certificate authority (CA) is a security bottleneck in the Satellite-Earth (S&E) Integration Network. A mobile-certificate authority Method based on PKI (Public Key Infrastructure) is proposed in this paper. Based on the RSA algorithm, the method realizes the anonymous authentication between neighboring nodes. It distributes the private key component to n nodes according to the threshold theory, and then realizes multi-node collaboration certification and dynamic permission-recovery. In addition, combining with distributed authentication and chain certification, this paper achieves distributed PKI and provides certification basis for high-speed nodes. This method frame, which is no-central, self-adaptive and traceable, could fulfill the quality of the S&E Integration Network. The experiments showed that the success rate and safety of it is much better than traditional distributed certification method. The method has far-reaching military value

Keywords: Network Security, Distributed Certificate Authority, Satellite-Earth Integration Networks.

1 Introduction

The S&E Integration Network is a typical heterogeneous interconnection information network system which is composed of the ground network, star network and interplanetary Internet. It is easy to be stolen and interference because of the large coverage area of network and the strong mobility, which seriously hindered its popularization and application. The PKI system is currently believed to be the solution to the security problems in a large open network environment. N.Asokan and P.Ginzboorg[1] have proposed a key exchange mechanism system, which is used to build trusted relationship between peers in the environment where there is not a trusted third party. The solution is actually an extension of Diffie-Hellman key exchange method [2]. It's suitable for small network which can rapidly establish

* Corresponding author.

group communication, but cannot for point to point secure communications, and it cannot provide non-repudiation services either. D.Balfanz[3]and the others proposed the verifiable identification mechanism. It's used in the local small-scale PZP secure communication between members. This solution cannot provide anti-repudiation services, and the relationship between the entities and their public keys cannot be certificated. And the solution cannot adapt to the changes of the scale of network. L. Zhou and Z. Hass [4] proposed the key management service system conceptual model and raise the idea of not trust the distributing principle. But in this manner, the single dealer has the whole private key information of CA. That information being exposed will endanger the entire network.

Therefore, the traditional authentication methods not only have complex node access logic, but also are not suitable for the S&E Integration Network. This paper designs a set of Mobile-Certificate Authority Method. Based on the RSA algorithm, the method realizes the anonymous authentication between neighboring nodes. It distributes the private key component to n nodes according to the threshold theory. The function of CA requires multiple nodes. The damage of one node will do not affect the whole system to continue to run. The method realizes multi-node collaboration certification and dynamic permission-recovery. In addition, combining with distributed authentication and chain certification, the method achieves distributed PKI and provides certification basis for high-speed nodes. This method frame, which is no-central, self-adaptive and traceable, can fulfill the qualities of the S&E Integration Network as no-central and node equivalence.

2 Anonymous Mobile-Certificate Authority Method

The basic idea of this new method: Firstly, a dealer generates the RSA parameters(N , e , d) in authentication system randomly, and then broadcasts the parameters (N , e) (public-key in RSA) in the communication field, then selects n trusted nodes and gives the private-key in RSA to them. These n nodes become the distributed CA nodes of the authentication system. Before the communication starts, do authentication using the fake name. It will not do the later works such as exchanging certificates until the fake name verify passed. If a new node needs to join, the new node should find a proximate distributed CA node, and then apply a physical certification to the CA node. If the new node passes the physical certification, it will get a passing authentication certificate (hereinafter referred to as a certificate) issued by the CA node who undertakes the physical certification. After receiving a certificate, the new node sends requests for certificate application to other neighboring nodes. If it can receive at least t distributed CA's signatures, then it can synthesize a legitimate certificate.

The connection process of star network and terrestrial network is divided into three stages: detection, authentication and association in the S&E Integration Network as shown in Fig. 1.

In the detection stage, star network can receive ground radio passively, and add to ground station network automatically. In the second stage, the star network receives the request of certificate message sent by ground station according to the authentication mechanism. After successful authentication, star network sends

association request to ground station. The ground station records the star network to association table after receiving the association request .Usually the ground station can establish the association with multiple star networks simultaneously, while star network can only establish the association with one ground station at the same time.

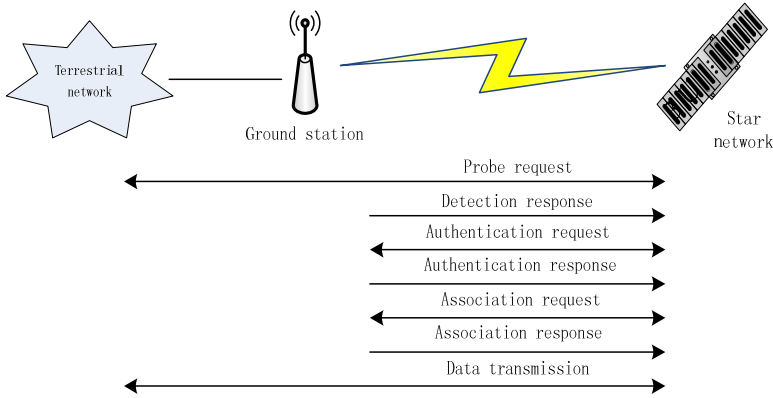


Fig. 1. The process of the network establishes the association with ground station

The initialization of this method are as follows:

Step1: The dealer generates the RSA parameters: N, e, d randomly and the broadcasts parameters: N, e (public-key, PK) to the field.

Step2: The nodes in the communication field send networking requests to the off-line manager. The off-line manager does physical authentication for those nodes that have sent requests. After the physical authentication, the off-line manager select n trusted nodes, and divides the key d into n components, and then stores the components into the n nodes. These nodes will be the distributed CA nodes in the communication field. Detailed description:

Randomly select a group of numbers: $a_1, a_2 \dots a_{t-1}$, generate a $t-1$ polynomial function $f(x)$ with the private-key: $f(x) = d + a_1x + a_2x^2 + \dots + a_{t-1}x^{t-1} \text{ mod } N$

Mark the nodes as V_i , and substitute the variables in the above equation, $y_i = f(v_i) \text{ mod } N (i = 1, 2, \dots, n)$.

By Lagrange formula: $f(x) = \sum_{i=1}^t y_i \prod_{j \neq i} (x - v_j) / (v_i - v_j) \text{ mod } N$

When $x = 0$: $f(0) = d = \sum_{i=1}^t y_i \prod_{j \neq i} v_j / (v_i - v_j) \text{ mod } N$

And then get the additional shared key : $d_i = y_i \prod_{j \neq i, j=1}^t v_j / (v_i - v_j) \text{ mod } N$

Save d_i into n nodes respectively. Those nodes which have the partial key d_i become certified nodes. Any t signatures of certificated nodes can cooperate in generating a complete and valid certificate.

Step 3: Those n trusted nodes generate their own PK_i / SK_i , and send the public-key PK_i to the manager to request for certificates. The manager will assign certificates which contain the public-key PK_i and the expiration information to the distributed CA nodes using the private-key SK after receiving the requests for certificates. Meanwhile it selects several bits of the public-key PK_i to calculate the fake name for node i , and records the relationship between the fake names and the real names for later use.

Step 4: The manager calculates the value of CBF (Compressed Bloom Filter) [6] according to the fake name. Each node assigns the trust certificates to the trust nodes around according to its public-key and private-key (PK_i / SK_i) pair. If node u thinks that public-key PK_v belongs to node v , u will bind PK_v and node v together, signs it with its private-key and gets a certificate like $(v, PK_v, T_{issue}, T_{expire})_{SK_u}$ and sends this certificate to node v . Each private certificate has a value called TTL whose initial value is $h-1$. The value decreases by 1 after each exchange. Node u exchanges certificates with its trust neighbours, and saves those certificates that it doesn't have. If there's not enough space, the node deletes the oldest certificates according to their expiration information.

Step 5: The manager broadcasts the final value of Compressed Bloom Filter to the communication field. After finishing the work, the off-line manager quits the communication field. The manager will save all the information if it can ensure security, otherwise it will destroy the entire secret information.

Anonymous mutual authentication between the inter-nodes of this method are as follows: Before the communication starts, the two nodes A and B need to do the anonymous mutual authentication which means they have to confirm the other one's identity.

Assuming that A starts the communication and B is being called. Here are the steps:

Step 1: A sends a request for the anonymous mutual authentication to B. If there's no response in a certain time, A resends the request or ends the request. B first validates the fake name of A by the Compressed Bloom Filter method after receiving the request from A. If the validation passed, B sends a confirm message which contains the fake name of B to A, and then begins to wait for the certificate of A, else B will discard the message.

Step 2: A first validates the fake name of B by the Compressed Bloom Filter method after receiving the confirm message from B. If the validation passed, A sends its certificate which contains the public-key of A and the expiration information to B, and then begins wait (wait for receiving B's certificate), else A ends the communication.

Step 3: B validates the certificate after receiving it from A. If the validation passed, B sends its own certificate to A and keep A's certificate as well which means B sets up a trust relationship with A, else B discards the certificate.

Step 4: A validates the certificate after receiving it from B. If the validation passed, A keeps B's certificate which means A sets up a trust relationship with B, else A discards the certificate.

They will only verify the fake name in order to improve the communication efficiency after the trust relationship be set up.

The steps of assigning the certificates to a new node are as follows :(assume that N is a new node and M is the distributed CA node)

Step 1: The distributed CA node M sends its certificates to new node N. Node N will first verify the certificate with the public key PK which is in the communication domain after receiving it from M. If the verification passed, node N randomly generates a symmetric session key K and sends K to M after encrypted K with M's public key PK_M which is in M's certification.

Step 2: After receiving a random symmetric session key k which is encrypted with M's public key PK_M , M decrypts it with M's private key SK_M and then gets the random symmetric session key K. At the same time, M signs a digital signature with its own private key SK_M on the certificate which is going to be issued to N, and puts a time-stamp (i.e., time tag) on it. Finally, M encrypts the signed certificate with K and then sends the certificate to N.

Step 3: N decrypts it with K after receiving the encrypted certificate, and then gets a passing authentication certificate.

Issue system certificate to the new node is as follows:

Step 1: the new node N broadcasts packets to its neighboring nodes to send joining request. The distributed CA node M which has received request packets sends its own certificate to N.

Step 2: N verifies the validity of M's certificate. At the same time, N generates a random symmetric session key K.

Step 3: After the validation of M's certificate being passed, N first encrypts the certificate with K, and then encrypts K with M's public key PK_M which is in M's certificate. Finally, N also encrypts its certificate which is applying for signature with K and sends things all above to M.

Step 4: after receiving the packets, M decrypts them with M's private key SK_M and gets the symmetric key K. M decrypts packets with K and gets its own certificate and N's certificate which is to be signed.

Step 5: M encrypts a random number sequence with N's public key PK_N which is in the certificate that is going to be signed. And then sends the number sequence to N to make sure that N actually holds the private key SK_N matching with the public key PK_N .

Step 6: after receiving the random number sequence which has be encrypted, N decrypts it with N's own private key SK_N and then encrypts the number sequence with the symmetric key K and last sends the number sequence back to M.

Step 7: after having confirmed that the sequence received from the N is the same as the previous one, M signs N's certificate with private keys of all distributed CA held by M and attaches its own digital signature and time-stamp. Finally M encrypts them together with the symmetric key K and then sends them to N.

Step 8: N decrypts it and gets a signed certificate and then verifies if the certificate has been modified during transmission with digital signature. Finally, N gets a partly signed certificate.

Repeat the above steps until N gets at least t copies of partly signature, and then we can synthesize a legitimate certificate.

When new node broadcasts message and then finds that it cannot link directly to achieve the threshold number of authentication nodes, it will look for authentication nodes through passing the trust. Node N sends proxy searching request to neighboring distributed CA nodes, and lets these neighboring nodes complete the searching job for it. These neighboring nodes will accept the request if they trust the new node (this trust relationship is also guaranteed by the certificate).

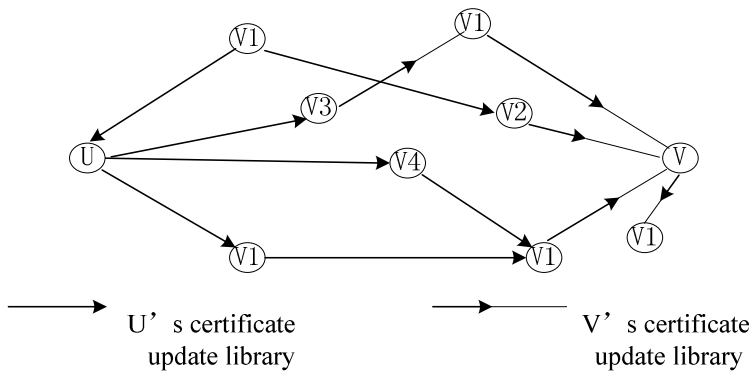


Fig. 2. Supposing that node u is the new node N's neighbor (u must be one of distributed CA nodes) and accepts the proxy request

As shown in Fig. 2, node u finds a certificate chain leading to authentication node from the intersection between updating certificate library S_u of u and the intended authentication node's updating certificate library S_v . The certificate chain must satisfy the following conditions. the first certificate can be verified directly by node u. The public key binding to the last certificate must belong to node v. The legitimacy of each certificate can be verified with the public key included in the previous certificate in the certificate chain. As shown in Fig. 2, the path (U, V_3, V_8, V) , (U, V_4, V_6, V) and (U, V_5, V_6, V) is the certificate chain to find out. If node u cannot find an appropriate path in $S_u \cup S_v$, node u searches them in the non-updating certificate library of the two nodes, and verifies if these certificates are legitimate. If the searching for trust chain is failed, node u needs to send message to inform the new node N. When the new node N receives a partly signed certificate; it saves the certificate and waits for the next

partly signed certificate. After N having received t copies of legitimate partly signed certificate, N can synthesize an effective system certificate immediately and join the network. At the same time, N sends message to inform neighboring proxy nodes about stopping searching authentication nodes for it.

3 Analysis and Verification of the System

In this paper, the simulation scene will be shown on the simulation platform ONES. The default parameter values of simulation are: Mobile model is RWP, Minimum speed=0.5, Maximum speed=1.5, Pause time=0s, Simulation region $E=5000 \times 5000 \text{KM}^2$, Message size=512-1024KB, Bandwidth=250kBps, Traffic load=5000-6000, Node number $M=200$. Target delay $Dt = 11\text{k}/\text{unit}$ simulation time, $TTL=20\text{k}/\text{unit}$ simulation time, Node transmission range $R=200\text{m}$, the simulation duration= $2 \times Dt$.

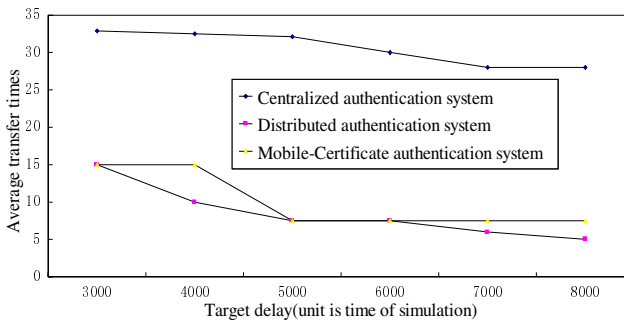


Fig. 3. The comparison of performance under different delay constraint is shown

This fig. 3 tests the methods' average end-to-end delay and average transmission times in different delay constraint Dt . Simulation time is $2 \times Dt$. The simulation results show that although the centralized authentication has the smallest end-to-end transmission delay, the average transmission times is substantially higher than the distributed method's because of the multi-hop routing. Traditional distributed method has longer delay due to combine certificate by serial mode. It is unable to meet the strict delay constraint. The extreme speed of traditional distributed method is nearly the same with the method in this paper, because the method's attack resistance is greatly reduced although distributed method interactive times are similar to this method system when t quantity in distributed system is small enough.

In the S&E Integration Network environment, this method not only could tolerate a number of verified nodes being broken, but also prevent the interference to verification from the malicious nodes.

Firstly, the difficulties of calculating of discrete logarithm problem in limited field ensure the security of the system. The signature function which uses RSA algorithm is hidden with index transforming. What transports in the network is the deformation of the key. In the most extreme case, even the attacker has collected part of the returning

signatures from every verified nodes, it is hard to calculate the SK. The difficulty is just like choosing plaintext attack to RSA algorithm.

Secondly, there are malicious nodes distributing the key component arbitrarily to pretend to be verified nodes to destroy the reconstruction of the key. In the mutual anonymous among domain nodes scheme that is designed in this paper, if a new node v_j find out the signature of some verified node sent by its neighboring node cannot go through the verification, it will give it up directly without judging after synthesizing and broadcast information with its signature to point out which node has the cheating actions.

Lastly, this paper calculate the CBF(Compressed Bloom Filter) value with fake name .The consumptions of the source are just some calculations of Hash function, low consumption and high efficiency. It can prevent moving attacker from attacking the node with DOS attack efficiently.

This system defines successful percentage as follows: It is the average percentage of receiving parts of threshold signature certificate and adding to the network successfully when the normal node hands out a request from the un-added network after network’s initialization.

Assuming verified nodes are moving randomly in the network, and the number of them is n, it needs t verified nodes to compound a legal certificate. The area net covers is S, the communication radius of normal node is r.

Thus, if we take distributed certification plan, the successful rate would be

$$R_{DCA} = \sum_{i=t}^n \left\{ C_n^i \times \left[\frac{r^2}{S} \right]^i \times \left[1 - \frac{r^2}{S} \right]^{n-i} \right\}.$$

Since the initialization of anonymous certification system , the trust certificates between the nodes are put in a settable TTL, we can know that the trust link of every node could be h-1 jump at most. Joining the jump from the new one to neighbor, the

mixed certification’s successful rate could be $R = \sum_{i=t}^n \left\{ C_n^i \times \left[\frac{(hr)^2}{S} \right]^i \times \left[1 - \frac{(hr)^2}{S} \right]^{n-i} \right\}.$

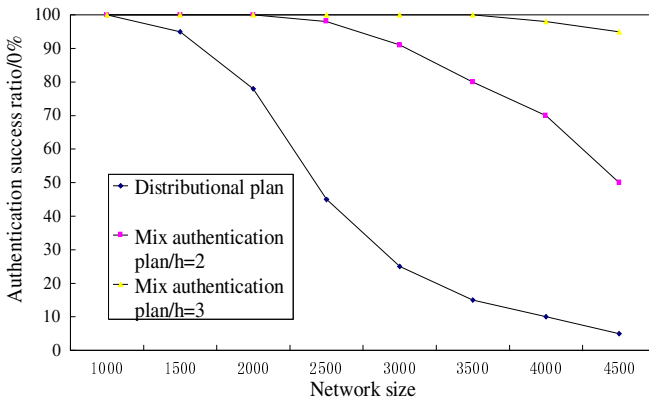


Fig. 4. Success rate of the distributed and hybrid authentication is shown

Because of the restriction of the length of trust link, the rate of successful searching of certification node is related to h .

In the case of the radius of nodes communication r is 200, network covered area is 5000m*5000m, threshold is (3,200), the successful rate of distributed and mixed are in the fig. 4 below:

According to this fig. 4, in terms of the successful rate, this authentication method is obviously better than distributed one even though the number n of certification nodes in this Mobile-Certificate Authority Method system is less than distributed ones. It can ensure a higher successful rate of certification. With the increasing of trust link h and nodes' moving speed, the advantage is getting greater.

4 Conclusions

In this paper, combining with distributed authentication scheme and chained authentication scheme, the Mobile-Certificate System which fulfill the quality of the S&E Integration Network is proposed. This method is an infrastructure based on PKI which is a fusion of threshold theory and RSA algorithm. And it has the characteristics of no-central, self-adaptive and traceable and applies to the harsh environment in which Information is easy to leak and Individual nodes cannot ensure safety. It can guarantee the communication security and robustness and has wide application prospect no matter in military or in commercial.

Acknowledgments. The research for this paper has been funded by the Jiangsu 973 Scientific Project (BK2011023, BK2011022), the National Natural Science Foundation of China (60903027), the Jiangsu Natural Science Foundation (BK2011370), the Aerospace Innovation Fund (CALT201102) and the Lianyungang Science & Technology Project (CG1124).

References

1. Asokan, N., Ginzboorg, P.: Key agreement method in Ad Hoc networks. *Computer Communications* 23(7), 413–422 (2006)
2. Diffie, W., Hellman, M.: New directions in cryptography. *IEEE Transactions on Information Theory* 22(6), 305–316 (2004)
3. Balfanz, D.: Authentication in Ad hoc Wireless Networks. In: *Proceedings of NDSS 2002 Conference*, pp. 1652–1671. American Institute of Aeronautics and Astronautics, New Orleans (2002)
4. Zhou, L., Hass, Z.: Securing Ad Hoc Networks. *IEEE Network Magazine* 13(6), 751–765 (2003)
5. Guo, S., Shen, A.-N.: A compromise-resilient pair-wise rekeying method in hierarchical wireless sensor networks. *Computer Systems Science and Engineering* 25(6), 397–405 (2011)
6. Li, Q.-M., Liu, F.-Y.: Strategic Internet risk detection and fault analysis method. *Journal of Computer Research and Development* 45(10), 1718–1723 (2008)

7. Huang, J., Zuo, M.J., Fang, Z.: Multi-state Consecutive k-out-of-n systems. *IIE Transactions on Quality and Reliability Engineering* 16(6), 527–534 (2003)
8. Huang, C., Wang, G.-L.: Energy-efficient beaconless real-time routing method for wireless sensor networks. *Computer Systems Science and Engineering* 26(3), 173–186 (2011)
9. Li, Q.: Multiple QoS Constraints Finding Paths Algorithm in TMN. *Information: An International Interdisciplinary Journal* 14(3), 731–738 (2011)
10. Qi, Y., Li, Q., Hou, J.: A Method to Solving Cyberspace Security-model WSN Security-model Equation. *Procedia Engineering* 15, 2052–2056 (2011)
11. Li, Q., Wang, R., Yin, J., Hou, J.: The design of data security synchronization in the network of satellite and ground security. *Key Engineering Materials* 439-440, 208–213 (2010)

Duration Modeling for Emotional Speech

Wen-Hsing Lai and Siou-Lin Wang

Dept. of Computer and Communication Engineering, National Kaohsiung First University
of Science and Technology, No.2, Jhuoyue Rd., 81164 Kaohsiung, Taiwan
{lwh,u0015901}@ncku.edu.tw

Abstract. Human interaction involves exchanging not only explicit content, but also implicit information about the affective state of the interlocutor. In recent years, researchers attempt to endow the computers or robots with humanity. Various affective computing models have been proposed, which covers the areas of emotion recognition, interpretation, management and generation. Therefore, to analyze and predict the prosodic information of different emotions is very important for the future applications. In this article, a duration modeling approach for emotional speech is presented. Seven kinds of emotion including natural, scare, angry, elation, sadness, surprise, and disgust are adopted. According to the statistics performed on a corpus with seven emotions, a question set considering acoustic and linguistic factors is designed. Experimental results show that the root mean squared errors (RMSEs) of syllable are 0.0725s and 0.0802 s for training and testing sets correspondingly. From the results, the impact of factors related to different emotions can be explored.

Keywords: Emotion, Duration, Binary Decision Tree.

1 Introduction

Human interaction involves exchanging not only explicit content, but also implicit information about the affective state of the interlocutor. In recent years, researchers attempt to endow the computers or robots with humanity. Various affective computing models have been proposed, which covers the areas of emotion recognition [1]-[7], interpretation, management and generation [8][9]. Therefore, to analyze and predict the prosodic information of different emotions is very important for the future applications.

In the studies of emotions, some researchers used rule-based method by collecting rules from literatures [10] or corpus [11]. Data-driven method [12] is also proposed. In this paper we proposed a duration modeling method by binary decision tree. The corpus we used is recorded in seven emotions and the text material is balanced sentences. In order to be able to design a prosodic model to analyze the features of sentences, a question set is designed according to the linguistic knowledge. After the completion of the training, the estimated duration can be generated according to the established tree.

The sections of this paper are summarized as follows. Section 2 discusses the seven emotions selected; Section 3 introduces the corpus; Section 4 describes the method of binary decision tree; Section 5 presents the experimental results; Finally, Section 6 summarizes this paper and describes the future works.

2 Emotions

Emotion is an important way of communication among people. The definition of emotion can be interpreted from the point of view of biology, sociology, psychology, theology, and so on [13, 14]. In our corpus, seven representative emotions are adopted, namely:

1. Neutral: When people talk with no emotion, we call it neutral.
2. Scared: It is the emotion expressed when confronted with dreadful things. A similar definition is fear.
3. Angry: It can also be called furious. Generally it is classified as a negative emotion because it is a reaction of offensive, hostile, or violent behavior. However, in the point of view of some psychologist [14], anger confers the individuals the vitality of defense. It can, therefore, also be interpreted as a positive emotion.
4. Elation: The similar are happiness and joy. Elation, happiness, and joy cannot be clearly distinguished. The features of elation include self-confidence and the feeling that brings human sensory satisfaction.
5. Sadness: It is also known as sorrow. A similar emotion is grief. The sadness is generally defined as the sentiments expressed when confronted with things that make one in a low mood. Although it is a negative emotion, it can also be used as a positive one.
6. Surprise: The surprise can be either positive or negative.
7. Disgust: It is a kind of negative emotion and is mainly the physical or spiritual rejection generated when confronted with certain things or personal bias.

3 Corpus

Phonemic balanced sentences are chosen as recording material. There are totally 168 sentences including 1616 syllables. The sentence length is between 7 and 10 syllables and the average length is 9.6 syllables.

The speaker is a 12-year-old elementary school student whose mother language is Mandarin Chinese. The student was asked to playact the seven kinds of emotions in voices. With each sentence interpreted with seven emotions, the total number of sentences is 1176. The recording data format is 96-KHz and 24-bit. Syllable boundary annotations are labeled manually and syllable duration can be determined easily from the labels.

The corpus was divided into two sets, the Training set and the Testing set. Including 7 emotions, the former contains 770 sentences including 6111 syllables and the latter contains 406 sentences including 5201 syllables.

4 Binary Decision Tree

A binary decision tree is adopted for prosodic modeling. The tree relies on a Yes-and-No question set to determine the path of branch. In this section, first, the question set is discussed and then the binary decision tree is introduced.

4.1 Question Set

A Yes-and-No question set for the binary decision tree is designed based on the characteristics of Mandarin Chinese. The questions can be grouped as five categories as follows.

1. Syllable position in Sentence: It is classified as {the first syllable, the last syllable, the first 1/3 position with the first syllable excluded, the second 1/3 position, the last 1/3 position with the last syllable excluded}.
2. Initial/Final Class: Chinese syllable can be divided into Initial and Final, and Finals can be further divided into medial, nucleus and coda. Final is classified as {monophthong, diphthong, Final with nasal ending}, and Initial is classified as {voiceless, voiced}. The voiceless Initial is further classified as {aspirated, un-aspirated}, and the voiced Initial is further classified as {nasal, un-nasal}.
3. Tone: Chinese syllable contains five tones, namely Tone 1, Tone 2, Tone 3, Tone 4, and Tone 5 (also known as Light Tone). For the Tone 5, the duration is generally the shortest.
4. Word length and syllable position in word: The syllable position in word is classified as {single-syllable word, first syllable, last syllable, others}. In addition, the word length is classified as {one-syllable word, two-syllable word, three-syllable word, four-syllable word, five-syllable word, others}.
5. Part of Speech: It is classified as gerund, determiner, quantifiers, numerals, auxiliary verb, expletive, preposition, pronoun, adverb, conjunctions, adjectives, verbs, and nouns.

4.2 Binary Decision Tree

A top-down binary decision tree used the following criterion to determine whether a node (cluster) was to be split into two son nodes (sub-clusters) based on a specific question:

Split based on the question with maximum $|\mu_1 - \mu_2|$. If $(n < \text{Threshold})$, then stop. Here (n, μ, v) , (n_1, μ_1, v_1) and (n_2, μ_2, v_2) are, respectively, sample counts, means and variances of the node and the two son nodes split based on a question.

By the above process, a generated binary decision tree of neutral emotion by using the data in Training set is shown in Fig. 1 as an example. A total of 182 leaf nodes

were generated. The value shown on each node is (sample counts, mean, variance) of the node. In prediction stage, the corresponding mean of the node which the sample falls into is assigned as the estimated duration.

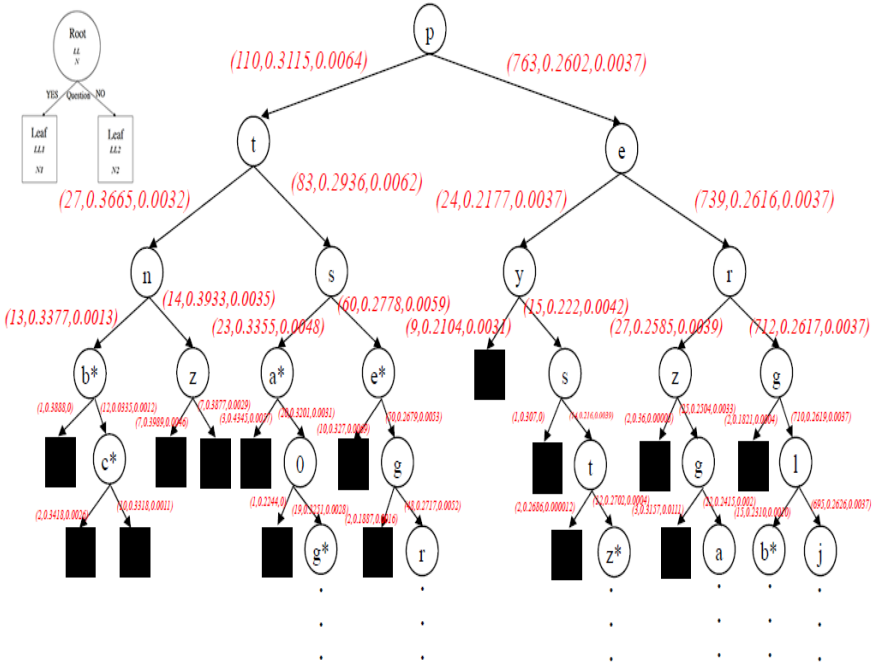


Fig. 1. Generated binary decision tree of neutral emotion

5 Experimental Results

Syllable is used as the basic unit for statistics and prosodic modeling since Mandarin Chinese is a syllable-based language. First, the distributions of duration of seven emotions are observed and compared and then, the results of binary decision tree are shown.

5.1 Statistics

The distributions of duration of the seven emotions are compared in Fig. 2. The 25-75 percent quartiles are drawn using a box. From Fig. 2, the emotions of angry, surprise, and disgust show the shorter duration by observing the center of box. However, some of the samples of the angry emotion are very long, and some are quite short, though the 25-75 percent quartiles concentrate on shorter duration.

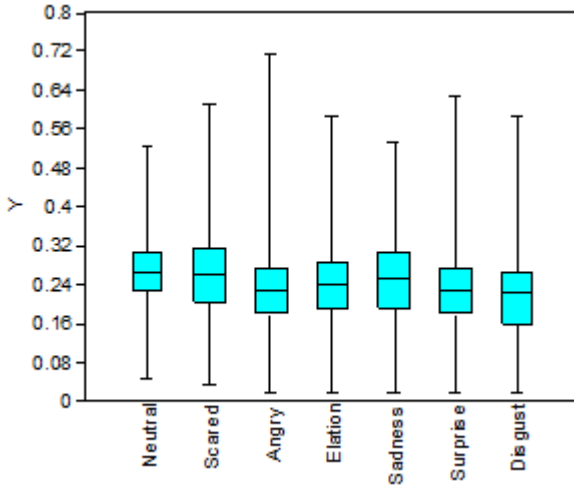


Fig. 2. The box plot of the duration (unit: second) for the seven emotions in Training set

5.2 The Results of Binary Decision Tree

After the establishment of the duration binary tree, the root mean square errors (RMSEs) of the Training and Testing set can be calculated. The RMSEs of training and testing for different emotions are shown in Table 1. The average RMSEs of syllable are 0.0725s and 0.0802 s for training and testing sets correspondingly.

Table 1. The RMSEs of the Training and Testing for the seven kinds of emotions

Emotion	Training	Testing
Neutral	0.045167	0.073784
Scared	0.077812	0.072924
Angry	0.076462	0.083847
Elation	0.068556	0.071624
Sadness	0.076486	0.080897
Surprise	0.093894	0.096007
Disgust	0.069738	0.082648

6 Conclusions and Future Works

This study analyzed the duration for different emotions and established the binary decision tree. Experimental results show that the RMSEs of syllable are 0.0725s and 0.0802s for training and testing sets correspondingly. From the tree, the impact of factors related to different emotions can be explored.

Besides duration, the prosodic features influenced by emotion include pause, pitch, and energy. Modeling and integrating these features will be very helpful for the

analysis of prosodic patterns of emotional speech and can be applied to systems like emotional speech identification and emotional speech synthesis.

Acknowledgments. This work was supported by NSC, Taiwan under Contract NSC100-2410-H-327-037-MY3.

References

1. Wu, C.H., Liang, W.B.: Emotion Recognition of Affective Speech Based on Multiple Classifiers Using Acoustic-Prosodic Information and Semantic Labels. *IEEE Trans. on Affective Computing* 2(1) (2011)
2. Koolagudi, S.G., Kumar, N., Rao, K.S.: Speech Emotion Recognition Using Segmental Level Prosodic Analysis. In: *ICDeCom* (2011)
3. Luengo, I., Navas, E., Hernández, I.: Feature Analysis and Evaluation for Automatic Emotion Identification in Speech. *IEEE Trans. on Multimedia* 12(6) (2010)
4. Lee, C.C., Mower, E., Busso, C., Lee, S., Narayanan, S.: Emotion Recognition Using a Hierarchical Binary Decision Tree Approach. *Speech Communication* 53 (2011)
5. Schuller, B., Batliner, A., Steidl, S., Seppi, D.: Recognising Realistic Emotions and Affect in Speech: State of the Art and Lessons Learnt from the First Challenge. *Speech Comm.* 53, 9–10 (2011)
6. Zeng, H.Z., Tu, J.L., Pianfetti Jr., B., Huang, T.S.: Audio–Visual Affective Expression Recognition Through Multistream Fused HMM. *IEEE Trans. on Multimedia* 10(4) (2008)
7. Slaney, M., McRoberts, G.: BabyEars: A Recognition System for Affective Vocalizations. *Speech Communication* 39 (2003)
8. Iida, A., Campbell, N., Higuchi, F., Yasumura, M.: A Corpus-based Speech Synthesis System with Emotion. *Speech Communication* 40 (2003)
9. Schröder, M.: Expressing Degree of Activation in Synthetic Speech. *IEEE Trans. on Audio, Speech, and Language Processing* 14(4) (2006)
10. Murray, I.R., Amott, J.L.: Synthesizing Emotions in Speech: Is It Time to Get Excited. In: *Fourth International Conference on Spoken Language*, vol. 3 (1996)
11. Al-Dakkak, O., Ghneim, N., Abou Zliekha, M., Al-Moubayed, S.: Prosodic Feature Introduction and Emotion Incorporation in an Arabic TTS. In: *2nd Information and Communication Technologies* (2006)
12. Jiang, D.N., Zhang, W., Shen, L.Q., Cai, L.H.: Prosody Analysis and Modeling for Emotional Speech Synthesis. In: *ICASSP* (2005)
13. Vidya Sagar, T., Sreenivasa Rao, K., Prasanna, S.R.M., Dandapat, S.: Characterization and Incorporation of Emotions in Speech. In: *IEEE INDICON* (2006)
14. Strongman, K.T.: *The Psychology of Emotion - Theories of Emotion in Perspective*. Wu-Nan Book Inc., Taipei (1992)

Research on Image Retrieval Based on Color and Shape Features

Hongwei Zhao, Xiao Chen, Wei Huang, Pingping Liu, and Lingjiao Ma

College of Computer Science and Technology, Jilin University,
Changchun, Jilin Province, China
xiaoxiaocctv5@163.com

Abstract. According to the robustness of color and shape feature extraction, a multi-feature matching algorithm which is combined of color features and shape features is proposed. In the respect of color feature, a new color histogram method based on main colors is proposed. By combining the major color retrieving method and the color histogram computing, two rapid elective filter are carried out, scope of the search is reduced and the retrieval efficiency is improved. In the respect of shape feature, the use of Fourier shape descriptor, an improved contour-based description method is proposed. According to the tangential angle of contours(curvature) is highlighted and factors such as complex coordinates and center distance are ignored, within a reasonable range, the accuracy is lowered appropriately and the query speed is improved significantly. Experiments in traffic signs image library, show that the proposed method of recognition accuracy is better than traditional methods, and efficiency has improved.

Keywords: image retrieval, color features, shape features, color histogram.

1 Introduction

With the rapid development of science and technology, high-capacity storage devices and digital information technology equipment is filled with people's lives, multimedia technology and network technology also enables the rapid development of information and data presented geometric growth trend [1-3]. In order to store vast amounts of data, we have to build a large database to store the data. However, the attendant problem is that it gives us a heavy burden on the search and retrieval. How to solve this difficult problem, recent research has been discussed in the scientific community and the hotspot and difficult. With the rapid development of multimedia technology, although people in image compression, video restoration, image processing, image storage, etc[4]. have made great progress, but has been overlooked aspect of image retrieval[5]. To address this increasingly acute problem, automated, intelligent, humane way of image retrieval, image retrieval, a new technology came into being, and it is Content-Based Image Retrieval,CBIR. This paper proposed based on color and shape of the characteristics of the technology and application of image retrieval precisely this area of research content.

2 Color Feature Extraction

2.1 Color Feature Extraction Method

Color space refers to a different wavelength of electromagnetic spectrum with different materials interaction that forms the chromatographic space. We are to do image feature extraction, often depends on the color in images of color space representation and the understanding. Color space to the color mathematical representation is a very important component, at the same image speaking, different color and color space transformation between the features of quantitative are deciding factor. Now most of the color model adopted is facing the hardware or for application. Color space from pose to now has hundreds of, many improved space model mostly only local change or special in a certain area. This paper will respectively introduced two of the most common color space, by comparing the differences between them, thus inferiority and selecting a suitable this paper color space[6].Color image engineering of quantitative is basic and important technology. it is an image segmentation and object extraction foundation. Along with the computer display system development, the current mainstream format for storing true color. True color images of 224 colors, and is in the actual sampling process, if we study of all the colors, often is not realistic. How to select the representative several kinds of color, and the various colors are incorporated into the delegates color, is color quantitative problems need to be solved.

Color histogram is refers to the use of statistics reflect the histogram form in the composition of the image color, namely all sorts of color in the whole image the proportion of the probability and occurrence. As first suggested using color histogram image expression method of color characteristics of Swain and is Ballard. Color histogram in many image retrieval systems is widely used in color characteristics. It is described in the whole image of different color the proportion of each color, without regard to the space location place, namely cannot describe the object or object image. Color histogram is especially suitable for describing the very difficult automatic segmentation images [7]. Color histogram based on image retrieval techniques-color feature is one of the important methods, it has the following advantages:

On feature extraction and similarity calculation is simple and efficient;

Do not suffer the scale or rotate the constraint condition such as in space, robustness.

But color histogram also exist three weaknesses:

Color histogram is a kind of mathematical statistics, it can distinguish the result of various colors of the whole image, and the proportion of should be expressed image space information distribution, easy to cause the result is same, but color histogram image difference is large;

Due to the color, the diversification of extracting features, is very easy to create the dimension disasters; 3, due to a lack of unified standards, image in quantitative process so, hard to avoid can appear error.

Generally, in comparison with 2 image eye, good at holding their main color. The so-called mass-tone, refers to the main body color, general image object surface is the image color or background color, occupy larger area. The number of mass-tone general not only have a, but may have multiple main color. Image mass-tone CBIR widely used in. Yang and others propose a MPEG - 7 main color extraction methods, this method has not adopted complex clustering method to extract image mass-tone[8]. The basic idea is first will RGB space of roughly divided into different interval, and then used each interval in the middle as between quantitative results. Mass-tone histogram method: considering the above problems will quantitative histogram produced the mass-tone attune histogram method. For an image, often a few colors cover most of the pixel image, and different colors in images appear probability. Because the images of the color information tonal reflect, mainly through Sun top were used to such person is tonal histogram and mass-tone attune to describe the image histogram color features. Tonal H is expected by color to identify, it USES $^{\circ}$ $\sim 360^{\circ}$ Angle 0 to measure. Close to 0° or 360° tonal is red, 120° nearby tonal it is green, and blue attune is near $^{\circ}$ in 240. In the color wheel, the main color along a round uniformly distributed, some other minor color is located in between main color. In order to obtain an independent of checkpoint image color descriptors, need a kind of color characteristics, in realizing considering shadow, cover and brightness changes under the influence of such factors, it can still independent of object surface shape and observation Angle. It has been found that eye for color tonal particularly sensitive, and color tonal is an independent of check.

2.2 An Improved Algorithm in This Paper

An image color concentration changes the image to create the gap with the original image is not great, but the image color value of each pixel is changed, then two images color histogram difference will be very big. Through the traditional color histogram of retrieval method is difficult to change the color value the two images positioning similar. Therefore, in view of the road traffic signs of semantic particularity, this article first through the extraction of image accents and times mass-tone, will RGB space of roughly divided into different interval, and then used each interval in the middle as between quantitative results, that is to say, first through the form of global histogram image quantification, and then only care about image color frequency appear highest and record highs in color, the image based on both a single screening; Then through the hue histograms, will road traffic signs, the red, yellow, blue as third tonal mass-tone, the most prominent by grasping tonal color as feature extraction, and then, a fast and effective image extraction, and finally reach the purpose of second screening. This paper is using the algorithm combining ideas two reasons: one, take pair of main color histogram extraction is to avoid the when color proportional distribution when the average caused compared by mistake examining phenomenon; Second, take mass-tone histogram method is to avoid image by Lord the influence of regional background region.

3 Shape Feature Extraction

3.1 Image Segmentation and Description Method Based on Contour

The objective is to image segmentation image space division becomes some meaningful area, thus further study separately. Image segmentation normally with adjacent pixels (or "block") the similarity between the judgement basis for major, and according to the specific segmentation problem concerns for other related factors, thus adding in segmentation method for this aspect factor improvement methods (such as resistance noise, brightness extraction, etc.)[9]. Adjacent pixels (or "block" refers to the similarity between the pixels (or "block") to some of the information between hopping degree, this information includes not only the color value, brightness, also including texture information, structure information, etc.)The description method based on contour extraction shapes of contour information only, this kind of description method commonly there are two kinds of forms, one kind is a continuous (namely global type), this description method description method of vector characteristics extracted from a target global, and no target outline for segmentation or block processing; Another kind is discrete (i.e. structural) describe methods. Discrete approach is description method with continuous contrary, namely first will outline is divided into many clips and corresponding feature extraction paragraphs respectively [10].

Fourier shape descriptor

Fourier shape descriptors describing aspects in the image outline used widely, the basic idea is the border to target the discrete Fourier transform (DFT), will transform as the result of its shape after describing due to a Fourier transform is reversible, therefore the nondestructive transformation in conversion process description method will cause the loss of information. Given the Bach J.R.[11] etc of Fourier descriptor description. Contour line of arbitrary pixel defined as the point where the curvature changes in tangent Angle profile of differential arc length change. $K(s)$ set for curvature function, it is:

$$K(s) = \frac{d}{ds} \theta(s) \quad (1)$$

$\theta(s)$ is the tangential point view of the contour

$$\left. \begin{aligned} \theta(s) &= \arctan \left(\frac{y'_s}{x'_s} \right) \\ y'_s &= \frac{dy_s}{ds} \\ x'_s &= \frac{dx_s}{ds} \end{aligned} \right\} \quad (2)$$

Distance from the object boundary point defined as objects center (x_s, y_s) is

$$R(s) = \sqrt{(x_s - x_c)^2 + (y_s - y_c)^2} \tag{3}$$

After coordinates function is to use a plural represent for pixel coordinates

$$Z(s) = (x_s - x_c) + j(y_s - y_c) \tag{4}$$

For complex coordinates functions of Fourier transform the shape of the image will be high frequency and low frequency information is divided into two parts, one with low frequency part portrays the shapes of macroscopic properties, and the high frequency part portrays the shapes of detail features. From the descriptive information, we may safely draw a Fourier shape descriptor. In addition, can ignore the phase information and retain only size information, such doing can keep target's rotating irrelevance. Due to the function of Fourier transformation is symmetrical, that is $|F_{-i}| = |F_i|$, therefore, as for curvature and centric distance, can consider only real part and ignore plural parts. To sum up, based on the curvature of Fourier shape descriptors used the following formula says:

$$f_k = \left[|F_1|, |F_2|, \dots, |F_{M/2}| \right] \tag{5}$$

F_i is Fourier transform parameter of the first component. Based on cancroids distance of Fourier shape descriptor can be expressed as:

$$f_R = \left[\frac{|F_1|}{|F_0|}, \frac{F_2}{F_0}, \dots, \frac{|F_{M/2}|}{|F_0|} \right] \tag{6}$$

For complex coordinates function, due to the positive and negative frequency component was adopted and skip and position of the important parameters related, after the normalized coordinate functions, complex shape the Fourier descriptor can be expressed as:

$$f_Z = \left[\frac{|F_{-(M/2-1)}|}{|F_1|}, \dots, \frac{|F_{-1}|}{|F_1|}, \frac{|F_2|}{|F_1|}, \dots, \frac{|F_{M/2}|}{|F_1|} \right] \tag{7}$$

Wavelet descriptor

By using wavelet transform, the image outline describes first to define the wavelet function clan. Wavelet function definition of family:

$$\psi_{mn}(t) = 2^{-m/2} \psi(2^{-m}t - n) \tag{8}$$

Assuming the outline of image for $f(t)$ function, its wavelet transform coefficient

$$c_{mn}(t) = \int_{-\infty}^{\infty} f(t) \psi_{mn}(t) dt \tag{9}$$

Using the wavelet coefficients can be rebuilt $f(t)$, the reconstruction process for

$$f(t) = \sum_{m=m_0+1}^{\infty} \sum_{n=-\infty}^{\infty} c_{mn} \psi_{mn}(t) + \sum_{m=-\infty}^{m_0} c_{mn} \psi_{mn}(t) \tag{10}$$

M_0 and the accuracy of truncated coefficient related needed. Assume scale functions for, $S_{mn}(t) = 2^{-m/2} S(2^{-m}t - n)$ Combined with wavelet reconstruction formula:

$$f(t) = \sum_{n=-\infty}^{\infty} c_{mn} S_{mn}(t) + \sum_{m=-\infty}^{m_0} \sum_{n=-\infty}^{\infty} c_{mn} S_{mn}(t) \tag{11}$$

In the above one, $c_{mn} S_{mn}(t)$ is called scale coefficients, called the wavelet coefficients of wavelet coefficients, by all composed of contour called wavelet contour descriptor.

Wavelet contour descriptors describing the shape of the high frequency part of overall information, the low frequency part describes the shape of the detail information [11]. Because of using multi-resolution wavelet transform, so in recognition representation method in the process of input image can be analyzed and dynamic adjustment of wavelet transform concrete parameters. But wavelet descriptor dependent on target curve, therefore, the starting point of the same object of wavelet said sampling curve probably because of different starting point to produce very big difference [12, 13].

Use the contour extraction algorithm improved

Aiming at the road traffic signs special semantic features, considering its shape features, using an improved Fourier method, this algorithm descriptive clauses ideas are: to highlight contour cutting Angle of the important position in the retrieval, and ignore other factors in the role of retrieval. $K(s)$ for curvature function, it is as formula (1). $\theta(s)$ is the tangential point view of the contour, which is as formula(2).

That is to say, we only pay attention to theta (s) values, when its y component and X-ray component content, $x^2 + y^2 = R^2$ Judge as round at this time; When theta (s) when the value is 60° , when for the triangle; When theta (s) when the value is 90° , when for square.

4 Experimental Result

Simulation experiment based on Visual c++ 6.0 for development platform, design and realize the road traffic signs recognition simulation system, this system can images low-level image processing work [11].

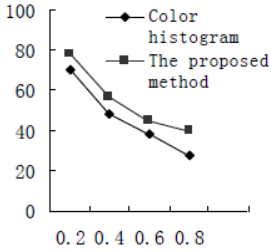


Fig. 1. Two color histogram result contrast chart

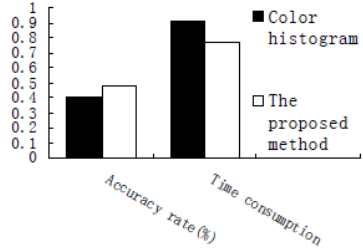


Fig. 2. Two color histogram retrieval precision and response time contrast chart

From Figure 1 and Figure 2, the proposed Color feature extraction method is the higher retrieval result and the lower time consuming in both methods. And shape feature extraction proposed is better than traditional Fourier method from Figure 3 and Figure 4. It can be implementation of an image retrieval system.

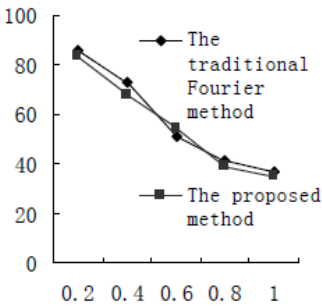


Fig. 3. Two Fourier descriptor retrieval result contrast chart

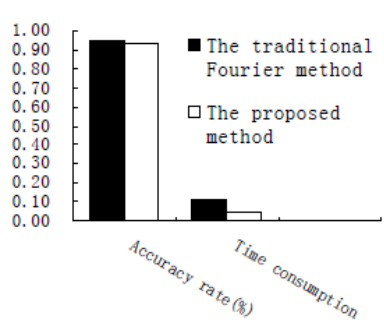


Fig. 4. Two Fourier descriptor retrieval precision and response time contrast chart

5 Conclusion

Based on the semantic features of road traffic signs, improved the general color histogram color feature extraction algorithm, and put forward a kind of based on double mass-tone color histogram improved algorithm; Experiments show that this algorithm in the recall and precision are significantly better than traditional methods in. According to the semantic features of road traffic signs, an improved method based on contour characterization, due to highlight the importance of the parameters such as the curvature, and ignore the importance of the other parameters, and pay attention to the details of consideration, thus puts forward an improved based on Fourier describe the son shape description method of feature extraction. Experiments show that this algorithm in inquires on traditional methods as precision slightly, but in the query efficiency has improved significantly.

Acknowledgment. The corresponding author is Liu Pingping. The authors are grateful to the anonymous reviewers for their insightful comments which have certainly improved this paper. This work is supported by National Natural Science Foundation of China (61101155), Science Foundation of Jilin Educational Committee (2009604) and Plan for Scientific and Technology Development of Jilin Province (20101504).

References

1. Cutsuridis, V.: A Cognitive Model of Saliency, Attention, and Picture Scanning. *Cogn. Comput.*, 292–299 (2009)
2. Doshi, A., Trivedi, M.M.: On the Roles of Eye Gaze and Head Dynamics in Predicting Driver's Intent to Change Lanes. *IEEE Transactions on Intelligent Transportation Systems* 3 (2009)
3. Quattoni, Torralba, A.: Recognizing Indoor Scenes. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR* (2009)
4. Berman, A., Shapiro, L.: Efficient image retrieval with multiple distance measures. In: *SPIE*, vol. 3022, pp. 12–31 (1997)
5. Niblack, W., Barber, R., Equitz, W.: The QBIC project: querying images by content using color texture, and shape. In: *Proceedings of the SPIE Conference on Storage and Retrieval for Image and Video Databases*, San Jose, CA, February 2-3, pp. 173–187 (1993)
6. Flickner, M., Sawhney, H., Niblack, W.: Query by image and video content: the QBIC System. *IEEE Computer*, 23–32 (1995)
7. Bach, J.R., Fuller, C., Gupta, A.: The Virage image search engine: an open framework for image management. In: *Proc. SPIE Storage and Retrieval for Image and Video Database*, pp. 76–87 (1996)
8. Pentland, A., Rosalind, W., Stanley, S.: Photobook: content-based manipulation of image databases. *International Journal of Computer Vision*, 233–254 (1996)
9. Smith, J.R.: Integrated spatial and feature image systems: retrieval, compression and analysis. PhD thesis, Graduate School of Arts and Sciences, Columbia University (1997)
10. Ma, W.Y., Manjunath, B.S.: NETRA: A toolbox for navigating large image database. In: *Proc. of IEEE International Conference on Image Processing*, Santa Barbara, California, USA, pp. 925–928 (1997)
11. Huang, T.S., Mehrotra, S., Ramlchandran, K.: Multimedia analysis and retrieval system (MARS) project. In: *Proc. of 33rd Annual Clinic on Library Application of Data Processing Digital Image Access and Retrieval* (1996)
12. Tang, J., Zhao, J., Xie, Y., Lei, X., Sun, C.: Research of Image Retrieval Based on Affinity Propagation Clustering Algorithm. In: *2010 APCID*, Beijing, China (2010)
13. Yang, H.J., Wang, W.J., Han, J.D.: Image Retrieval Based on Contourlet Texture and Scalable Color Descriptor. In: *PACIA 2010*, Beijing, China (2010)

Existence and Simulations of Periodic Solution for Impulsive Predator-Prey System with Stage Structure for the Predator

Kaihua Wang, Wenxiang Zhang, and Zhanji Gui*

School of Mathematics and Statistics, Hainan Normal University,
Haikou, Hainan, 571158
kaihuawang@qq.com, zhanjigui@sohu.com

Abstract. The principle aim of this paper is to explore the existence of periodic solution of a predator-prey model with stage structure for the predator and impulsive perturbations. Sufficient and realistic conditions are obtained by using Mawhin's continuation theorem of the coincidence degree. Further, some numerical simulations show that our model can occur in many forms of complexities including periodic oscillation and chaotic strange attractor.

Keywords: periodic solution, stage structure, impulses, coincidence degree theory.

1 Introduction

Predator-prey models with stage structure for the predator have received considerable attention in recent years [1-9]. In [1, 8, 9], predator-prey model with stage structure for the predator has been considered as follows:

$$\begin{cases} \dot{x}(t) = x(t)(r - ax(t)) - \frac{bx(t)}{1+mx(t)}y_2(t), \\ \dot{y}_1(t) = k\frac{bx(t)}{1+mx(t)}y_2(t) - (D + d_1)y_1(t), \\ \dot{y}_2(t) = Dy_1(t) - d_2y_2(t). \end{cases}$$

Here $x(t)$, $y_1(t)$, $y_2(t)$ are the densities of prey, immature and mature predators at time t respectively. r is intrinsic growth rate, a is the rate of intra-specific competition, $bx/(1+mx)$ represents the Holling II functional response of the mature predator, which describes how the consumption rate of the predator depends on prey density, k represents the conversion coefficient under the assumption that the reproduction rate of the mature predators is directly proportional to the amount of prey consumed. d_1 and d_2 represent the death rates of immature and mature predators, and D denotes the rate at which immature predators become mature predators.

* Corresponding author.

In this paper, we will consider the following predator-prey model with periodical coefficients and impulsive effects:

$$\begin{cases} \dot{x}(t) = x(t)(r(t) - a(t)x(t)) - \frac{b(t)x(t)}{1+m(t)x(t)}y_2(t), & t \neq t_k, \\ \dot{y}_1(t) = k\frac{b(t)x(t)}{1+m(t)x(t)}y_2(t) - (D(t) + d_1(t))y_1(t), & t \neq t_k, \\ \dot{y}_2(t) = D(t)y_1(t) - d_2(t)y_2(t), & t \neq t_k, \\ \Delta x(t_k) = x(t_k^+) - x(t_k^-) = p_k^1 x(t_k), & t = t_k, \\ \Delta y_1(t_k) = y_1(t_k^+) - y_1(t_k^-) = p_k^2 y_1(t_k), & t = t_k, \\ \Delta y_2(t_k) = y_2(t_k^+) - y_2(t_k^-) = p_k^3 y_2(t_k), & t = t_k, \quad k = 1, 2, \dots, \end{cases} \tag{1}$$

In system (I), we give two hypotheses as follows:

- (H1) $r(t), a(t), b(t), m(t), D(t), d_1(t), d_2(t)$ are continuous positive T -periodic functions;
- (H2) $1 + p_k^i > 0$ are constants and there exists a positive integer m such that $t_{k+m} = t_k + T, p_{k+m}^i = p_k^i (i = 1, 2, 3)$.

2 Basic Concepts and Lemma

Let $J \subset \mathbb{R}$, denote by $PC(J, \mathbb{R})$ the set of functions $\psi : J \rightarrow \mathbb{R}$, which are piecewise continuous in $[0, T]$, and have points of discontinuity $t_n \in [0, T]$, where they are continuous from the left. Let $PC^1(J, \mathbb{R})$ denote the set of functions ψ with derivative $\dot{\psi}(t) \in PC(J, \mathbb{R})$. Throughout this paper we deal with the Banach space of T -periodic functions

$$PC_T = \{\psi \in PC([0, T], \mathbb{R}) \mid \psi(0) = \psi(T)\}$$

with the supremum norm:

$$\|\psi\|_{PC_T} = \sup\{|\psi(t)| : t \in [0, T]\}$$

and

$$PC_T^1 = \{\psi \in PC^1([0, T], \mathbb{R}) \mid \psi(0) = \psi(T)\}$$

with the supremum norm:

$$\|\psi\|_{PC_T^1} = \max\{\|\psi\|_{PC_T}, \|\dot{\psi}\|_{PC_T^1}\}.$$

we will also consider the product space $PC_T \times PC_T$ which is also a Banach space with the norm

$$\|(\psi_1, \psi_2)\|_{PC} = \|\psi_1\|_{PC} + \|\psi_2\|_{PC}.$$

Moreover, for any $y \in C_T$ or $y \in PC_T$, define average value of y as follows: $\bar{y} := \frac{1}{T} \int_0^T y(t) dt$ and the minimum, maximum of y respectively are: $y^L := \min_{t \in [0, T]} y(t), y^M := \max_{t \in [0, T]} y(t)$.

In order to obtain the existence of T -periodic solution to system (II), we must use the following lemma, named as the continuation theorem of coincidence degree theory [10].

Let X, Z be normed vector spaces, $L : \text{Dom}L \subseteq X \rightarrow Z$ be a linear mapping, $N : X \rightarrow Z$ be a continuous mapping. If $\dim \text{Ker}L = \text{comdim} \text{Im}L < +\infty$ and $\text{Im}L$ is closed in Z , then the mapping L will be called a Fredholm mapping of index zero. If L is a Fredholm mapping of index zero, there exist continuous projects $P : X \rightarrow X$ and $Q : Z \rightarrow Z$ such that $\text{Im}P = \text{Ker}L$, $\text{Im}L = \text{Ker}Q = \text{Im}(I - Q)$. It follows that $L|_{\text{Dom}L \cap \text{Ker}P} : (I - P)X \rightarrow \text{Im}L$ has an inverse which is denoted by K_P . If Ω is an open bounded subset of X , the mapping N will be called L -compact on $\overline{\Omega}$ provided that $QN(\overline{\Omega})$ is bounded and $K_P(I - Q)N : \overline{\Omega} \rightarrow X$ is compact. Since $\text{Im}Q$ is isomorphic to $\text{Ker}L$ there exists an isomorphism $F : \text{Im}Q \rightarrow \text{Ker}L$.

Lemma 1. *Let L be a Fredholm mapping of index zero and N be L -compact on $\overline{\Omega}$. Suppose that*

- (a) *For each $\lambda \in (0, 1)$, every solution x of $Lx = \lambda Nx$ such that $x \notin \partial\Omega$;*
- (b) *$QNx \neq 0$ for each $x \in \text{Ker}L \cap \partial\Omega$;*
- (c) *$\text{deg} \{FQN, \Omega \cap \text{Ker}L, 0\} \neq 0$.*

Then the equation $Lx = Nx$ has at least one solution lying in $\text{Dom}L \cap \overline{\Omega}$.

3 Existence of Periodic Solution

In this section, we study the existence of positive periodic solution to (II).

Theorem 1. *If system (I) satisfies*

- (H3) $kb^M D^M - (D^L + d_1^L)d_2^L m^L > 0$
- (H4) $\min_{x \in [B_1, A_1]} \{-a^M m^L x^2 + (m^L r^L - a^M)x + r^L\} > 0$

Then system (I) has at least one positive T -periodic solution.

Here $A_1 = \ln \frac{r^M}{a^L}$, $B_1 = F_1 - I_1 - \sum_{k=1}^m |\ln(1 + p_k^1)|$

$F_1 = \ln \frac{d_2^L D^L + d_2^L d_1^L}{kb^M D^M - (D^L + d_1^L)d_2^L m^L}$,

$I_1 = 2\bar{r}T + \sum_{k=1}^m \ln(1 + p_k^1) + \sum_{k=1}^m |\ln(1 + p_k^1)|$.

Proof. Let $x(t) = e^{u_1(t)}$, $y_1(t) = e^{u_2(t)}$, $y_2(t) = e^{u_3(t)}$ then system (II) is reformulated as

$$\begin{cases} \dot{u}_1(t) = r(t) - a(t) \exp\{u_1(t)\} - \frac{b(t) \exp\{u_3(t)\}}{1 + m(t) \exp\{u_1(t)\}}, \\ \dot{u}_2(t) = k \frac{b(t) \exp\{u_1(t) + u_3(t) - u_2(t)\}}{1 + m(t) \exp\{u_1(t)\}} - D(t) - d_1(t), \\ \dot{u}_3(t) = D(t) \exp\{u_2(t) - u_3(t)\} - d_2(t), \\ \Delta u_i(t_k) = \ln(1 + p_k^i), \quad i = 1, 2, 3. \end{cases} \quad (2)$$

If system (2) has a T -periodic solution $(u_1(t), u_2(t), u_3(t))^T$, then

$$(x_1^*(t), x_2^*(t), y^*(t))^T = (e^{u_1(t)}, e^{u_2(t)}, e^{u_3(t)})^T$$

is a positive T -periodic solution to system (1). So, in the following, we discuss the existence of T -periodic solution to system (2).

In order to use Lemma 1, we set $\mathbf{u} = (u_1(t), u_2(t), u_3(t))^T$. Define $X = \{x \in PC(R, R^3) : x(t + T) = x(t)\}$, $Z = X \times R^{3m}$, then it is standard to show both X and Z are Banach space when they are endowed with the norms $\|x\|_c = \sup_{t \in [0, T]} |x(t)|$ and $\|(x, c_1, c_2, c_3)\| = (\|x\|_c^2 + |c_1|^2 + |c_2|^2 + |c_3|^2)^{1/2}$.

Let $\text{Dom}L \subset X = \{x \in C^1 [0, T; t_1, \dots, t_m] \mid x(0) = x(T)\}$, $L: \text{Dom}L \rightarrow Z$, $L\mathbf{u} = (\mathbf{u}', \Delta\mathbf{u}(t_1), \dots, \Delta\mathbf{u}(t_m))$; $N : X \rightarrow Z$, $N: \text{Dom}L \rightarrow Z$, $N\mathbf{u} = (\mathbf{u}', \Delta\mathbf{u}(t_1), \dots, \Delta\mathbf{u}(t_m))$. It is easy to prove that L is a Fredholm mapping of index zero.

Consider the operator equation

$$L\mathbf{u} = \lambda N\mathbf{u}, \quad \lambda \in (0, 1). \tag{3}$$

Suppose that $\mathbf{u}(t) = (u_1(t), u_2(t), u_3(t))^T$ is a periodic solution of (3) for certain $\lambda \in (0, 1)$. Integrating (3) over the interval $[0, T]$, we obtain

$$\begin{cases} \bar{r}T + \sum_{k=1}^m \ln(1 + p_k^1) = \int_0^T a(t) \exp\{u_1(t)\} + \frac{b(t) \exp\{u_3(t)\}}{1+m(t) \exp\{u_1(t)\}} dt, \\ \bar{D}T + \bar{d}_1T - \sum_{k=1}^m \ln(1 + p_k^2) = \int_0^T k \frac{b(t) \exp\{u_1(t)+u_3(t)-u_2(t)\}}{1+m(t) \exp\{u_1(t)\}} dt, \\ \bar{d}_2T - \sum_{k=1}^m \ln(1 + p_k^3) = \int_0^T D(t) \exp\{u_2(t) - u_3(t)\} dt. \end{cases} \tag{4}$$

From (2) and (4), we have

$$\begin{cases} \int_0^T |\dot{u}_1(t)| dt \leq 2\bar{r}T + \sum_{k=1}^m \ln(1 + p_k^1) + \sum_{k=1}^m |\ln(1 + p_k^1)| = I_1, \\ \int_0^T |\dot{u}_2(t)| dt \leq 2\bar{D}T - \sum_{k=1}^m \ln(1 + p_k^2) + \sum_{k=1}^m |\ln(1 + p_k^2)| = I_2, \\ \int_0^T |\dot{u}_3(t)| dt \leq 2\bar{d}_2T - \sum_{k=1}^m \ln(1 + p_k^3) + \sum_{k=1}^m |\ln(1 + p_k^3)| = I_3. \end{cases} \tag{5}$$

Since $u_i(t) \in PC_T$, there exist $\xi_i, \eta_i \in [0, T]$, ($i = 1, 2, 3$) such that

$$u_i(\xi_i) = \min_{t \in [0, T]} u_i(t), \quad u_i(\eta_i) = \max_{t \in [0, T]} u_i(t).$$

It is clear that $\dot{u}_i(\xi_i) = 0, \dot{u}_i(\eta_i) = 0$. From this and system (2), we obtain

$$\begin{cases} r(\eta_1) - a(\eta_1) \exp\{u_1(\eta_1)\} - \frac{b(\eta_1) \exp\{u_3(\eta_1)\}}{1+m(\eta_1) \exp\{u_1(\eta_1)\}} = 0, \\ k \frac{b(\eta_2) \exp\{u_1(\eta_2)+u_3(\eta_2)-u_2(\eta_2)\}}{1+m(\eta_2) \exp\{u_1(\eta_2)\}} - D(\eta_2) - d_1(\eta_2) = 0, \\ D(\eta_3) \exp\{u_2(\eta_3) - u_3(\eta_3)\} - d_2(\eta_3) = 0 \end{cases} \tag{6}$$

and

$$\begin{cases} r(\xi_1) - a(\xi_1) \exp\{u_1(\xi_1)\} - \frac{b(\xi_1) \exp\{u_3(\xi_1)\}}{1+m(\xi_1) \exp\{u_1(\xi_1)\}} = 0, \\ k \frac{b(\xi_2) \exp\{u_1(\xi_2)+u_3(\xi_2)-u_2(\xi_2)\}}{1+m(\xi_2) \exp\{u_1(\xi_2)\}} - D(\xi_2) - d_1(\xi_2) = 0, \\ D(\xi_3) \exp\{u_2(\xi_3) - u_3(\xi_3)\} - d_2(\xi_3) = 0 \end{cases} \tag{7}$$

From (6)₁, we have

$$a^L \exp\{u_1(\eta_1)\} \leq a(\eta_1) \exp\{u_1(\eta_1)\} \leq r(\eta_1) \leq r^M.$$

then

$$u_1(\eta_1) \leq \ln \frac{r^M}{a^L} = A_1. \tag{8}$$

From (7)₁, we have

$$\frac{b^L \exp\{u_3(\xi_3)\}}{1 + m^M \exp\{A_1\}} \leq \frac{b(\xi_1) \exp\{u_3(\xi_1)\}}{1 + m(\xi_1) \exp\{u_1(\xi_1)\}} \leq r(\xi_1) \leq r^M,$$

then

$$u_3(\xi_3) \leq \ln \frac{r^M + r^M m^M \exp\{A_1\}}{b^L} = E_3.$$

Thus we get

$$u_3(t) \leq u_3(\xi_3) + \int_0^T |\dot{u}_3(t)| dt + \sum_{k=0}^m |\ln(1 + p_k^3)| \leq E_3 + I_3 + \sum_{k=0}^m |\ln(1 + p_k^3)| = A_3. \tag{9}$$

From (7)₃, we have

$$\exp\{u_2(\xi_2)\} \leq \exp\{u_2(\xi_3)\} = \frac{d_2(\xi_3) \exp\{u_3(\xi_3)\}}{D(\xi_3)} \leq \frac{d_2^M \exp\{A_3\}}{D^L},$$

then

$$u_2(\xi_2) \leq \ln \frac{d_2^M \exp\{A_3\}}{D^L} = E_2.$$

Thus we get

$$u_2(t) \leq E_2 + I_2 + \sum_{k=0}^m |\ln(1 + p_k^2)| = A_2. \tag{10}$$

From (6)₃, we have

$$\exp\{u_3(\eta_2) - u_2(\eta_2)\} \leq \exp\{u_3(\eta_3) - u_2(\eta_3)\} \leq \frac{D^M}{d_2^L}.$$

From (6)₂, we have

$$\begin{aligned} D^L + d_1^L &\leq \frac{kb(\eta_2) \exp\{u_1(\eta_2) + u_3(\eta_2) - u_2(\eta_2)\}}{1 + m(\eta_2) \exp\{u_1(\eta_2)\}} \leq \frac{kb^M \exp\{u_1(\eta_1) + u_3(\eta_2) - u_2(\eta_2)\}}{1 + m^L \exp\{u_1(\eta_1)\}} \\ &\leq \frac{kb^M D^M \exp\{u_1(\eta_1)\}}{d_2^L (1 + m^L \exp\{u_1(\eta_1)\})}, \end{aligned}$$

Because of (H3), we have

$$u_1(\eta_1) \geq \ln \frac{d_2^L D^L + d_2^L d_1^L}{kb^M D^M - (D^L + d_1^L) d_2^L m^L} = F_1.$$

Thus we get

$$u_1(t) \geq u_1(\eta_1) - \int_0^T |\dot{u}_1(t)| dt - \sum_{k=1}^m |\ln(1 + p_k^1)| \geq F_1 - I_1 - \sum_{k=1}^m |\ln(1 + p_k^1)| = B_1. \tag{11}$$

From (6)₁, we have

$$\frac{b^M \exp\{u_3(\eta_3)\}}{1 + m^L \exp\{u_1(\eta_1)\}} \geq \frac{b(\eta_1) \exp\{u_3(\eta_1)\}}{1 + m(\eta_1) \exp\{u_1(\eta_1)\}} \geq r^L - a^M \exp\{u_1(\eta_1)\},$$

that is

$$b^M \exp\{u_3(\eta_3)\} \geq -a^M m^L \exp\{2u_1(\eta_1)\} + (m^L r^L - a^M) \exp\{u_1(\eta_1)\} + r^L.$$

F denotes the minimum of function $f(x) = -a^M m^L x^2 + (m^L r^L - a^M)x + r^L$ as $x \in [B_1, A_1]$. Because of (H)₄, we have

$$u_3(\eta_3) \geq \ln \frac{F}{b^M} = F_3.$$

Thus we get

$$u_3(t) \geq F_3 - I_3 - \sum_{k=1}^m |\ln(1 + p_k^3)| = B_3. \tag{12}$$

From (6)₃, we have

$$\exp\{u_2(\eta_2)\} \geq \frac{d_2(\eta_3) \exp\{u_3(\eta_3)\}}{D(\eta_3)} \geq \frac{d_2^L \exp\{B_3\}}{D^M}.$$

then

$$u_2(\eta_2) \geq \ln \frac{d_2^L \exp\{B_3\}}{D^M} = F_2.$$

Thus we get

$$u_2(t) \geq F_2 - I_2 - \sum_{k=1}^m |\ln(1 + p_k^2)| = B_2. \tag{13}$$

Now, we have

$$\begin{aligned} \sup_{t \in [0, T]} |u_1(t)| &\leq \max\{|A_1|, |B_1|\} = N_1, \\ \sup_{t \in [0, T]} |u_2(t)| &\leq \max\{|A_2|, |B_2|\} = N_2, \\ \sup_{t \in [0, T]} |u_3(t)| &\leq \max\{|A_3|, |B_3|\} = N_3. \end{aligned}$$

Obviously, there exists a constant $N_4 > 0$ such that $\max\{|u_1|, |u_2|, |u_3|\} < N_4$. Take $r > N_1 + N_2 + N_3 + N_4$, $\Omega = \{x \in X \mid \|x\|_c < r\}$, then Ω is L -compact on $\overline{\Omega}$. So, for $\forall \mathbf{u} = (u_1, u_2, u_3)^T \in \partial\Omega \cap \text{Ker}L$, we have $QN\mathbf{u} \neq 0$. Let $J : \text{Im}Q \rightarrow x, (d, 0, \dots, 0) \rightarrow d$. When $\mathbf{u} \in \Omega \cap \text{Ker}L$, in view of the assumptions in Mawhin's continuation theorem, one obtains, $\text{deg}\{FQN, \Omega \cap \text{Ker}L, 0\} \neq 0$. By now we have proved that Ω satisfies all the requirements in Mawhin's continuation theorem. Hence, (II) has at least one T -periodic solution in $\text{Dom}L \cap \overline{\Omega}$. \square

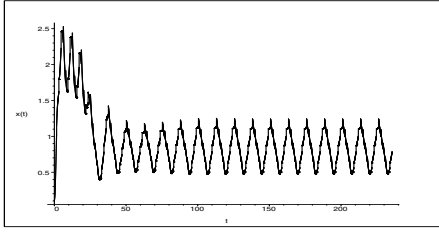


Fig. 1. Time-series of $x(t)$ evolved in system (II) with $\omega = \pi/2$

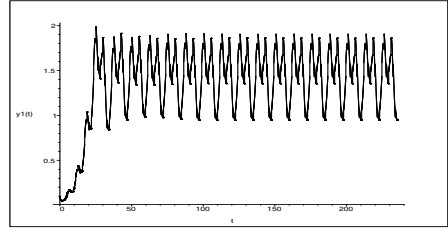


Fig. 2. Time-series of $y_1(t)$ evolved in system (II) with $\omega = \pi/2$

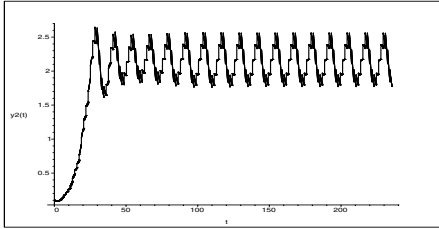


Fig. 3. Time-series of $y_2(t)$ evolved in system (II) with $\omega = \pi/2$

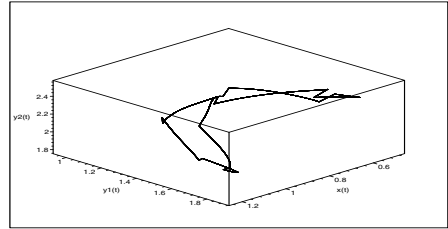


Fig. 4. Phase portrait of 2π -periodic solution of system (II) with $\omega = \pi/2$

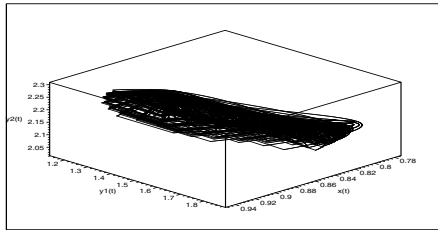


Fig. 5. Phase portrait of chaotic strange attractor of system (II) with $\omega = 2$

4 Some Simulations

In this section, we shall discuss an example to illustrate main results. For system (II), we take: $t_n = n\omega$, $r(t) = 2 + 0.2 \cos(t)$, $1 + 0.2 \sin(t)$, $b(t) = 1.2 - 0.2 \cos(t)$, $m(t) = 1.4 + 0.4 \sin(t)$, $D(t) = 0.8 + 0.1 \sin(t)$, $d_1(t) = 0.6 + 0.1 \cos(t)$, $d_2(t) = 0.6 + 0.1 \sin(t)$, $p_k^i = -0.1$. Obviously, all conditions of Theorem 1 are satisfied.

If $\omega = \pi/2$, then system (II) under the above conditions has a unique 2π -periodic solution (see Fig.1-Fig.4, where we take $[x_1(0), x_2(0), y(0)]^T = [0.1, 0.1, 0.1]^T$). We find the occurrence of sudden changes in the figures of the time-series and phase portrait. The influence of pulse is obvious.

If $\omega = 2$, then (H2) is not satisfied. Periodic oscillation of system (II) under the above conditions will be destroyed by impulsive effect. Numeric results (see Fig.5) show that system (II) under the above conditions has Gui chaotic strange attractor [11].

5 Conclusion

In this paper, we consider the existence of periodic solution of a predator-prey model with stage structure for the predator and impulsive perturbations. Sufficient and realistic conditions are obtained by using Mawhin's continuation theorem of the coincidence degree. Further, some numerical simulations show that our model can occur in many forms of complexities including periodic oscillation and chaotic strange attractor.

Acknowledgements. This work is supported jointly by the Natural Sciences Foundation of China under Grant No. 60963025, Natural Sciences Foundation of Hainan Province under Grant No. 110007 and the Start-up fund of Hainan Normal University under Project No. 00203020201.

References

1. Georgesc, P., Hsien, Y.H.: Global dynamics of a predator-prey model with stage structure for the predator. *Society for Industrial and Applied Mathematics* 67, 1379–1395 (2007)
2. Chen, L.J., Chen, F.D.: A stage-structure and harvesting predator-prey system. *Annals of Differential Equations* 03, 293–301 (2011)
3. Zhao, M., Cheng, R.F.: Positive Periodic Solution in a Ratio-Based Predator-Prey System with Stage Structure for Predator. *Journal of Biomathematics* 25, 88–96 (2010)
4. Lai, W.Y., Zhan, Q.Y.: Permanence and Global Stability of a Predator-prey System with Stage Structure and Time Delays. *College Mathematics* 28, 50–57 (2012)
5. Ma, Z.H., Wang, S.F., Wang, W.T., Li, Z.Z.: Permanence of a stage-structured predator-prey system with a class of functional responses. *Comptes Rendus Biologies* 334, 851–854 (2011)
6. Zhang, Z.Q., Wang, Z.C.: Periodic Solutions of a Predator-Prey System with Stage-Structure for Predator. *Acta Mathematica Sinica* 48, 541–548 (2005)
7. Ling, L.S.: Stability and Hopf Bifurcation in a Predator-prey Model with Stage Structure and Time Delay. *Mathematica Applicata* 25, 131–139 (2012)
8. Wang, W.: Global dynamics of a population model with stage structure for predator. In: *Advanced Topics in Biomathematics*, pp. 253–257. World Scientific, River Edge (1997)
9. Wang, W., Chen, L.: A predator-prey system with stage structure for predator. *Comput. Math. Appl.* 33, 83–91 (1997)
10. Gaines, R.E., Mawhin, J.L.: *Coincidence degree and nonlinear differential equations*. Springer, Berlin (1977)
11. Zhang, J., Gui, J.: Periodic solutions of nonautonomous cellular neural networks with impulses and delays. *Nonlinear Analysis: Real World Applications* 19, 1891–1903 (2009)

Dynamics and Simulations of Multi-species Competition-Predator System with Impulsive

Yan Yan, Kaihua Wang, and Zhanji Gui*

School of Mathematics and Statistics, Hainan Normal University,
Haikou, Hainan, 571158, P.R. China
oishi19840923@163.com, zhanjigui@sohu.com

Abstract. We investigate the dynamics of a class of multi-species competition predator interaction models with Beddington-DeAngelis functional response. Sufficient conditions for existence of a positive periodic solution are given and sufficient criteria are established for the global stability and the globally exponential stability of the system by using the comparison principle and the Lyapunov method. In addition, some numerical simulation shows that our models can occur in many forms of complexities including periodic oscillation and strange chaotic strange attractor.

Keywords: Global stability, Simulations, Competition-predator.

1 Introduction

It is well known that the traditional predator-prey systems with prey-dependent functional response fail to model the interference among predator. To overcome the shortcoming Arditi and Ginzburg proposed the ratio-dependent predator-prey model which is depicted as follows [1]:

$$\begin{cases} x' = x(a - bx) - cxy/(my + x), \\ y' = y(-d + fx)/(my + x), \end{cases}$$

which incorporates mutual interference by predators. However, it has somewhat singular behaviors at low densities and has been criticized on other grounds. See [2] for a mathematical analysis and the references in [3] for some aspects of the debate among biologists about ratio-dependence. The Beddington-DeAngelis form of functional response has some of the same qualitative features as the ratio-dependent models form but avoids some of the same behavior of ratio-dependent models at low densities. Hence it seems worth further study. For a thorough biological background to the model, we refer to [3-8].

Although much progress has been seen in the studied of predator-prey models with the Beddington-DeAngelis functional response. In the real world, any biological or environmental parameters are naturally subject to fluctuation in

* Corresponding author.

time, so it is reasonable to study the corresponding nonautonomous system. In this paper, we will study the predator-prey system with impulsive perturbations responses, we obtain the system:

$$\begin{cases} \dot{x}_i(t) = x_i(t) \left[b_i(t) - \sum_{k=n+1}^{n+m} \frac{c_{ik}(t)y_k(t)}{\alpha_{ik}(t)+\beta_{ik}(t)x_i(t)+\gamma_{ik}(t)y_k(t)} \right. \\ \quad \left. - \sum_{k=1}^n a_{ik}(t)x_k(t) \right], \\ \dot{y}_j(t) = y_j(t) \left[-r_j(t) + \sum_{k=1}^n \frac{d_{jk}(t)x_k(t)}{\alpha_{jk}(t)+\beta_{jk}(t)x_k(t)+\gamma_{jk}(t)y_j(t)} \right. \\ \quad \left. - \sum_{k=n+1}^{n+m} \delta_{jk}(t)y_k(t) \right], \\ t \neq t_k (k \in \mathbb{N}^+), i = 1, \dots, n, j = n + 1, \dots, n + m, \\ \Delta x_i(t) = x_i(t^+) - x_i(t^-) = (b_{ik} + h_{ik})x_i(t), \\ \Delta y_j(t) = y_j(t^+) - y_j(t^-) = (b_{jk} + h_{jk})y_j(t), \\ t = t_k (k \in \mathbb{N}^+), i = 1, \dots, n, j = n + 1, \dots, n + m, \end{cases} \quad (1)$$

where

- $x_i(t) (i = 1, 2, \dots, n)$ denote the densities of prey species at time t , respectively;
- $y_j(t) (j = n + 1, n + 2, \dots, n + m)$ denote the density of predator species at time t , respectively;
- b_{ik} and b_{jk} represent the birth rate of $x_i(t)$ and $y_j(t)$ at time t , respectively;
- h_{ik} and h_{jk} represent the harvesting (stocking) rate of $y_j(t)$ at time t , respectively. When $h_{ik}, h_{jk} > 0$, it stands for harvesting, while $h_{ik}, h_{jk} < 0$ means stocking;
- $x_i(t^+)$ and $x_i(t^-)$ represent the right and left limits of $x_i(t)$ at t , $y_j(t^+)$ and $y_j(t^-)$ represent the right and left limits of $y_j(t)$ at t .

The ranges of the indices $i \in \{1, 2, \dots, n\}$ and $j \in \{n + 1, \dots, n + m\}$ are used in this paper unless otherwise stated. Throughout the paper, we give the hypothesis as follows.

- (H₁) For any $t \in \mathbb{R}$, $b_i(t), a_{ik}(t), r_j(t), d_{jk}(t), \alpha_{jk}(t), \beta_{jk}(t), \gamma_{jk}(t), (k = 1, \dots, n), \delta_{jk}(t), c_{ik}(t), \alpha_{ik}(t), \beta_{ik}(t), \gamma_{ik}(t), (k = n + 1, \dots, n + m)$, are nonnegative continuous T-periodic functions.
- (H₂) $b_{ik}, b_{jk} > 0, 1 + b_{ik} + h_{ik} > 0, b_{jk} > 0, 1 + b_{jk} + h_{jk} > 0, b_{ik}, h_{ik}, b_{jk}, h_{jk}, (k \in \mathbb{N}^+)$ are constants. There exists a positive integer q , such that $t_{k+q} = t_k + T, b_{i(k+q)} = b_{ik}, b_{j(k+q)} = b_{jk}, h_{i(k+q)} = h_{ik}, h_{j(k+q)} = h_{jk} (k \in \mathbb{N}^+)$. Without loss of generality, we also suppose that $t_k \neq 0$ and $[0, T] \cap \{t_k | k \in \mathbb{N}^+\} = \{t_1, t_2, \dots, t_s\}$, then it follows that $q = s$.
- (H₃) $x_i(t), y_j(t)$ is left-continuous at t_k , i.e., the following relations are satisfied:
 $x_i(t_k^-) = x_i(t_k), x_i(t_k^+) = (1 + b_{ik} + h_{ik})x_i(t_k), k \in \mathbb{N}^+,$
 $y_j(t_k^-) = y_j(t_k), y_j(t_k^+) = (1 + b_{jk} + h_{jk})y_j(t_k), k \in \mathbb{N}^+.$
- (H₄) $t_1 < t_2 < \dots$ and $\lim_{k \rightarrow \infty} t_k = \infty$.
- (H₅) $x_i(t_0^+) > 0, y_j(t_0^+) > 0$.

2 Existence of Positive Periodic Solution and an Illustrative Example

In this section, we denote $\Delta_i = \bar{b}_i + \frac{1}{T} \sum_{k=1}^q \ln(1 + b_{ik} + h_{ik})$, $\Delta_j = -\bar{r}_j + \frac{1}{T} \sum_{k=1}^q \ln(1 + b_{jk} + h_{jk})$, $\mathbf{B} = \begin{pmatrix} \bar{a}_{i^*j^*}, & \mathbf{0} \\ \mathbf{0}, & \bar{\delta}_{i'j'} \end{pmatrix}$, $\Delta = \begin{pmatrix} \Delta_i \\ \Delta_j \end{pmatrix}$,

$$H_j = \max \left\{ \ln \left[\frac{\Delta_j + \sum_{k=1}^n \frac{\bar{d}_{jk}/\beta_{jk}}{\delta_{jk}} \right] \right\} + \sum_{k=1}^q \ln(1 + b_{jk} + h_{jk}) + 2\bar{r}_j T,$$

$$H_i = \max \left\{ \ln \left[\frac{\Delta_i}{\bar{a}_{ii}} \right] \right\} + 3 \sum_{k=1}^q \ln(1 + b_{ik} + h_{ik}) + 2\bar{b}_i T,$$

where $i^*, j^* \in 1, \dots, n$, $i', j' \in n + 1, \dots, n + m$. We also denote by $\mathbf{B}_k (k \in 1, \dots, n + m)$ the matrix obtained by replacing the k th column of with Δ .

Now, we can obtain the sufficient conditions for existence of a positive periodic solution of system (II) by using a continuation theorem in coincidence degree theory (see [9]). Use the method similar to [10–12], we can easily get the following sufficient conditions:

Theorem 1. *If system (I) satisfies (H_1) – (H_5) as well as the following conditions*

$$(H_6) \quad \Delta_i - \sum_{k=n+1}^{n+m} \bar{c}_{ik} \gamma_{ik} > \sum_{k=1, k \neq i}^n \bar{a}_{ik} \exp\{H_i\},$$

$$(H_7) \quad \Delta_j > \sum_{k=n+1, k \neq j}^{n+m} \bar{\delta}_{jk} \exp\{H_j\},$$

$$(H_8) \quad \mathbf{B} > 0, \mathbf{B}_k, k \in (1, \dots, n + m).$$

Then system (I) has at least one positive T -periodic solution.

We shall discuss an example to illustrate this result. In (II), we take $b_1 = 6 + \sin t$; $a_{11} = 0.4 - 0.2 * \cos t$; $c_{12} = 0.4 - 0.18 * \sin t$; $r_1 = 0.5 + 0.1 * \sin t$; $r_2 = 0.2 + 0.1 * \cos t$; $d_{11} = 2 + \sin t$; $d_{21} = 3 + \cos t$; $\delta_{12} = 2 + 0.1 * \sin t$. $\delta_{22} = 5 + 0.1 * \cos t$; $\alpha_{11} = 4 + \cos t$; $\alpha_{12} = 8 - \cos t$; $\alpha_{11} = 4 + \cos t$; $\beta_{11} = 3 + \cos t$; $\beta_{12} = 5 + \sin t$; $\beta_{21} = 2 + \cos t$; $\gamma_{11} = 3 + \sin t$; $\gamma_{12} = 3 + \sin t$; $\gamma_{21} = 4 + \sin t$; $p_1 = 0.1$; $q_1 = 0.2$; $q_3 = 0.3$.

If $T = 2\pi$, all the sufficient conditions $(H_1) - (H_8)$ are satisfied. Then system (II) under the above conditions has a unique periodic solution (See Fig.1-Fig.4, where $[x(t), y_1(t), y_2(t)]^T = [0.1, 0.1, 0.1]^T$). The influence of pulse is obvious.

If $T = 2$, Periodic oscillation of system (II) will be destroyed by impulsive effect. Numeric results (see Fig.5) show that system (II) has Gui chaotic strange attractors [15].

3 Global Stability and Globally Exponential Stability of Solutions

Let $\mathbf{x}_0 = (x_{10}, \dots, x_{n0}, y_{n+1,0}, \dots, y_{n+m,0})$ and $x_{i0}, y_{j0} \in \mathbb{R}_+$. We denote by $\mathbf{x}(t) = \mathbf{x}(t; t_0, \mathbf{x}_0) = (x_1(t), \dots, x_n(t), y_{n+1}(t), \dots, y_{m+n}(t))$ the solution of system (II) satisfying the initial conditions $\mathbf{x}(t_0 + 0; t_0, \mathbf{x}_0) = \mathbf{x}_0$ and by $J = J(t_0, \mathbf{x}_0)$ the maximal interval of type $[t_0, \eta)$ in which the solution $\mathbf{x}(t; t_0, \mathbf{x}_0)$ is defined.

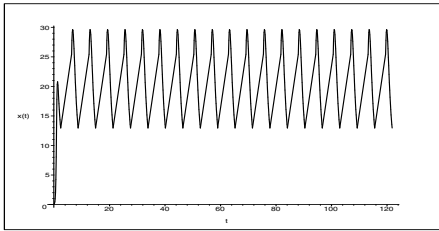


Fig. 1. Time-series of $x(t)$ evolved in system (I)

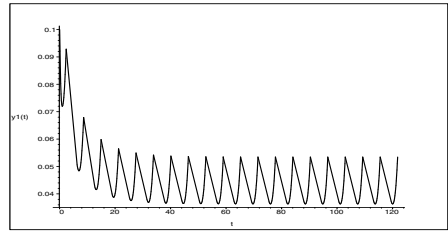


Fig. 2. Time-series of $y_1(t)$ evolved in system (I)

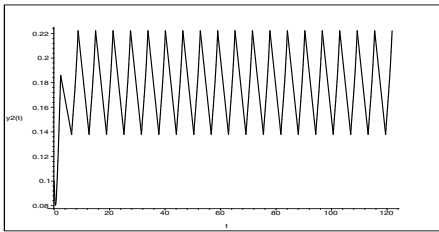


Fig. 3. Time-series of $y_2(t)$ evolved in system (I)

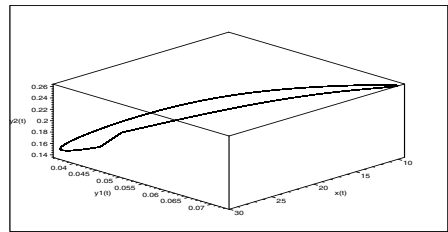


Fig. 4. Phase portrait of a 2π periodic solution of system (I)

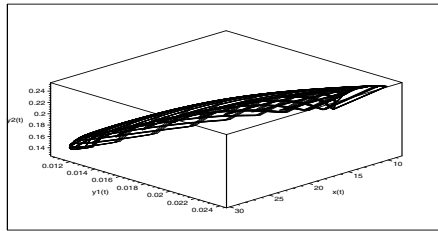


Fig. 5. Phase portrait of a 2 periodic solution of system (I)

Let $\mathbf{x}(t; t_0, \mathbf{x}_0), \mathbf{x}^*(t) = (x_1(t), \dots, x_n(t), y_{n+1}(t), \dots, y_{n+m}(t))$ and $\mathbf{x}^*(t; t_0, \mathbf{x}_0^*), \mathbf{x}^*(t) = (x_1^*(t), \dots, x_n^*(t), y_{n+1}^*(t), \dots, y_{n+m}^*(t))$ be any two solutions of (I) with initial conditions $\mathbf{x}(t_0+0; t_0, \mathbf{x}_0) = \mathbf{x}_0, \mathbf{x}^*(t_0+0; t_0, \mathbf{x}_0^*) = \mathbf{x}_0^*, t_0 \in \mathbb{R}_+$.

We put forward two definitions in [13], [14].

Definition 1. (Ahmad, [13]). The system (I) is said to be

- (a) globally stable if for all $\varepsilon > 0$, there exists $\delta = \delta(\varepsilon, t_0) > 0$ such that if $\mathbf{x}(t), \mathbf{x}^*(t) \in \mathbb{R}_+^{n+m}$, with $\|\mathbf{x}_0 - \mathbf{x}_0^*\| \leq \delta$ then for all $t \geq t_0, \|\mathbf{x}(t; t_0, \mathbf{x}_0) - \mathbf{x}^*(t; t_0, \mathbf{x}_0^*)\| < \varepsilon$;
- (b) globally asymptotically stable if it is globally stable and $\lim_{t \rightarrow \infty} \|\mathbf{x}(t; t_0, \mathbf{x}_0) - \mathbf{x}^*(t; t_0, \mathbf{x}_0^*)\| = 0$.
- (c) globally exponentially stable if for all $\alpha > 0$, there exists $\gamma = \gamma(\alpha) > 0$ such that $\mathbf{x}(t), \mathbf{x}^*(t) \in \mathbb{R}_+^{n+m}$, with $\|\mathbf{x}_0 - \mathbf{x}_0^*\| \leq \alpha$ then for all $t \geq t_0, \|\mathbf{x}(t; t_0, \mathbf{x}_0) - \mathbf{x}^*(t; t_0, \mathbf{x}_0^*)\| < \gamma \|\mathbf{x}_0 - \mathbf{x}_0^*\| \exp\{-\varphi(t - t_0)\}$.

Definition 2. (Ahmad, [13], [14]). We say that the function $V(t, \mathbf{x})$, $V : [t_0, \infty) \times \mathbb{R}_+^{n+m}$, belongs to the class V_0 if the following conditions are satisfied:

- (1) The function V is continuous in $\bigcup_{i=1}^{\infty} G_i$ and $V(t, 0) = 0$ for $t \in [t_0, \infty)$.
- (2) The function V satisfies locally the Lipschitz condition with respect to \mathbf{x} on each of the sets G_i .
- (3) For each $k \in \mathbb{N}^+$ there exist the finite limits $V(t_k - 0, \mathbf{x}) = \lim_{t \rightarrow t_k, t < t_k} V(t, \mathbf{x})$, $V(t_k + 0, \mathbf{x}) = \lim_{t \rightarrow t_k, t > t_k} V(t, \mathbf{x})$.
- (4) For each $k \in \mathbb{N}^+$ the following equalities are valid: $V(t_{k-1}, \mathbf{x}) = V(t_k, \mathbf{x})$.

We can easily prove the following two Lemmas, which will be used to prove our main theorems.

Lemma 1. Suppose the hypotheses (H_1) - (H_5) hold. There exist functions $P_i, Q_i, P_j, Q_j \in \mathbb{R}^{n+m}$ such that $P_i(t) \leq x_i(t) \leq Q_i(t)$, $P_j(t) \leq y_j(t) \leq Q_j(t)$ for all $t \geq t_0$.

Lemma 2. Suppose the hypotheses (H_1) - (H_5) hold.

$\mathcal{X}(t) = \mathcal{X}(t; t_0, \mathbf{x}_0) = (x_1(t), \dots, x_n(t), y_{n+1}(t), \dots, y_{n+m}(t))$ is a solution of (1), then there exist positive constants $\tau_i, \theta_i, \tau_j, \theta_j$ such that $\tau_i \leq x_i(t) \leq \theta_i$, $\tau_j \leq y_j(t) \leq \theta_j$, for all $t \in (t_{k-1}, t_k], k \in \mathbb{N}^+$ and if in addition $0 < 1 + b_{ik} + h_{ik} \leq 1$, $0 < 1 + b_{jk} + h_{jk} \leq 1$, then $\tau_i \leq x_i(t) \leq \theta_i$, $\tau_j \leq y_j(t) \leq \theta_j$, for all $t \in J$.

We introduce the following notations: $G_k = (t_{k-1}, t_k) \times \mathbb{R}_+^{n+m}$, $k \in \mathbb{N}^+$, $G = \bigcup_{k=1}^{\infty} G_k$. Let $V \in V_0$, for any $(t, \mathcal{X}) \in [t_{k-1}, t_k) \times \mathbb{R}_+^{n+m}$, the right-hand derivative $D^+V(t, \mathcal{X}(t))$ along the solution $\mathcal{X}(t; t_0, x_0)$ of (1) is defined by $D^+V(t, \mathcal{X}(t)) = \lim_{h \rightarrow 0^+} \inf \frac{1}{h} [V(t+h, \mathcal{X}(t+h)) - V(t, \mathcal{X}(t))]$.

Define $m(t) = \sum_{i=1}^n |x_i(t) - x_i^*(t)| + \sum_{j=n+1}^{n+m} |y_j(t) - y_j^*(t)|$ and consider a Lyapunov function

$$V((\mathcal{X}(t), \mathcal{X}^*(t))) = \sum_{i=1}^n \left| \ln \frac{x_i(t)}{x_i^*(t)} \right| + \sum_{j=n+1}^{n+m} \left| \ln \frac{y_j(t)}{y_j^*(t)} \right|. \tag{2}$$

Theorem 2. Let the following conditions hold:

- (1) The hypotheses (H_1) - (H_5) hold.
- (2) There exist non-negative continuous functions δ, δ_i such that

$$(H_9) \quad a_{ii}^l + \sum_{k=n+1}^{n+m} \frac{c_{ik}^l \beta_{ik}^l \tau_j^*}{(\alpha_{ik}^u + \beta_{ik}^u \theta_i + \gamma_{ik}^u \theta_j)^2} - \sum_{k=1, k \neq i}^n a_{ik}^u > \delta_i - \sum_{j=n+1}^{n+m} \frac{a_{ji}^u \beta_{ji}^u \theta_j^*}{(\alpha_{ji}^l + \beta_{ji}^l \tau_i + \gamma_{ji}^l \tau_j)^2},$$

$$(H_{10}) \quad \delta_{jj}^l - \sum_{k=n+1, k \neq j}^{n+m} \delta_{jk}^u - \sum_{k=1}^n \frac{d_{jk}^u \alpha_{jk}^u + d_{jk}^u \beta_{jk}^u \theta_i^*}{(\alpha_{jk}^l + \beta_{jk}^l \tau_i + \gamma_{jk}^l \tau_j)^2} > \delta_j(t) + \sum_{i=1}^n \frac{(c_{ij}^u \alpha_{ij}^u + c_{ij}^u \beta_{ij}^u \theta_i^*)}{(\alpha_{ij}^l + \beta_{ij}^l \tau_i + \gamma_{ij}^l \tau_j)^2}.$$

- (3) $0 < 1 + b_{ik} + h_{ik} \leq 1, 0 < 1 + b_{jk} + h_{jk} \leq 1$.

Then the solution $\mathcal{X}(t)$ of (1) is globally stable.

Proof. Consider the upper right derivative $D^+V(\mathcal{X}(t), \mathcal{X}^*(t))$ along the solution of system (III). For $t \geq t_0$ and $t \neq t_k, k \in \mathbb{N}^+$, we derive the estimate

$$\begin{aligned}
 & D^+V(\mathcal{X}(t), \mathcal{X}^*(t)) \\
 &= \sum_{i=1}^n \operatorname{sgn}(x_i(t) - x_i^*(t)) [-a_{ik}(t)(x(t) - x^*(t)) \\
 &\quad - \sum_{k=n+1}^{n+m} \left(\frac{c_{ik}(t)y_k(t)}{A_1} - \frac{c_{ik}(t)y_k^*(t)}{A_2} \right)] \\
 &\quad + \sum_{j=n+1}^{n+m} \operatorname{sgn}(y_j(t) - y_j^*(t)) \left[\sum_{k=n+1}^{n+m} -\delta_{jk}(t)(y_k - y_k^*) \right. \\
 &\quad \left. - \sum_{k=1}^n \left(\frac{d_{jk}(t)x_k(t)}{A_3} - \frac{d_{jk}(t)x_k^*(t)}{A_4} \right) \right] \\
 &\leq \sum_{i=1}^n \left[-a_{ii}^l + \sum_{k=1, k \neq i}^n a_{ik}^u - \sum_{k=n+1}^{n+m} \frac{c_{ik}^l \beta_{ik}^l \tau_j^*}{(\alpha_{ik}^u + \beta_{ik}^u \theta_i + \gamma_{ik}^u \theta_j)^2} \right. \\
 &\quad \left. + \sum_{j=n+1}^{n+m} \frac{d_{ji}^u \beta_{ji}^u \theta_j^*}{(\alpha_{ji}^l + \beta_{ji}^l \tau_i + \gamma_{ji}^l \tau_j)^2} \right] |x_i(t) - x_i^*(t)| \\
 &\quad + \sum_{j=n+1}^{n+m} \left[-\delta_{jj}^l + \sum_{k=n+1, k \neq j}^{n+m} \delta_{jk}^u \right. \\
 &\quad \left. + \sum_{k=1}^n \frac{d_{jk}^u \alpha_{jk}^u + d_{jk}^u \beta_{jk}^u \theta_i^*}{(\alpha_{jk}^l + \beta_{jk}^l \tau_i + \gamma_{jk}^l \tau_j)^2} + \sum_{i=1}^n \frac{(c_{ij}^u \alpha_{ij}^u + c_{ij}^u \beta_{ij}^u \theta_i^*)}{(\alpha_{ij}^l + \beta_{ij}^l \tau_i + \gamma_{ij}^l \tau_j)^2} \right] |y_j(t) - y_j^*(t)|.
 \end{aligned}$$

Where

$$\begin{aligned}
 A_1 &= \alpha_{ik}(t) + \beta_{ik}(t)x_i(t) + \gamma_{ik}(t)y_k(t); \\
 A_2 &= \alpha_{ik}(t) + \beta_{ik}(t)x_i^*(t) + \gamma_{ik}(t)y_k^*(t); \\
 A_3 &= \alpha_{jk}(t) + \beta_{jk}(t)x_k(t) + \gamma_{jk}(t)y_j(t); \\
 A_4 &= \alpha_{jk}(t) + \beta_{jk}(t)x_k^*(t) + \gamma_{jk}(t)y_j^*(t).
 \end{aligned}$$

Thus in view of hypothesis (H₉), we obtain

$$D^+V(\mathcal{X}(t), \mathcal{X}^*(t)) \leq -\delta(t)m(t), \tag{3}$$

for $t \geq t_0, t \neq t_k (t \in \mathbb{N}^+)$, where $\delta(t) = \min \delta_i, \delta_j$, For $t \geq t_0, t = t_k (t \in \mathbb{N}^+)$ we have

$$\begin{aligned}
 V(\mathcal{X}(t_k^+), \mathcal{X}^*(t_k^+)) &= \sum_{i=1}^n \left| \ln \frac{x_i(t_k^+)}{x_i^*(t_k^+)} \right| + \sum_{j=n+1}^{n+m} \left| \ln \frac{y_j(t_k^+)}{y_j^*(t_k^+)} \right| \\
 &= \sum_{i=1}^n \left| \ln \frac{(1 + b_{ik} + h_{ik})x_i(t_k)}{(1 + b_{ik} + h_{ik})x_i^*(t_k)} \right| + \sum_{j=n+1}^{n+m} \left| \ln \frac{(1 + b_{jk} + h_{jk})y_j(t_k)}{(1 + b_{jk} + h_{jk})y_j^*(t_k)} \right| \\
 &= V(\mathcal{X}(t_k), \mathcal{X}^*(t_k)).
 \end{aligned} \tag{4}$$

Then the inequality

$$V(\mathcal{X}(t_k), \mathcal{X}^*(t_k)) \leq V(\mathcal{X}(t_k^+), \mathcal{X}^*(t_k^+)) - \int_{t_0}^t \delta(s)m(s) dt, \quad t \geq t_0, \text{ holds.}$$

By the Mean Value Theorem and by Lemma 2 it follows that for any closed interval contained in $t \in (t_{k-1}, t_k], k \in \mathbb{N}^+$, there exist positive numbers r and R such that for every i, j $r \leq x_i(t), y_j(t), x_i^*(t), y_j^*(t) \leq R$ and

$$\begin{aligned} \frac{1}{R}|x_i(t) - x_i^*(t)| &\leq |\ln x_i(t) - \ln x_i^*(t)| \leq \frac{1}{r}|x_i(t) - x_i^*(t)|, \\ \frac{1}{R}|y_j(t) - y_j^*(t)| &\leq |\ln y_j(t) - \ln y_j^*(t)| \leq \frac{1}{r}|y_j(t) - y_j^*(t)|. \end{aligned} \tag{5}$$

Hence we obtain

$$\begin{aligned} V(\mathbf{x}_0, \mathbf{x}_0^*) &= \sum_{i=1}^n |\ln x_i(t_0^+) - \ln x_i^*(t_0^+)| + \sum_{j=n+1}^{n+m} |\ln y_j(t_0^+) - \ln y_j^*(t_0^+)| \\ &\leq \frac{1}{r} \|\mathbf{x}_0 - \mathbf{x}_0^*\|. \end{aligned} \tag{6}$$

Further, from (3) and (4) we have

$$D^+V(\mathcal{X}(t), \mathcal{X}^*(t)) \leq 0, \quad t \geq t_0, t \neq t_k,$$

and

$$\Delta V(\mathcal{X}(t), \mathcal{X}^*(t)) \leq V(t_0^+, \mathbf{x}_0, \mathbf{x}_0^*), \tag{7}$$

for all $t \geq t_0$. Given $0 < \varepsilon < R$, choose $\delta(t) = \frac{\varepsilon r}{2R}$. Then from (5)–(7) it follows that

$$\sum_{i=1}^n |x_i(t) - x_i^*(t)| + \sum_{j=n+1}^{n+m} |y_j(t) - y_j^*(t)| < \varepsilon,$$

for all $t \geq t_0$, whenever $\|\mathbf{x}_0 - \mathbf{x}_0^*\| \leq \delta$ and $t_0 \in \mathbb{R}^+$. Since $t_0 \in \mathbb{R}^+$ is arbitrary, by Definition 1(a), the system (1) is globally stable. This proves the theorem. \square

Theorem 3. *In addition to the assumptions of Theorem 2, suppose that there exists a constant such that*

$$\int_{t_0}^t \delta(s) ds = c(t - t_0). \tag{8}$$

Then the system (1) is globally exponentially stable.

Proof. We consider again the Lyapunov function (2), from (3) and (5) we obtain

$$D^+V(\mathcal{X}(t), \mathcal{X}^*(t)) \leq \delta(t)m(t) \leq \delta(t)rV(\mathcal{X}(t), \mathcal{X}^*(t)). \tag{9}$$

From the above estimate and (6), we have

$$V(\mathcal{X}(t), \mathcal{X}^*(t)) \leq V(\mathbf{x}_0, \mathbf{x}_0^*) \exp \left\{ -r \int_{t_0}^t \delta(s) ds \right\}, \quad (10)$$

for all $t \geq t_0$. Then, from (5), (6), (7) and (10) we deduce the inequality $\sum_{i=1}^n |x_i(t) - x_i^*(t)| + \sum_{j=n+1}^{n+m} |y_j(t) - y_j^*(t)| \leq \frac{R}{r} \|\mathbf{x}_0 - \mathbf{x}_0^*\| e^{-rc(t-t_0)}$, for $t \geq t_0$. This shows that the system (1) is globally exponentially stable. This proves the theorem. \square

Acknowledgments. This work is supported jointly by the Natural Sciences Foundation of China under Grant No. 60963025, Natural Sciences Foundation of Hainan Province under Grant No. 110007 and the Start-up fund of Hainan Normal University under Project No. 00203020201.

References

1. Beddington, J.R.: Mutual interference between parasites or predators and its effect on searching efficiency. *J. Animal Ecol.* 44, 331–340 (1975)
2. Kuang, Y., Baretta, E.: Global qualitative analysis of a ratio-dependent predator-prey system. *J. Math. Biol.* 36, 389–406 (1998)
3. Cosner, C., DeAngelis, D.L., Ault, J.S., Olson, D.B.: Effects of spatial grouping on the functional response of predators. *Theoret. Pop. Biol.* 56, 65–75 (1999)
4. Abrams, P.A., Ginzburg, L.R.: The nature of predation: Prey-predator. Ratio-dependent or neither. *Trends Ecol. Evol.* 15, 337–341 (2000)
5. Fan, M., Kuang, Y.: Dynamics of a nonautonomous predator-prey system with the Beddington-DeAngelis functional response. *J. Math. Anal. Appl.* 295, 15–39 (2004)
6. Rui, X.: Global etability and Hopf bifurcation of a predator-prey model with stage structure and delayed predator response. *J. Math.* 67, 1683–1693 (2012)
7. Xiaohu, W., Shuyong, L., Xu, D.: Globally exponential stability of periodic solutions for impulsive neutral-type neura networks with delays. *J. Math.* 64, 65–75 (2011)
8. Lantang, M., Gliu, X.: Positive periodic soution for ratio-dependent n -specieses discrete time system. *J. Math.* 56, 577–589 (2011)
9. Gaines, R.E., Mawhin, J.L.: *Coincidence Degree and Nonlinear Differential Equations.* Springer, Berlin (1977)
10. Gaines, R.E., Mawhin, J.L.: *Coincidence Degree and Nonlinear Differential Equations.* Springer, Berlin (1977)
11. Zhang, J., Gui, Z.J.: Existence and stability of periodic solutions of high-order Hopfield neural networks with impulses and delays. *Journal of Computational and Applied Mathematics* 224, 602–613 (2009)
12. Yan, Y., Wang, K.H., Gui, Z.J.: Periodic Solution of Impulsive Predator-Prey Models with the Beddington-DeAngelis Functional Response. In: 5th International Congress on Mathematical Biology, pp. 86–91. World Academic Press (2011)
13. Anokhin, A., Berezansky, L., Braverman, E.: Exponential stability of linear delay impulsive differential equations. *J. Math. Anal. Appl.* 193, 923–941 (1995)
14. Samoilenko, A.M., Perestyuk, N.A.: *Impulsive Differential Equations.* World Scientific Series on Nonlinear Sciences. Ser. A, Singapore (1995)
15. Lin, Z., Liu, J., Pedersen, M.: Periodicity and blowup in a two-species cooperating model. *Nonlinear Analysis. Real World Applications* 12, 479–486 (2011)

Research on the Distal Supervised Learning Model of Speech Inversion

Ying Chen and Shaobai Zhang

Computer Department, Nanjing University of Posts and Telecommunications,
210003 Nanjing, Jiangsu, China
adzsb@163.com

Abstract. To the problem that articulatory information is not readily available in typical speaker-listener situations, a method that estimates such information from the acoustic signal was proposed, namely speech inversion. Distal supervised learning (DSL) was selected as one of machine learning strategies for speech inversion to study. Eight tract variables were used as articulatory information to model speech dynamics, and the experiment's background and theoretical foundation of distal supervised learning also were analyzed. Besides a global optimization approach was proposed and the results when speech signal is parameterized as acoustic parameters (APs) were compared with as mel-frequency cepstral coefficients (MFCCs). The results showed that distal supervised learning has a good estimation performance for tract variables.

Keywords: articulatory information, speech-inversion, distal supervised learning, tract variables.

1 Introduction

At present, the performance of the advanced automatic speech recognition (ASR) systems is affected in casual or spontaneous speech. This problem is mainly due to the great variability of spontaneous speech, and the variability is mainly caused by contextual variation (i.e. coarticulation). Many different studies have claimed that articulatory information can be used to improve the performance of automatic speech recognition systems. Unfortunately, such articulatory information is not so easy to get in typical speaker-listener situations. Therefore, we need to use a method to estimate such information, and this method is usually termed "speech-inversion".

The process of generating the acoustic speech signal by the organs in the human vocal tract can be represented by a function f like $f: t \rightarrow x$, where x is a vector that represents the acoustic speech signal, t is a vector representing the configuration of the articulators, and f is the function that defines the forward mapping from the articulatory domain to the acoustic domain. Thus, given a vector t_a , representing a specific articulatory configuration, we can obtain a specific speech output x_a , given f is known. In recognition tasks, we know acoustic speech signal but not the articulatory data. If we define a function g such that $g: x \rightarrow t$, then the articulatory

configuration t_b can be obtained from the speech signal sample x_b using the function g . Thus g is the inverse of function f and the formula represents the task of acoustic to articulatory speech inversion.

Several machine learning techniques have been implemented for the task of speech inversion, such as artificial neural network (ANN), support vector regression (SVR) and auto regressive artificial neural network (AR-ANN). ANN has versatility in nonlinear[1] regression problems but it falls short in ill-posed regression problems where the ill-posedness is due to one-to-many mapping; the estimated articulatory trajectories from SVR model are found to be corrupted by estimation noise; the feedback loop in the AR-ANN architecture helps to maintain the inherent smoothness of the articulatory trajectories but at the same time can be a source of progressive error introduction and it has a high computational cost. Supervised learning with distal teacher or distal supervised learning (DSL) not only can address the one-to-many mapping problem, but also has a better estimation performance than ANN, SVR and AR-ANN. The following section of this article will focus on DSL.

2 Tract Variables (TVs)

This article uses eight tract variables (TVs)[2-3] as articulatory information to model speech dynamics. Tract variables describe the constriction degree and location of distinct organs along the vocal tract (as shown in Fig.1). Each TV involves its own set of associated articulators and they are lip aperture (LA), lip protrusion (LP), tongue tip constriction degree (TTCD), tongue tip constriction location (TTCL), tongue body constriction degree (TBCD), tongue body constriction location (TBCL), velum (VEL) and glottis (GLO).

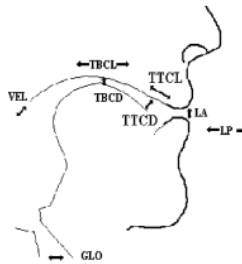


Fig. 1. Tract variables (TVs) from different constriction locations

Tract variables[4-6] describe geometric features of the shape of the vocal tract tube in terms of constriction degree and location; an active gesture is specified by activation on set and offset times and parameter values for a set of critically damped, second-order differential equations[7], shown in formula (1), where M , B , and K are mass, damping coefficient, and stiffness parameters of each TV (represented by z) and z_0 is the target position of the gesture:

$$M\ddot{z} + B\dot{z} + K(z - z_0) = 0 \quad (1)$$

There are three advantages for using TVs. First, the TVs specify the salient features of the vocal tract area functions directly. Second, the constrictions in TV space are controlled by articulatory gestures directly and they embody the speaker's phonological goals. Also TVs can be expected to bear a relation to speech acoustics that is closer to one-to-one than does the complete area function, and help to reduce the non-uniqueness of speech inversion. So it is the TV value that is informative in terms of phonological category and lexical access. Finally, incorporating TV information which estimated from the acoustic signal not only improves the performance of gesture recognition but also can help in improving noise-robustness[8] of ASR systems.

3 Dataset and Signal Parameterization

The database used in the research comes from the literature [9]. The database is generated by task dynamic and applications model along with HLSyn, and it contains synthetic speech along with their articulatory specifications. The synthetic database was generated by inputting the text for the 420 unique words. The output synthetic speech[4] was sampled at 10 kHz and the TV time functions and gestural scores were sampled at 200 Hz; seventy-five percent of the data were used for training, ten percent for validation, and the rest for testing.

Speech signal is parameterized as acoustic parameters (APs)[10-11] and mel-frequency cepstral coefficients (MFCCs). According to the relevance of acoustic parameters, 40 different APs were selected, while for the MFCCs, 13 cepstral coefficients were extracted. The APs were measured using a 10-ms window with a frame rate of 5 ms, and MFCCs' each acoustic features was measured at a frame rate of 5 ms (time-synchronized with the TVs) with window duration of 10 ms. The acoustic features and the target articulatory information (i.e. the TVs) were z-normalized and then scaled such that their dynamic range is confined within [-0.95,+0.95]. According to the existing observation we can know that, incorporating dynamic information helps to reduce the non-uniqueness problem for the speech inversion task. Hence, the input features are contextualized in the experiments reported in this article. The feature contextualization is defined by the context-window parameter \hat{C} , where the current frame (with feature dimension d) is concatenated with \hat{C} frames from before and after the current frame (with a frame shift of 2 or time shift of 10 ms), generating a concatenated feature vector of size $(2\hat{C}+1)d$. Some prior researches[12] have identified that the optimal context parameter \hat{C} for the MFCCs is 8 (context duration of 170 ms) and for the APs is 9 (context duration of 190 ms) and such values are used in the experiments presented in the rest of the article.

4 Distal Supervised Learning (DSL)

To address the issues with conventional supervised learning architectures for one-to-many mapping cases, Jordan et al.[13], proposed supervised learning with a distal teacher or DSL. In the DSL[4] paradigm, there are two models placed in cascade with one another: 1) the forward model which generates acoustic features given the articulatory trajectories, hence M-to-1 mapping and 2) inverse model which generates the articulatory trajectories from acoustic features, hence 1-to-M mapping. Given a set of $[x_b, y_b]$ pairs, DSL first learns the forward model, which is unique but not necessarily perfect. DSL learns the inverse model by placing it in cascade with the forward model as shown in Fig.2. The DSL[4] architecture can be interpreted as an “analysis-by-synthesis” approach, where the forward model is the synthesis stage and the inverse model is the analysis stage; in the DSL approach, the inverse model is trained (its weights and biases updated) using the error that is backpropagated through the forward model whose previously learned weights and biases are kept constant.

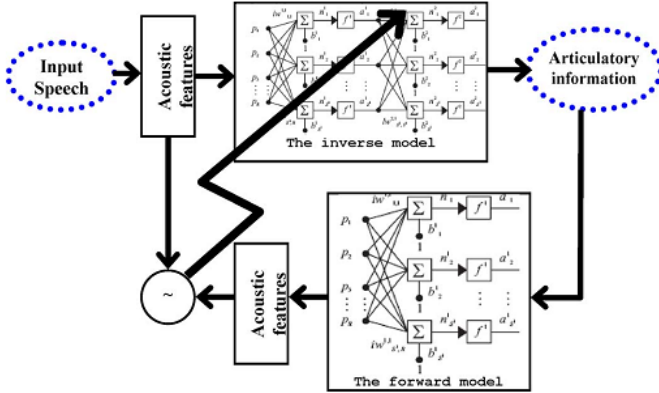


Fig. 2. Distal supervised learning approach for obtaining acoustic to TV mapping

Considering a forward mapping between an input vector x and an output vector y , using a vector of network weights and biases, w and b , the relationship can be expressed as:

$$\hat{t} = g(x, w, b) \tag{2}$$

Learning the forward model is based on the following cost function[10]:

$$L = \frac{1}{2} E[(t - \hat{t})^T (t - \hat{t})] \tag{3}$$

where t is the desired target for a given input. For the inverse model, [13] defined two different approaches, a local optimization approach and an optimization along the trajectory approach. The local optimization approach[4] necessitates using an online

learning rule, whereas the optimization along trajectory requires recurrency in the network (hence, error minimization using backpropagation in time), both of which significantly increase the training time and memory requirements. In this article, a global optimization approach is used which is similar with local optimization approach, and it uses the tools of DSL as proposed in [13], but instead uses batch training in the feedforward network. The cost function that the DSL tries to minimize is represented as:

$$J = \frac{1}{2N} \sum_{k=1}^N [(t_k^* - t_k)^T (t_k^* - t_k)] \quad (4)$$

where N is the total number of training samples, t_k is the target vector for the k th training sample and t_k^* is the actual target output from the network. The weight update rule is as follows:

$$w[n+1] = w[n] - \eta \nabla_w J_n \quad (5)$$

where η is the learning rate, $w[n]$ represents the weights of the network at time index n . The gradient can be obtained from formula (4) using the chain rule:

$$\nabla_w J_n = \frac{1}{N} \sum_{k=1}^N \left(- \frac{\partial x_k^T}{\partial w} \frac{\partial t_{k,n}^*}{\partial x_k} (t_k - t_{k,n}^*) \right) \quad (6)$$

where $t_{k,n}^*$ is the estimated target vector for the k th training sample at the n th time instant.

5 Experiment, Results and Discussion

The DSL architecture was trained for all the eight TV trajectories for each acoustic feature. The forward models were created using single hidden-layer feedforward artificial neural networks and trained using SCG algorithm, and the number of neurons in the hidden layer was optimized using the root mean squared error over the validation set. The inverse models were built using a 3-hidden-layer network and the number of neurons in each layer was optimized using the root mean squared error on the validation set. The DSL models were trained using gradient descent learning algorithm (with a variable learning rate), momentum learning rule (momentum=0.9) and mean squared predicted performance error with regularization as the optimization criteria (regularization parameter=0.4). The number of neurons in the forward model was 350 and 400 and in the inverse model respectively were 150-100-150 and 250-300-250 for MFCC and AP.

What the experiment needs to demonstrate is that given a speech signal, tract variables can be estimated with a high accuracy. Two quantitative measures were used in the experiment to compare the shape and dynamics of the estimated articulatory trajectories with the actual ones. The two measures are the root mean squared error (RMSE) and the Pearson product-moment correlation (PPMC) coefficient. The RMSE[4] gives the overall

difference between the actual and the estimated articulatory trajectories, whereas the PPMC gives a measure of amplitude and dynamic similarity between them; the RMSE and the PPMC are defined as follows:

$$RMSE = \sqrt{\frac{1}{N}(e-t)^T(e-t)} \quad (7)$$

$$\gamma PPMC = \frac{N \sum_{i=1}^N e_i t_i - [\sum_{i=1}^N e_i][\sum_{i=1}^N t_i]}{\sqrt{N \sum_{i=1}^N e_i^2 - (\sum_{i=1}^N e_i)^2} \sqrt{N \sum_{i=1}^N t_i^2 - (\sum_{i=1}^N t_i)^2}} \quad (8)$$

where e represents the estimated TV vector and t represents the actual TV vector having N data points, and the N for TVs is 8. Some of the TVs have a different measuring unit (e.g., TBCL and TTCL are measured in degrees), thus to better summarize the inversion performance for all articulatory trajectories, we use the non-dimensional mean normalized RMSE, RMSE_{nm} [5] and its average, RMSE_{nm_avg} defined by:

$$RMSE_{nm,i} = \frac{RMSE_i}{\sigma_i} \quad (9)$$

$$RMSE_{nm_avg} = \frac{1}{N} \sum_{i=1}^N RMSE_{nm,i} \quad (10)$$

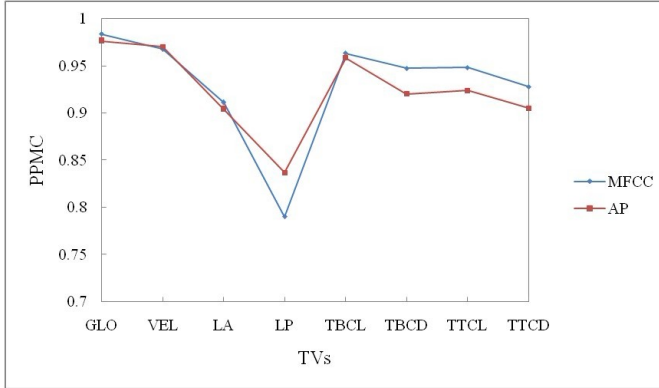


Fig. 3. PPMC for TV estimation using MFCC and AP

The TV estimation results from DSL for both APs and MFCCs are shown in Fig.3 and Fig. 4.

It can be observed from the plots of Fig. 3 that comparing the PPMC of six TVs (except the velum and the lip protrusion) when speech signal is parameterized as

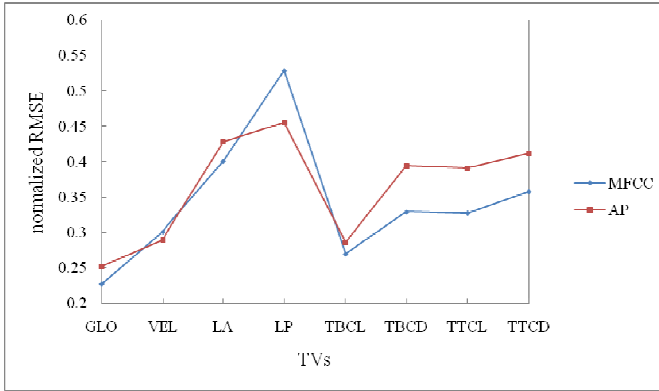


Fig. 4. Normalized RMSE for TV estimation using MFCC and AP

MFCCs and APs, the value of former is higher. Also, it can be observed from the plots of Fig.4 that comparing the RMSE of six TVs (except the velum and the lip protrusion) when speech signal is parameterized as MFCCs and APs, the value of former is lower. Note that lower RMSE and higher PPMC indicate better performance of the estimation. So the plots of Fig.3 and Fig.4 shows that the APs overall offered better accuracy for velum and the lip protrusion, whereas for the other TVs, the MFCCs provided better results. It also can be observed from Fig.3 and Fig.4 that the estimation performance for glottis is best while for lip protrusion is worst among all TVs from DSL. Since almost all the TVs have high PPMC and low RMSE, DSL has a good estimation performance.

6 Conclusion

In Cartesian coordinates, many different sets of articulatory location may represent the same vocal tract constriction, but the tract variables specification is unique. McGowan pointed out that the non-uniqueness problem with speech inversion is ameliorated by the use of TVs, hence, for TVs we can expect a further reduction in non-uniqueness for the speech inversion task. Besides as pointed out before, DSL can address the problem of ANN in ill-posed regression issue and has a better estimation performance for TVs than SVR and AR-ANN. Therefore DSL plays an important role in speech inversion. As pointed out before[4], the DSL topology is more like an analysis-by-synthesis architecture, where the performance of synthesis part entirely depends upon the accuracy of the forward model; and to ensure a highly accurate forward model, exhaustive data is typically required to ensure the forward model has examples of all possible pairs of articulatory data and acoustic observation. It is the shortcoming of DSL, but it has little influence on DSL's application value.

Acknowledgment. This work is supported by the National Natural Science Foundation of China (No. 61073115).

References

1. Neiberg, D., Ananthakrishnan, G., Engwall, O.: The Acoustic to Articulation Mapping: Non-linear or Non-unique. In: Proc. Interspeech, pp. 1485–1488 (2008)
2. Zhuang, X., Nam, H., Hasegawa-Johnson, M., Goldstein, L., Saltzman, E.: The Entropy of Articulatory Phonological Code: Recognizing Gestures from Tract Variables. In: Proc. Interspeech, pp. 1489–1492 (2008)
3. Zhuang, X., Nam, H., Hasegawa-Johnson, M., Goldstein, L., Saltzman, E.: Articulatory Phonological Code for Word Classification. In: Proc. Interspeech, pp. 2763–2766 (2009)
4. Mitra, V., Nam, H., Espy-Wilson, C.Y., Saltzman, E., Goldstein, L.: Retrieving Tract Variables from Acoustics: a Comparison of Different Machine Learning Strategies. *IEEE Journal of Selected Topics in Signal Processing* 4, 1027–1045 (2010)
5. Katsamanis, A., Papandreou, G., Maragos, P.: Face Active Appearance Modeling and Speech Acoustic Information to Recover Articulation. *IEEE Trans. Audio, Speech, Lang. Process.* 17(3), 411–422 (2009)
6. Mitra, V., Özbek, I., Nam, H., Zhou, X., Espy-Wilson, C.: From Acoustics to Vocal Tract Time Functions. In: Proc. ICASSP, pp. 4497–4500 (2009)
7. Byrd, D., Saltzman, E.: The Elastic Phrase: Modeling the Dynamics of Boundary-Adjacent Lengthening. *J. Phonetics* 31(2), 149–180 (2003)
8. Mitra, V., Nam, H., Espy-Wilson, C., Saltzman, E., Goldstein, L.: Noise Robustness of Tract Variables and their Application to Speech Recognition. In: Proc. Interspeech, U.K., pp. 2759–2762 (2009)
9. Nam, H., Goldstein, L., Saltzman, E., Byrd, D.: Tada: An Enhanced, Portable Task Dynamics Model in Matlab. *J. Acoust. Soc. Amer.* 115(5-2), 2430 (2004)
10. Juneja, A.: Speech Recognition Based on Phonetic Features and Acoustic Landmarks. Ph. D. dissertation, Univ. of MD, College Park (2004)
11. He, X., Deng, L.: Discriminative Learning for Speech Processing. In: Juang, G.H. (ed.). Morgan & Claypool, San Mateo (2008)
12. Mitra, V., Nam, H., Espy-Wilson, C.: A Step in the Realization of a Speech Recognition System Based on Gestural Phonology and Landmarks. In: Proc. 157th Meeting ASA, Portland, vol. 125, p. 2530 (2009); *J. Acoust. Soc. Amer.*
13. Jordan, M.I., Rumelhart, D.E.: Forward Models–Supervised Learning with a Distal Teacher. *Cogn. Sci.* 16, 307–354 (1992)

Low Power Pulse Width Modulation Design for Class D Audio Amplifier Systems

Ruei-Chang Chen¹, Shih-Fong Lee¹, and Yeong-Chau Kuo²

¹ Department of Electrical Engineering, Da-Yeh University, Changhua 51591, Taiwan

² Department of Electronic Engineering, National Kaohsiung First University of Science
and Technology, Kaohsiung City, 81164, Taiwan

b9204007@yahoo.com.tw, sflee@mail.dyu.edu.tw,
yckuo@nckust.edu.tw

Abstract. This paper presents the design and implementation of a novel pulse width modulation (PWM) chip. With low-power, high-performance, small area, and high speed, these circuits are employed in portable computer systems, such as the power circuits, electronic circuits, video and music amplifiers circuits, communications and control circuits, wireless communication and high-frequency circuit systems. This PWM chip followed the chip implementation center advanced design flow, and then was fabricated using Taiwan Semiconductor Manufacture Company 0.35- μm 2P4M mixed-signal CMOS process. The chip supply voltage is 3.3 V which can operate at a maximum frequency of 100 MHz. The total power consumption is 3.0268 mW, and the chip area size is 1.016 mm \times 1.016 mm. Finally, the PWM chip was tested and the experimental results are discussed. From the excellent performance of the chip verified that it can be applied to audio amplifiers, communications control, etc.

Keywords: pulse width modulation, Communications technology, systems applications, computer chip technology.

1 Introduction

Since PWM circuit [1], [2] was successfully fabricated with integrated circuits (ICs), widespread attention has been attracted. PWM system feature high power efficiency, and have been successfully applied to various audio/video products. However, most of the studies up to date have focused on half bridge converters presumably because they are easier to implement.

Recently, the advance of semiconductor fabrication techniques furthers the need of various control and amplification systems for audio/video electronic products has been roared. In particular, the fabrication techniques for complementary metal-oxide-semiconductor (CMOS) field effect transistors used in this design has been developed quite maturely, suitable for new generation ICs. In addition, CMOS technology has the advantages of high efficiency, fast turn out, low power consumption, and low cost over other technologies. This paper describes the design and implementation of a low power PWM for Class D [3], [4] audio amplifier systems.

2 Differential Control PWM Chip Design

The proposed block diagram is shown in Fig. 1. PWM chip use differential control design. Because of PWM have the advantage of high power efficiency and less power consumption. Oscillator [5], [6] can produce a square waveform that resembles a clock which drives a subsequent voltage ramp generator [7]. Voltage ramp generator output a triangular reference signal V_+ at a much higher frequency carrier. The input stage of the op-amp [8], [9] is a differential amplifier. The original audio signal V_{in} can be op-amp Enhanced. The audio signal V_- amplified by an op-amp is the input to a comparator [10], [11] which compares it with a triangular reference signal V_+ . The comparator is a switching amplifier. The output of the comparator gives the PWM signal. This signal can be amplified and modulation to higher frequency. An LC filter is a novel full-bridge converter 2nd-order low-pass filter employed between the PWM output stage and the load to attenuate the high frequency carrier, thereby recovering the audio signal.

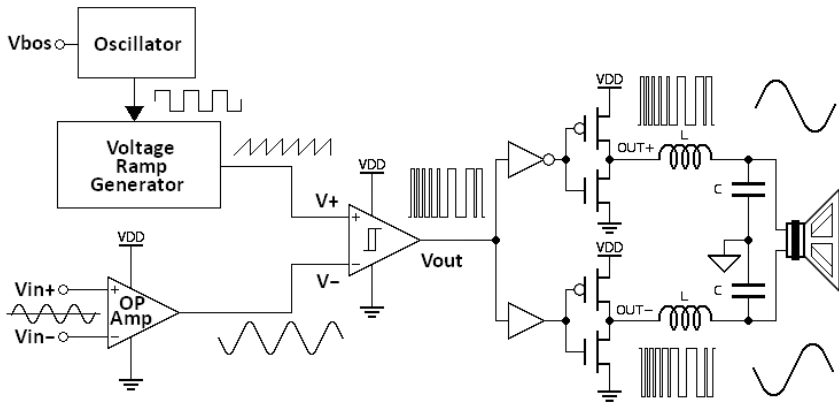


Fig. 1. The proposed block diagram of PWM system

2.1 Operational Amplifier, Op-amp

The two-stage op-amp schematic is shown in Fig. 2(a). The input stage of the op-amp is a differential amplifier. This type of circuit amplifies the difference between two input voltages, V_{in+} and V_{in-} , and produces an output voltage V_- which is proportional to the difference between these inputs. The audio signal is applied to the inverting input of the first stage op-amp. The amplified audio signal is fed into the inverting input of the second stage op-amp. Note that the input to M1 is identified as the inverting input V_{in-} , and the input to M2 is identified as the non-inverting input V_{in+} .

The metal-oxide-semiconductor field effect transistor, or MOSFET, is a voltage controlled device, and input impedance is infinitely great, so the gate current I_G is zero. In the op-amp circuit, the inputs V_{in+} and V_{in-} are both dc biased at 1.65 V. We

apply an ac audio signal of 1mV to the inverting input V_{in-} . When V_{in-} goes more positive than V_{in+} , the relative input polarities do not match the polarity signs, and the output V_o goes negative. When V_{in-} goes more negative than V_{in+} , the output V_o goes positive, the 180° phase shift between the op-amp input and output signals.

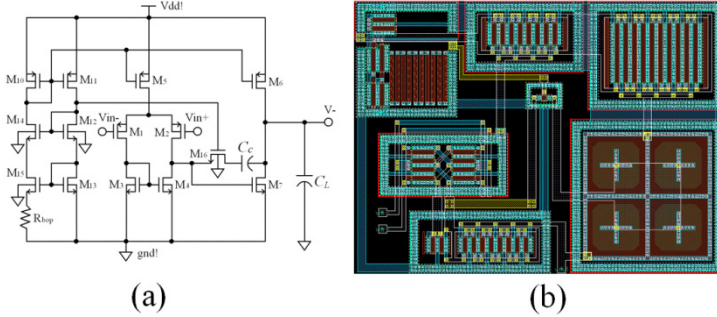


Fig. 2. The two-stage op-amp circuit. (a) Schematic. (b) Layout.

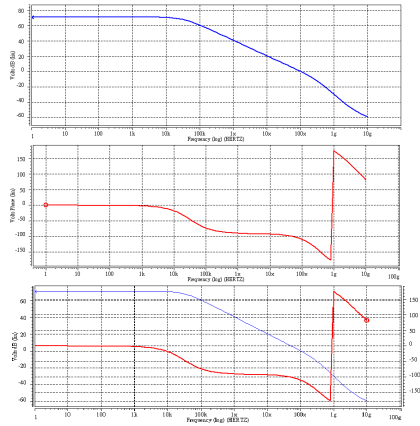


Fig. 3. Simulation results for the proposed two-stage op-amp

The layout for the two-stage op-amp circuit is shown in Fig. 2(b). All devices or circuits prone to produce electromagnetic interference or susceptible to interference are enclosed with double layer guard rings. Since a small tilt angle is used in the ion implantation process to prevent channeling effect, some devices may be slightly asymmetric depending on their location on the wafer. In addition, the line width of long interconnects may not be the same due to non-uniform etching. These are all likely to induce failure. Dummy devices are included to avoid these problems. In particular, the differential pair composed by op-amps M1 and M2 is interlaced lined up to obtain perfect symmetry and matching. In addition, source contacts and drain

contacts were closely laid in order to achieve higher current driving capacity. In addition, currents in all MOS transistors should flow in the same direction to prevent cross interference among transistors. At last, guard ring is placed along the outer peripheral of the layout to protect the chip from interference among systems.

We carried out HSPICE simulation for the circuit illustrated in Fig. 2(a) and obtained the curves shown in Fig. 3. The voltage gain is shown as the blue curve of Fig. 3. About 72 dB is gained for the voltage of the operation amplifier with a bandwidth of about 102 MHz. The phase response is shown as the red curve in the middle has a phase angle of 65°. The complete response is also shown for comparison. In this design, the gain margin of op-amp is determined first, and the phase margin is determined next. The purpose is to prevent instability caused by the rapid change of poles and zeros. In fact, this problem of instability can be easily overcome for phase response by increasing the bias voltage which can increase phase margin. From the simulation results shown in Fig.3, it is fully verified that the proposed design can satisfactorily meet the set requirements. The detail specification is summarized in Table 1.

Table 1. Simulated performance of the operational amplifier

Parameters	Values
Voltage Gain (A_v)	$A_v \geq 70$ dB
Gain Bandwidth (GB)	$GB \geq 100$ MHz
Slew Rate (SR)	$SR \geq 10$ V/ μ s
Phase Margin (PM)	$PM \geq 60^\circ$
Load Capacitor (CL)	$CL \leq 2$ pF
Power Supply Voltage	3.3 V
Output Swing	0~3.3 V
Power Dissipation	1.6012 m Watts
Total Area	93.60 μ m \times 68.75 μ m

2.2 Voltage Controlled Ring Oscillator, VCO

The proposed five-stage voltage controlled ring oscillator circuit is shown as the CMOS inverter in Fig. 4(a). It consists of two complementary MOS transistors, i.e., an NMOS and a PMOS. In the digital circuits design for the five-stage inverter circuit, the size of PMOS is twice that of NMOS. Its purpose is to make the N type and P type transistor the same as the working tempo. The oscillator circuit layout is shown in Fig. 4(b). The spontaneous oscillation frequency of ring oscillator circuit can be controlled with the V_{bos} voltage which supplies a voltage in the range of 0.55~3.3 V. As oscillation frequency is increased, oscillator can produce a square waveform that resembles the simulation results shown in Fig. 5(a). By changing the voltage V_{bos} , the oscillator circuit can output waveforms of different frequencies. Therefore, the purpose for the ring oscillator is to generate a clock which drives a subsequent voltage ramp generator circuit.

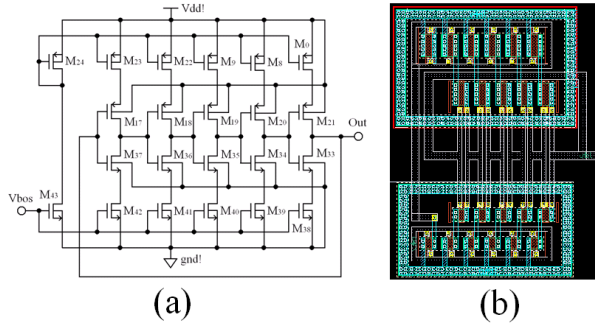


Fig. 4. The voltage controlled ring oscillator circuit. (a) Schematic. (b) Layout

2.3 Voltage Ramp Generator

The voltage ramp generator is shown in Fig. 5(b). The layout for the voltage ramp generator circuit is shown in Fig. 6(a). Inside the voltage ramp generator, signal amplification is accomplished by using inverting amplifiers. The output signal of preceding stage is fed into the input of succeeding stage whose channel width is doubled, and current capacity is amplified by 2. Four consecutive inverting amplifiers make up the non-inverting amplifying system which charges or discharges the capacitor in the succeeding stage. Simulation result for voltage ramp generator output voltage waveform is shown in Fig. 6(b). Ramp wave at the desired frequency can be achieved by adjusting the frequency of oscillator or the value of capacitance. Since the frequency range audible to human is around 20~20 KHz, the frequency for modulation carrier should be at least 10 times of the maximum audible frequency. In fact, the higher the frequency of modulation carrier, the finer the recovered signal. Fig. 7(a) shows the measured output voltage waveforms of the voltage ramp generator. The triangular ramp wave generated by the voltage ramp generator is used as a carrier signal and is fed into the non-inverting input of the comparator in the succeeding stage.

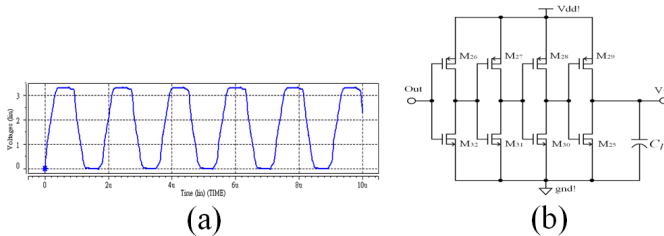


Fig. 5. (a) Simulation results for the VCO output voltage waveforms. (b) Schematic of the voltage ramp generator circuit.

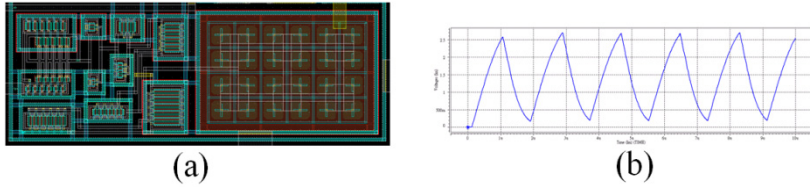


Fig. 6. (a) Layout for the voltage ramp generator circuit. (b) Simulation result for the output voltage waveform of voltage ramp generator.

2.4 Hysteresis Comparator

As shown in Fig. 7(b), the comparator circuit is a switching amplifier. The audio signal amplified by an op-amp is the input to a comparator, which compares it with a voltage ramp wave. The result is a pulse-width-modulated square wave whose period is equal to that of the voltage ramp, and its pulse width represents a sample of the audio signal. The voltage ramp frequency is set at a frequency that is very much higher than that of the audio signal. If the audio signal applied to the inverting input (V^-) is larger than the voltage ramp signal applied to the non-inverting input (V^+), the output of comparator is 0 (Low). On the contrary, if the audio signal is smaller than the voltage ramp signal, the output of comparator is V_{DD} (High).

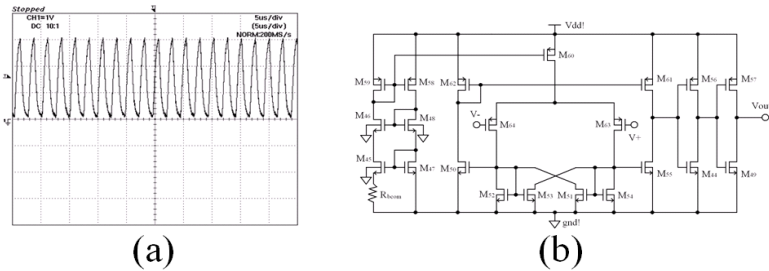


Fig. 7. (a) Measured output voltage waveforms of voltage ramp generator. (b) Schematic of the comparator circuit.

The resulting simulation results for the waveforms of comparator are shown in Fig. 8(a). The first curve (green) is the voltage ramp signal applied to the non-inverting input of comparator. The second curve (blue) is the audio signal applied to the inverting input of comparator. And the third curve (red) represents the output signal of comparator (V_{out}). These three curves are put together and is shown as the fourth curve in Fig. 13. Based on these simulation results, a layout for comparator was successfully designed and is shown in Fig. 8(b).

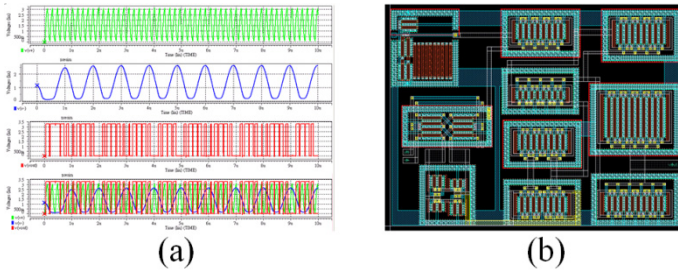


Fig. 8. (a) Simulation results at 1 MHz for the comparator. (b) Layout for the comparator.

3 Experimental Results

Fig. 9(a) shows the PWM chip layout. The PWM chip microphotograph is shown in Fig. 9(b). This PWM chip followed the chip implementation center advanced design flow, and then was fabricated using Taiwan Semiconductor Manufacture Company 0.35 μm 2P4M mixed-signal CMOS process. The chip supply voltage is 3.3 V which can operate at a maximum frequency of 100 MHz. The total power consumption is 3.0268 mW, and the chip area size is, including pads, 1.016 \times 1.016 mm². The inputs Vin+ and Vin- are both dc biased at 1.65 V. We apply an ac audio signal of 1mV to the inverting input Vin-. The simulation results for the waveforms of PWM output are shown in Fig. 10. Measured PWM output voltage waveforms are shown in Fig. 11(a).

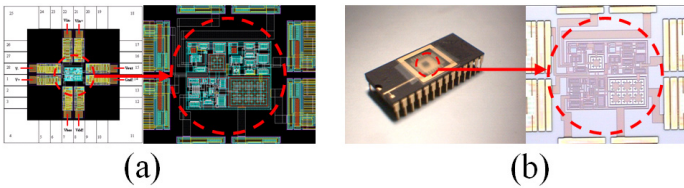


Fig. 9. (a) Layout of the PWM circuit. (b) Microphotograph of the PWM chip.

Fig. 11(b) shown the PWM output voltage waveforms, all meet we design of requirements, correctly of produced narrow different of square signals, circuit performance full show meet we all design of requirements.

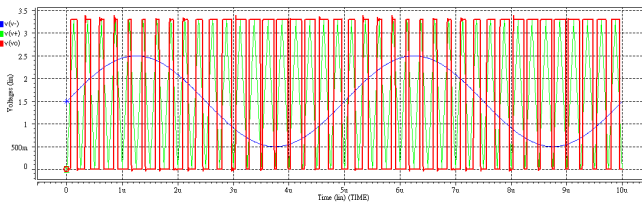


Fig. 10. Simulation results for PWM output voltage waveform

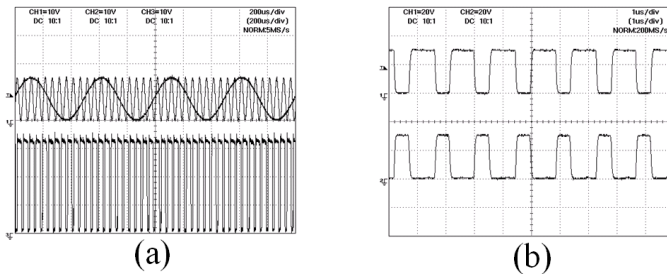


Fig. 11. (a) Measured PWM output voltage waveforms. (b) Experimental results of the PWM output voltage waveforms.

4 Conclusion

In conclusion, a novel PWM technique for Class D audio systems has been successfully designed and implemented. Circuit design, simulation, analysis, layout design, chip implementation, and measurement are all included in this study. The performance of PWM amplifier is enhanced in this study. As verified from theoretical analysis, simulation results, and measurement results, the chip possess super characteristics suitable for various applications such as audio amplification and modulation control.

References

1. Hava, A.M., Un, E.: A high-performance PWM algorithm for common-mode voltage reduction in three-phase voltage source inverters. *IEEE Trans. Power Electronics* 26(7), 1998–2008 (2011)
2. Patella, B.J., Prodic, A.: High-frequency digital PWM controller IC for DC–DC converters. *IEEE Trans. Power Electronics* 18(1), 438–446 (2003)
3. Rojas-Gonzalez, M.A., Sanchez-Sinencio, E.: Low-power high-efficiency Class D audio power amplifiers. *IEEE J. Solid-State Circuits* 44(12), 3272–3284 (2009)
4. Berkhout, M., Dooper, L.: Class-D audio amplifiers in mobile applications. *IEEE Trans. Circuits Syst. I* 57(5), 992–1001 (2010)
5. Leung, B.: A switching-based phase noise model for CMOS ring oscillators based on multiple thresholds crossing. *IEEE Trans. Circuits Syst. I* 57(11), 2858–2869 (2010)
6. Hajimiri, A., Limotyrakis, S., Lee, T.H.: Jitter and phase noise in ring oscillators. *IEEE J. Solid-State Circuits* 36(6), 790–804 (1999)
7. Li, M., Amaya, R.E.: Design of mM-W fully integrated CMOS standing-wave VCOs using low-loss CPW resonators. *IEEE Trans. Circuits Syst. II* 59(2), 78–82 (2012)
8. Mahattanakul, J.: Design procedure for two-stage CMOS operational amplifiers employing current buffer. *IEEE Trans. Circuits Syst. II* 52(11), 766–770 (2005)
9. Kurkure, G., Dutta, A.K.: A novel adaptive biasing scheme for CMOS Op-Amps. *J. Semiconductor Technology and Science* 5(3), 168–172 (2005)
10. Goll, B., Zimmermann, H.: A comparator with reduced delay time in 65-nm CMOS for supply voltages down to 0.65 V. *IEEE Trans. Circuits Syst. II* 56(11), 810–814 (2009)
11. Fiorenza, J.K., Sepke, T., Holloway, P., Sodini, C.G., Lee, H.-S.: Comparator-based switched-capacitor circuits for cabled CMOS technologies. *IEEE J. Solid-State Circuits* 41(12), 2658–2668 (2006)

A New Logic Method for Education Resource Software Guarantee

Guan Wei¹ and Lv Yuanhai²

¹ Information Center, Xi'an University of Posts and Telecommunications,
Xi'an 710121, China

² Information Center, Xi'an University of Posts and Telecommunications,
Xi'an 710121, China

guanw@xupt.edu.cn, lyh@xupt.edu.cn

Abstract. The research attempts to understand the reasons that affect education resource software guarantee. Furthermore, this paper attempts to understand the cultural influences on education resource software and the interaction between the social reason and other reasons that affect education resource software guarantee. Based on deep analysis of current literature on software guarantee reasons, a logic method is formulating. The reasons included in the method are divided into four categories: management reasons, software quality reasons, people reasons, and social reasons. Finally in this paper, a theoretical method is given.

Keywords: Education resource software, Guarantee, Software quality reasons, Social Reasons, Management Reasons.

1 Introduction

According to Jones [1], many organizations indicated that a number of their education resource software failed; and between one and two thirds of education resource software exceed their budget and time. Further, the expert argued that about half of the expensive education resource software at the end will be considered out of control and cancelled.

In a survey of 292,000 application projects in large, medium and small cross industry companies, the Standish group showed that about 25% of projects were cancelled before completion and 51% exceeded their budgets and time scales and had fewer features and functions than originally specified [3]. Another survey by Taylor [4] showed that out of 1064 education resource software studied, only 120 projects (12%) were successful.

The Standish group [3] distinguished three types of systems: successful, challenged and impaired. A successful system means that “the system is completed on time and on budget, with all features and functions as initially specified”, a challenged system means, “a system completed and operational but over-budget, over the time estimate, and offer fewer features and function than originally specified”. An impaired system means “a system that is cancelled at some point during the development cycle”.

A project is usually deemed as successful if meets requirements is delivered on time and delivered within budget [4]. Therefore software risk management is an approach that attempts to form risk oriented correlates of development guarantee into a readily applicable set of principles and practice [5]. More over the implementation of education resource software is a complex task involving the successful alignment of both the technical and social system within an organization [6]. Furthermore after decades of research, systems development and implementation projects remain notoriously hard to mange and many continue to end in failure [7].

The main objective of this research is to investigate the reasons that affects on education resource software guarantee of education resource software in Jordanian firms.

2 Related Work

This section provides an overview on related work on education resource software guarantee reasons. The classification provided below is based on deep analysis of the literature and the authors' opinions. The education resource software guarantee reasons are divided into four categories:

2.1 Management Reasons

The guarantee of education resource software is often related to management reasons [7]. Management reasons are divided into four reasons: the presence of formal methodology, clear business objective, executive support and minimized project scope.

2.2 Software Quality Reasons

The guarantee of education resource software is often related to software quality reasons [4]. Software quality reasons are divided into three reasons: the presence of standard software infrastructure, understanding requirements and managing requirements changes and reliable estimates.

2.3 People Reasons

The guarantee of education resource software is often related to people reasons [3]. People reasons are divided into two reasons: user involvement and the presence of an experienced project manager.

2.4 Social Reason

Denison [7] mentioned that the social reason is a critical guarantee reason in education resource software. None of these authors included this reason in their method or models.

3 Research Problem

Many researchers have become interested in researching the reasons that affect on education resource software guarantee. Executive support affects the process and progress of project and lack of it put the project at bad situation, [5]. If project lack user involvement it fail even if it developed on time and budget, project fails if it does not meet user needs or expectation, professionals of project concern and care on this part, that lead to fail to achieve project objectives. Past literatures showed that experienced education resource software manager can identify risk on project and lack of senior manager commitment seen as most critical risk on project [6]) Boehm and Ross pointed out that the present of two different team in software development with two different objective, one deal with user requirement and another deal with technical challenges this lead to misunderstanding of objectives for the project.

For any project time is enemy for it, since scope impact time, if we minimize it we can impact it within time so the chance to guarantee is increase [6]. In contrast to requirements which are in changeable state, infrastructure needs stability. Standish group found that about 75% of application code is infrastructure and by using standard infrastructure, the application team concentrate on business rule rather than technology standard infrastructure can shortcut applications integration that many developmental team fail to application it. When we create base level of requirement to our project and then develop those features we can reduce the requirement changes, help user and sponsors to see the result faster, and add benefits for project managers to prepare and link the need and criteria for the next phase of our project. The use of good formal methodology provides realistic pictures about the project and some step may be reusable so tendency to reinvent the wheel minimized and stability of project increased, also formal methodology gives manager the ability to estimate the real time so the risk is reduced. In order to develop project you need to make good and realistic estimation which is hard planning and through it you purchase the requirement and component of projects, managers must use their collective knowledge and experience to get good estimate that reflect the real effort needed. Budget and cost estimation is a crucial reason in education resource software guarantee. These reasons include small milestones, proper planning, competent staff and ownership and communication skills [6].

4 Content Analysis

Based on deep analysis of the above literature, the following are our findings:

- 1) Executive Support is an important guarantee reason for education resource software.
- 2) Clear Business Objective is an important guarantee reason for education resource software.
- 3) Formal Methodology is an important guarantee reason for education resource software.

4) Minimizing project Scope is an important guarantee reason for education resource software.

5) Standard software Infrastructure is an important guarantee reason for education resource software.

6) Understanding Requirements and Managing Requirements Changes are important guarantee reasons for education resource software.

7) Reliable Estimates is an important guarantee reason for education resource software.

8) User Involvement is an important guarantee reason for education resource software.

9) Experienced Project Manager is an important guarantee reason for education resource software.

10) Social is an important guarantee reason for education resource software as follow:

A. Social influences the behavior of all individuals within an organization, including how decisions are made, who makes them, how rewards are given, who is promoted, how people are treated, and how the organization responds to its environment.

B. To change the social of an organization, people need to be aware of what drives the thinking, feeling, and behavior of the organization.

C. Social provides stability and predictability as it gives direction for behavior, ideas, and how to respond in various situations.

D. Information Technology manager must not only manage operations, finances, and implementation for education resource software but also the social.

E. A positive relationship was found between the social reasons and other reasons that affect education resource software guarantee.

5 Research Methodologies

An empirical study as a combination of questionnaire survey and interview was applied in this research. Only 25 managers were interviewed because the others excused because of they were busy or in traveling.

6 Research Model

Research model in Figure 1 is built based on the combination of several past literatures. Based on these literatures the research has formed the following hypothesis: 1-There are positive relationships between (standard requirements, user involvement, executive support, clear business objective, minimized scope, reliable estimation, formal methodology, standard infrastructure, manager experience and other reasons) reasons And education resource software guarantee.

7 Samples

The most of the sample are males (n = 170) which consist (78.6%) of the sample where the females portion consists (25.1%) of the sample. For education variable, bachelor degree took the high portion (71.5%) whereas high certificates portion was (28.5%). For experience variable, the high portion went to (12-16 years) which consisted (19.6%); the lowest portion went to (6 years and less), (15.6%).

8 Results on Reasons That Contributed to Education Resource Software Guarantee

Based on the research model Figure 1, and the results that we have got from the surveys, we summarize all reasons that affect on education resource software guarantee according to severity degree and frequency degree.

The results show that means of influential reasons in guarantee of software engineering projects (scope, executive support, clear object, manager experience, social, reliable estimation, standard requirement, user involvement, formal methodology, and firm infrastructure) came at high degree and the total mean of influential reasons in guarantee of software engineering projects in terms of its severity (3.72), SD (0.52); Standard requirements dimension ranked first degree with mean reached(3.96) followed by user involvement (3.86), executive support at third degree (3.65), clear objective at fourth degree (3.62), while the dimension of other reasons ranked tenth rank with mean (3.45) at middle degree.

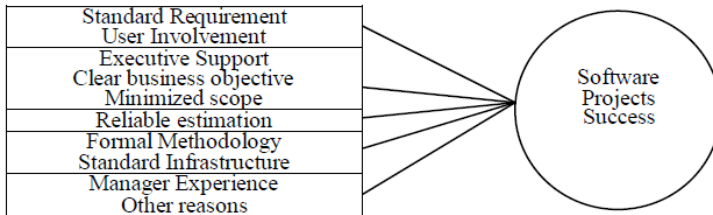


Fig. 1. Research Model

Also the results show that means of influential reasons in guarantee of software engineering projects (scope, executive support, clear object, manager experience, social, reliable estimation, standard requirement, user involvement, formal methodology, and firm infrastructure) came at high degree and the total mean of influential reasons in guarantee of software engineering projects in terms of its frequency (3.58), SD (0.56); social dimension ranked first degree with mean reached (3.82) followed by standard requirement (3.85), user involvement at third degree (3.56), executive support at fourth degree (3.61), while the dimension of manager experience ranked tenth rank with mean (3.36) at middle degree.

9 Result Discussions

9.1 Standard Requirement Reason

This reason came at first rank in terms of severity degree and frequency. From the researcher's view, weakness of workers in collecting requirements and disability to use the right methods in collection process is from the reasons that lead to make these requirements unclear. Each project has its privacy in requirements collection process by the user and analyzing these requirements; also if these requirements don't identify accurately, this may cause failure to the project; documentation mechanism considered as an important requirement that is through revising available documents in Jordanian institutions there was no documentation mechanism for requirements in order to revise it with systems users when finishing these systems at the end. User's misunderstanding of needed requirements is one of the reasons that lead to imperfection of requirements and changing them during working on the system. Identifying project method, clarifying and defining objective, and user involvement play an important role in defining the project requirements.

9.2 User Involvement Reason

This reason ranked the second rank in terms of severity and frequency degree to assure the interest in continuation of teamwork members when they doing their works perfectly without any imperfection may affect negatively quantities and qualitative outcomes. For the importance of the big role that user involvement plays in guarantee of software engineering projects, this dimension is very important to clarifying work objectives and making balance among teamwork members when they sharing their roles and enable them to know unclear things in the institution, so, involving the user within teamwork helps in resolving some of vagueness in specific items of the work which is to be achieved by software engineering team in the institution. User involvement process minimizing resistance change for new system that is the user who intended to be involved within teamwork should have a positive role to protect the project and tries to convince the others with the importance of the project for the institution in all; also this user should have a role in communication flexibility process between system development team and system users because he/she is the only one who can explain the teamwork and his /her co- workers views.

9.3 Executive Support Reason

In terms of severity and frequency degree in guarantee of software engineering projects, this dimension ranked the third rank. Executive administration will minimize difficulties that facing system development team. Support will be provided to the project if the expected benefits are big and if the project can be applied in more than one location in the institution. Software engineering manager should have the ability to attract executive administration in order to provision support and distributing the available resources to the project's phases.

10 Logic Method for Education Resource Software Guarantee

Based on deep analysis of the literature and our findings, we suggest the following logic method for software project guarantee (see Figure 2). We believe that Education resource software is a lengthy undertaking involving a set of complex activities that take a lot of time and cost. The guarantee in such an undertaking depends on good and reliable estimation. Based on the above analysis, we classify education resource software guarantee reasons into four categories: management reasons which are the presence of formal methodology, clear business objective, executive support and project scope, software quality reasons that are standard software infrastructure, understanding requirements and managing requirements changes and reliable estimates, people reasons that are user involvement and experienced project manager, and the cultural reasons which involve management social. Also in the analysis, a relationship was found between the management reasons, people reasons and software quality reasons.

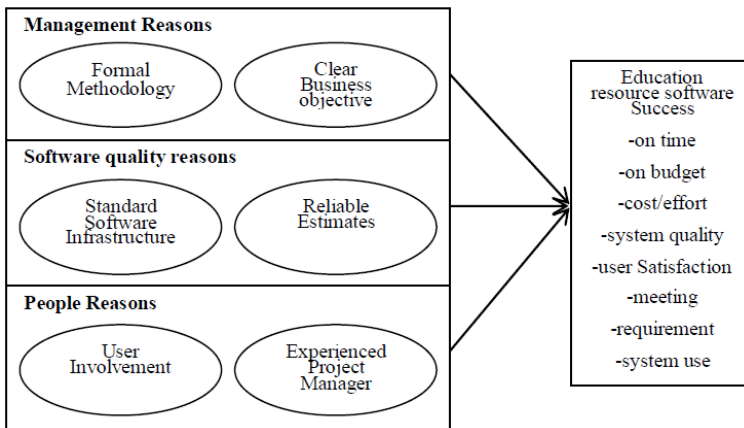


Fig. 2. The new method

11 Conclusions

The research attempts to understand the reasons that affect education resource software guarantee. Furthermore, this paper attempts to understand the cultural influences on education resource software and the interaction between the social reason and other reasons that affect education resource software guarantee. In addition, this paper provides the interested audience, either in private or academic sectors, with a short description of those reasons that result in education resource software guarantee. Finally in this paper, a theoretical method is given.

The deep analysis of the literature showed that project guarantee is dependent on many reasons including executive support, clear objectives, presence of a formal methodology, minimization of project scope, use of a standard software infrastructure,

understanding and managing requirements, making reliable estimates, user involvement, presence of experienced project manager, and last not least, taking cultural aspects into consideration. The work presented above resulted in the formation of a logic method for education resource software guarantee. This method is novel in the sense that it includes many reasons that are not found together in any other method. The frame-work stresses the importance of cultural reasons in project guarantee. The reasons included in the method are validated through content analysis of the literature. The new method is currently being calibrated by using it to evaluate any ongoing and completed education resource software. The results will be presented in a forthcoming paper.

References

1. Sarma, W.S., Rao, V.: A Rough-Fuzzy Method for Integration off Candidate Data unit for Software Reuse. *Pattern Recognition Letters* 24(6), 875–886 (2003), doi:10.1016/S0167-8655(02)00199-X
2. González-Calero, P.A.: Applying Knowledge Modeling and Case-Based Reasoning to Data unit integration. *IEE Proceedings-Software* 147(5), 169–177 (2000)
3. Lucrédio, D., et al.: Integration of information, Using Metric Indexing. In: *IEEE International Conference on Information Reuse and Integration*, Las Vegas, November 8-10, pp. 79–84 (2004)
4. Lucrédio, D., et al.: A Survey on Software Data unit Search and Integration. In: *30th IEEE Euromicro Conference*, Rennes, August 31-September 3, pp. 152–159 (2004)
5. Salton, G., Wong, A., Yang, C.S.: A Data feature model for Automatic Indexing. *Communications of the ACM* 18(11), 613–620 (1975), doi:10.1145/361219.361220
6. Sorumgard, L.S., Sindre, G., Stokke, F.: Experiences from Application of a Faceted Classification Scheme. In: *Advances in Data Unit Integration, Selected Papers from the 2nd International Workshop on Software Reusability*, Lucca, March 24-26, pp. 116–124 (1993)
7. Lovins, J.B.: Development of a Stemming Algorithm. *Mechanical Translation and Computational Linguistics* 11(1-2), 22–31 (1968)

Research on 3D Object Rounding Photography Systems and Technology

Zhenjie Hou¹, Junsheng Huang¹, and Jianhua Zhang²

¹ School of Information Science and Engineering, Changzhou University, Changzhou, China

² College of Computer Science and Engineering, Inner Mongolia Agricultural University, Hohhot, China

hhderek@163.com

Abstract. The 3D Rounding Photography system is a application of Desktop Virtual Reality. It takes the pictures around the objects with the computer, and then deal with them for the 360 degree display in the local computer or a web page. The 3D visual sense and image real-time amplification of this system is the key technical problem. Under researching on the system, it puts forward a new solution based on the Silverlight technology. At the same time, it researches the image interpolation application, and proposes an improved triangle interpolation based on the determination of the pixels edge s. The experimental results show that the solution meet the visual sense and interaction Requirements of the 3D rounding photography system, and better to maintain the image edge in the image of partial enlargement process, has some practical value.

Keywords: 3D Rounding Photography, SilverLight, Image Interpolation.

1 Introduction

Desktop virtual reality technology use a conventional computer monitor to display the virtual world, also be called the Window of the World. Compared to the "immersive" virtual reality technology, desktop virtual reality system has the advantages of low costing and easy to promote [1]. On the basis of the widely used desktop virtual reality, the rapid developments of the panoramic view technology become new and popular visual technologies. Suturing or processing the picture, it can realize the tasks of looking around landscape or dragging objects in the three-dimensional space freely [2]. It is called 3D Object rounding photography systems through accessing to the object image of 360 degree and processing the points or zooming the images.

In the early Object rounding photography production, it become the obstacles to large-scale promote because of the shortage of photography precision, production efficiency and application method. With the development of rich internet applications, 3D Object rounding photography, a low-cost, high efficiency, simple three-dimensional imaging modalities are used in large quantities. Rich internet applications take full advantage of the hardware capabilities of the client to improve their abilities, enhance interactivity and presentation layer logic to improve the user experience is to

fill the usability gap between local applications and internet applications [3]. Web development trends and directions have served as a carrier of 3D Object rounding photography system-wide dissemination and application. 3D Object rounding photography system as a business application development achieve its main approach based on actionscript scripting programming and simulated three-dimensional interactive flash applications. There are also other implementations, such as the dynamic GIF images, QuickTime rounding photography file. In recent years Silverlight technology as the main products of Microsoft Corporation in the rich internet applications develop rapidly and fully meet the requirements of 3D Object rounding photography system development.

Compared to the other implementations, silverlight technology not only developed 3D Object rounding photography system in maintaining small file size and compositing operation simply, but also has the ability in processing control and system scalability.

This paper presents the improved triangular linear interpolation algorithm during studying the 3D Object rounding photography systems for local image zoom functions. The experimental results show the algorithm can keep the edge sharpness of the image more clearly under the premise of maintaining the linear interpolation algorithm speed.

2 Silverlight-Based 3D Object Rounding Photography Systems

The basic materials for 3D Object rounding photography system are group photos of the object in 360 degrees. Spliced into a panoramic image or dynamic picture and display on a computer. Visitors with the mouse or keyboard control can control the direction and distance to meet the interactive desktop virtual reality technology requirements. In the rapid popularization of the internet, it is the preferred carrier for 3D Object rounding photography system.

Silverlight technology is a new generation of rich internet application technology, and is cross-browser, cross-platform implementation of the Net Framework, has extremely superior vector graphics, animation, multimedia and rich network communication function. Compared with traditional WPF, Java Swing, Delphi, silverlight has more lightweight runtime environment and better background language framework supports. Therefore, using silverlight technology to develop 3D Object rounding photography system is a completely new system solution and has great practical value.

2.1 System Analysis

3D Object rounding photography system is an image processing and display system, including image processing, image storage, and web presence. The design of the system will display in the internet browser, so the current internet connection speed must be considered, a small size and file format are the primary needs of the entire system. Secondly, a group of pictures obtained through professional equipment need to go through the process of image processing, the synthesis of specific file format

need real-time storage, which requires to create an image database. Rendering process in the later steps need obtain image data from the database and can be interactive played in the play-side. The web presence is part of the vision and interactive that the 3D Object rounding photography system users experience and need the clear playback control logic. After analyzing requirements of the above system, there are many solutions can be used. For example, any language supported by Windows system can realize the image processing, just meeting the requirement of file format and communicating with the database. Similarly, it has considerable freedom to choose the database. For example, if choosing direct storage of image data, we can use the SQL Server, the Oracle, etc. The rich internet applications are used mainly in the final step: rendering. Rich Internet applications mainly include Java FX, JavaScript / Ajax, the Microsoft ActiveX, Silverlight, Flash, etc.

2.2 System Design

In order to achieve 3D image rendering capabilities , we use silverlight to develop rich internet applications in 3D Object rounding photography system. The system developed the local image processing clients and establishing a database based on SQL server for data storage and data communication, the file format of the image information is swf file.

The system operation and transmission of image data based on SQL Server database have functions of real-time image storage and interaction between web client and database. Adapt Linq to control database, so the local database is simple and convenient to operate. This system has a very strong process control capability and potential for extensions. Because the Silverlight client uses a text-based XAML format, enhanced search system applications can be found.

2.3 Interactive Features

The interactive features of the 3D Object rounding photography system refers to the vision and interactive of 3D object on the client. How to implement the functions that user can see the objects left or right, near or far, zoom freely is the client application logic must be addressed. Based on the interactive features of the Silverlight technology, 3D Object rounding photography system is a client application that uses the C# language to meet the system logic requirements on the .Net frame. Our system achieves the interaction logic in the playback side different from the current widespread use of embedding logic directly into the Swf file. These methods we used will greatly reduce the amount of network data transmission. The experimental data indicate that a nearly 900k Swf file with logical data can be reduced to below 500k after extracting required picture materials our system needed. It is also prove directly that the system has the superiority to reduce the network burden..

3D Object rounding photography system we implemented shown in Fig. 1.

After testing, the memory usage shown in Table 1. The test computer information is: CPU, Intel dual-core E4500, Memory: 2G, Operating System: Microsoft Windows XP Professional, SP3.

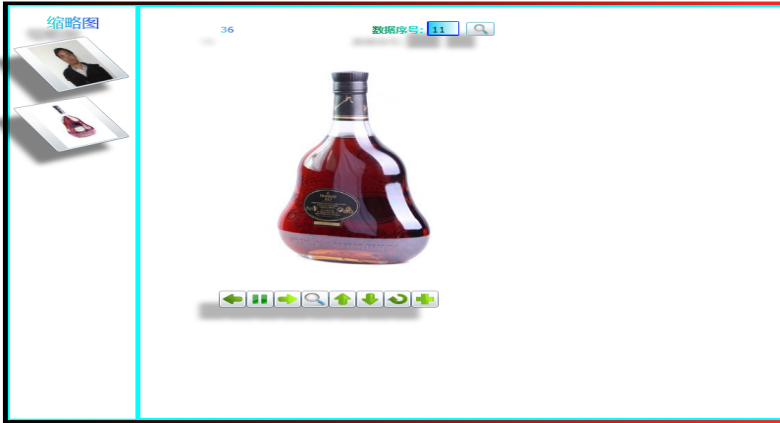


Fig. 1. The System of Running

Table 1. Instance of the Test Memory Usage

Swf file size	497K			
System the number of running instances	1	2	3	4
IE process memory footprint (K)	106448	174648	237580	305212

From the data in table 1 can be seen that the system occupied the memory of the IE process is not high. Especially the personal computer memory is generally more than 2G now, the system is provided with practical value.

3 Researches on Real Time Amplification of Images in 3D Object Rounding Photography Systems

Images in real time zoom feature are another concern for the design of the system. Image scaling technologies that the most widely used are nearest neighbor interpolation and bilinear interpolation taking into account the computational speed of the image interpolation. However, due to the linear interpolation algorithm is a low-pass filter, the enlarged image will have details degradation and loss some high frequency information, while the high-frequency information is contained in the edge of the image [4]. Therefore, some of the non-linear interpolation algorithm successively been proposed to overcome the inherent defects of the linear interpolation algorithm to preserve the image edge, for example the bilateral filtering interpolation, the local covariance interpolation, wavelet interpolation method, the neural network image interpolation and partial differential equations-based image interpolation method [5]. Although these methods are able to obtain a clearer edge, due to the algorithms of high complexity and computing time and other reasons ,it is difficult to be used in the actual project. Therefore, taking into account the 3D Object rounding photography system, this paper will adapt improved triangular linear

interpolation algorithm to maintain the computation speed advantage and preserve the image edge information better.

The triangular linear interpolation algorithm will judge the edge of four pixels that will be interpolated linearly to eliminate the point of the maximum pixel value compared with the other three pixels and interpolate bilinearly in the triangle. Compared with bilinear interpolation, the ability of triangular linear interpolation of saving the image edge and clarity has improved.

3.1 Improved Triangular Linear Interpolation Algorithm

In the adjacent four pixels, it is the focus of the triangular linear interpolation algorithm that how to determine the pixel edge direction in the fastest time [7-11]. In order to meet the requirements on the computing speed, the paper put forward the rapid determination method proposed by Dan Su, Philip Willis. Set pixels on the corner are a, c and b, d; just compare the size of the angle on the absolute value of the pixel values will be able to quickly distinguish the pixel edge. if $|a-c| > |b-d|$, it can be decided an edge where b, d on the diagonal [6]. That is to say that the quadrilateral pixel structure the original bilinear interpolation required is divided into two triangles, then take the triangular linear interpolation. Due to the linear interpolation using only three pixel values but can be achieved interpolation in the X direction and Y direction, it can be think as the degradation form of a bilinear interpolation. However, such marginal determine has an advantage in speed, when a diagonal pixel values close to another and the difference is very large between two pairs of diagonal pixels values, for example, a and c pixel value of 10, 11, and b and d pixel value of 178,180, we can not guarantee this determine correct. This paper proposed a new improvement determine on the basis of ensuring the interpolation speed.

For the above possibilities of justice miscarriage, this paper proposes the edge overlap determine method. If the grid as the original image to be to determine the four adjacent pixels neighborhood, we will determine the four large grid composed of 1245,2356,4578,5689, and two directions of marking with 0 and 1 will be stored. If the edge of the slope is positive, then stored as 1, otherwise store 0. Let f(5) is pending to determine the direction, if $f(1245) = f(5689) = 0$, then $f(5) = 0$, otherwise if $f(2356) = f(5678) = 1$, then $f(5) = 1$, if the above conditions are neither satisfied, then f(5) will directly compare the difference of four pixels on the angle.

3.2 Time Complexity Analysis

In this section will analyze the time complexity of expansion algorithm. Set the original image height m width m, the linear interpolation of the promoter region is divided into the number of $(m-1) * (m-1) * 2$. Firstly, determine four squares and diagonal connection need eight subtraction and six comparisons, and if they failed to yield conclusive results then take the two subtractions and one comparison. Set $n = (m-2) * (m-2)$. The improved triangular linear interpolation determining the edge required up to $14 * n + 3n = 17n$, so the computing speed still can be guaranteed for the complexity of proposed expansion algorithm is $O(n)$ same as the linear interpolation algorithm.

3.3 Interpolation Comparative Analysis and Results

In order to compare the improved algorithm and the original triangular interpolation algorithm, this paper adopts the lena image to test algorithm. Interpolation results in Figure 2.



Fig. 2. Lena Image Interpolation Effect (2 Times). (a) Original image. (b)Result through the original triangle interpolation algorithm. (c) Result through the improved algorithm.

From Figure 2 can be seen that the proposed algorithm improve the image clarity and edge transition, especially when the image interpolation magnification is largely. At the same time, the peak signal to noise ratio improved better than the original interpolation algorithm. Through the above analysis and experiments, the improved algorithm can be used instead of ordinary triangular linear interpolation algorithm in 3D Object rounding photography system.

4 Conclusions

3D Objects rounding photography system as an efficient 3D imaging modalities, widely used in the industry of e-commerce, virtual tour, virtual museum etc. With the further development of system technology and rich internet applications, 3D Object rounding photography system will surely achieve more spectacular results. This paper presents a design scheme based on Silverlight technology with good process control capabilities and potential for expansion. At the same time improve the triangular linear interpolation algorithm to meet the real-time system requirements and improve the interpolation results.

References

1. Fu, X., Guo, B.: Overview of image interpolation technology. *Computer Engineering & Design* 30, 141–142 (2009)
2. Tomasi, C., Manduchi, R.: Bilateral Filtering for Gray and Color Images, pp. 839–846. IEEE, Piscataway (1998)

3. Li, X., Orchard, M.T.: New Edge-Directed Interpolation, vol. 10, pp. 1521–1527. Institute of Electrical and Electronics Engineers Inc. (2001)
4. Derado, G., Bowman, F.D., Patel, R., Newell, M., Vidakovic, B.: Wavelet Image Interpolation (WII): A Wavelet-Based Approach to Enhancement of Digital Mammography Images, pp. 203–214. Georgia Institute of Technology (2007)
5. Ulo, L.: Numerical solution of evolution equations by the haar wavelet method. *Applied Mathematics and Computation* 185, 695–704 (2007)
6. Marsi, N., Carrato, S.: Neural network-based image segmentation for image interpolation, pp. 388–397. IEEE, Piscataway (1995)
7. Casciola, G., Montefusco, L.B., Morigi, S.: Edge-driven Image Interpolation using Adaptive Anisotropic Radial Basis Functions, vol. 36, pp. 125–126. Springer, Netherlands (2010)
8. Zhou, J., Xue, Z., Wan, S.: Survey of Triangulation Methods. *Computer and Modernization*, 75 (July 2010)
9. Dyn, N., Levin, D., Rippa, S.: Data Dependent Triangulations for Piecewise Linear Interpolation. *IMA J. Numerical Analysis* 10, 137–154 (1990)
10. Yu, X., Morse, B.S., Sederberg, T.W.: Image Reconstruction Using Data-Dependent Triangulation, vol. 21, pp. 62–63. Institute of Electrical and Electronics Engineers Computer Society (May 2001)
11. Su, D., Willis, P.: Image Interpolation by Pixel Level Data-Dependent Triangulation, vol. 23, pp. 190–191. Blackwell Publishing Ltd. (June 2004)

Facial Expression Feature Selection Based on Rough Set

Dong Li^{1,2}, Yantao Tian¹, Chuan Wan¹, and ShuaiShi Liu¹

¹ School of Communication Engineering, Jilin University, Changchun, 130025, China

² Jilin Institute of Metrology, Changchun, 130022, China

Lidong0726@126.com

Abstract. An improved reducing algorithm for rough set attributes has invented for answering the question of the excessive features vector dimensions. It obtains the local feature vector through geometric feature points. By introducing the rough set and improved reducing algorithm that it is able to select optimally among the existing expression features, also clipping the redundancy and useless information for the selection of expression feature. The experiment has showed that, this method has demonstrated high level of validity for its more convenience, higher recognition rate and more efficiency.

Keywords: Rough Set, Feature Selection, Facial Expression Recognition.

1 Introduction

Emotion recognition is an integral part of quantitative studies of human behavior. Emotionally-cognizant human-computer and human-robot interfaces promise a more responsive and adaptive user experience. Many applications can benefit from an accurate emotion recognizer. In many real world problems, reducing dimension is an essential step before any analysis of the data can be performed.

In pattern recognition and general classification problems, methods such as Principal Component Analysis (PCA), Independent Component Analysis (ICA) and Fisher Linear Discriminate Analysis (LDA) have been extensively used [1-7]. These methods find a mapping from the original feature space to a lower dimensional feature space.

The recognition of facial expression is the process of extraction, optimization and categorization of the facial expression information. It enables the computer to know the facial expression so as to predict the peoples' mental status. The interaction between human and computer is realized in this way. The mainstream method of obtaining the expression feature is the geometric feature points, because it can work without the consideration of skin color, partial block and lights. However, this method often leads to 2 extremes: feature extracted too rare to fully express valid facial information and category; feature extracted has so many redundancies that it caused repetition and waste, sometimes the main useful category is buried among the piled trash information which costs so many store spaces and processing time ----so called "the disaster of feature vector". The key is to find a balance tackling the clash caused by the high and low feature vector, the ideal method is to reduce the redundancy also

keep the decisive useful data. Branch –and –bound search algorithm, forward sequence, backward sequence and maximum, minimum choice method etc., they provide a feasible proposal[8-12].

Variable selection procedures have been used in different settings. Among them, the regression area has been investigated extensively. A multi layer perceptron is used for variable selection. Stepwise discriminant analysis for variable selection is used as inputs to a neural network that performs pattern recognition of circuitry faults. Other regression techniques for variable selection are described. In contrast to the regression methods, which lack unified optimality criteria, the optimality properties of PCA have attracted research on PCA based variable selection methods [13-18]. As will be shown, these methods have the disadvantage of either being too computationally expensive, or choosing a subset of features with redundant information. It is hard to analysis, discover and reason the relations between the data, whereas the target of rough set theory is to link knowledge with category. The introduction of rough set theory into the feature choice and the improved reducing algorithm make up this shortcoming, obtain the easiest feature set, and at the same time optimize the system recognition. It elevates the efficiency by applying the result to the facial expression recognition.

2 Rough Set

The rough set theory is a new mathematical tool for analysis and processing obscure, inaccurate, disagree and incomplete information and knowledge, offering effective technology for artificial processing information.

The rough set theory related is stated below:

Definition 1: A knowledge expression set T is 5 tuples $T = \langle U, C, D, V, f \rangle$, U as the object set, $C \cup D = R$ as the attribute set, finite nonempty set, C and D as the subset representing condition attributes and decision attributes, $D \neq \Phi$, $V = \cup Va, a \in R$, Va is the range of attribute a , f is information function, designating every attribute value x include in U pointed.

Definition 2: Supposing R as equivalence relation family, $r \in R$ if $IND(R) = IND(R - \{r\})$, r is the reduced knowledge in R ; if $P = R - \{r\}$ is independent, P is a reduction in R . If $Q \subseteq P$, supposing Q is independent, and $IND(Q) = IND(P)$, Q is a reduction of P , as $Q \in RED(P)$, $RED(P)$ representing all the reduction set of P , all the necessary knowledge set is called the core of P , as $CORE(P)$, $CORE(P) = RED(P)$.

Definition 3: Supposing $S = (U, R, V, f)$ as decision table $C \cup D = R, C \cap D = \Phi$, C as condition attributes set, D as decision attributes set, if $U/C = \{x_1, x_2, \dots, x_n\}$, $U/D = \{y_1, y_2, \dots, y_n\}$.

The support of C in D can be defined as:

$$K(D) = \frac{1}{|U|} \sum_{i=1}^n |CY_i| = \frac{1}{|U|} \sum_{i=1}^n |POS_C(Y_i)| \quad Y_i \in U/D$$

The $||$ representing the number of element geometry included, the support of decision attribute broadly measures categorization ability of the decision table, which is also called categorization quality.

3 Geometry Feature Selection Method Based on Rough Set

The online version of the volume will be available in LNCS Online. Members of institutes subscribing to the Lecture Notes in Computer Science series have access to all the pdfs of all the online publications. Non-subscribers can only read as far as the abstracts. If they try to go beyond this point, they are automatically asked, whether they would like to order the pdf, and are given instructions as to how to do so.

3.1 Feature Extraction

The effective extraction of expression feature is the key of expression recognition technology, which directly decides the result of expression recognition. Great process has been made in expression extraction method based on preset knowledge, profile, exterior and related information. How to improve to select the optimal feature which represents mostly facial expression is very important question to be answered.

As the main constituting organs in human face, browns, eyes, noses, mouth etc. builds the rich and varied facial expression. The various changes of its profile, size and relative position lead to multiple difference of facial expression, the geometry description of the changes in its size and structure can be treated as an important feature for facial expression. Using the salient feature dot to obtain a group of the recognition feature is a simple and easy method. Based on this extraction standard, , as shown in Fig.1.

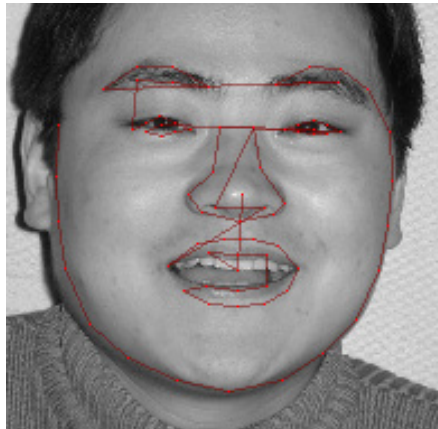


Fig. 1. Facial features points location

The description of the distance between the points demonstrates the changes of browns, eyes, nose and mouth's structures and shapes. Taken the commonness into the consideration, it cut the individual variance and retains the commons. There is no change in the distance between the inner eye comers; d -- the distance between the inner eye corner is the feature normalization factor. After the normalization, the Euclidean distance between i and j can be defined as:

$$Dis(i, j) = \frac{\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}}{d} \quad (1)$$

After the normalization, the vertical distance $Hei(i, j)$ from i to j can be defined as:

$$Hei(i, j) = \frac{|y_i - y_j|}{d} \quad (2)$$

After the normalization, the horizontal distance $Wid(i, j)$ can be defined as:

$$Wid(i, j) = \frac{|x_i - x_j|}{d} \quad (3)$$

x_i and x_j is the x-coordinate of i and j ; y_i and y_j is the y-coordinate of i and j .

3.2 Feature Selection

In the process of feature selection, the scarcity leads to the lack of detail; the over-abundance leads to redundancy and hard to recognize. So, the feature selection is very important in the expression recognition system. This paper is to reduce the redundancy and irrelevant information by using the feature selection algorithm, guaranteeing the strong recognition ability after the selection, saving the processing time.

Heuristic attribute reduction method acquires an optimal and near-optimal result through the adding of the heuristic information, cutting the research space. The former one merely starting from core attributes, obtains the core value of the decision table, then acquire the reduction set through reorganization of the attributes after the reduction in a certain order. After acquiring of an easier reduction, the mere consideration of the importance of attributes ignores the repetition of the attributes contributed into the information system. The improved one makes up the deficiency by firstly finding the minimum knowledge quantity attributes, clipping it from the condition attribute set to search any unit element among recognizable matrix, if any, it will be added into the reduction set, if not it will continue the processing the knowledge quantity; secondly, it continues the research of minimum quantity from the remained attribute reduction set, then clipping, search through the recognizable

matrix, if any adding, if not go on processing, the circle keeps going . This paper discrete the decision table by using the method of continuous decision table attributes discretization based on data distribution, it works through the improved heuristic attributes reduction method, the process can be described as:

Input: decision making system $T=(U, C \cap D)$, U as universe, the set of expression sample in decision table, C as condition attribute set (expression feature set), D decision attribute set.

Output: all satisfying condition attributes set C is the reduction RED

The processing as below:

Step 1: calculate the recognizable matrix M through a recognizable matrix algorithm based on decision table, work out the upper triangular (or the lower triangular) matrix

Step 2: initialization: $RED = \Phi$, $CORE = \Phi$, $W(a_i) = 0$.

Step 3: recheck whether there is any unit element in recognizable matrix acquired from step 2, if there any, adding all of them into the core set, $CORE = CORE \cup \{a_i\}$ transforming a recognizable matrix into an attribute set $CORE$.this step obtained might include element unit, or empty one.

Step 4: assign value for RED , $RED = CORE$, deleting all the core unit in set M , which invent the new recognizable matrix.

Step 5: calculate the knowledge quantity of every attribute according to the new matrix $W(a_i)$, calculating $|M|$.

Step 6: when $|M| \neq 0$, processing the following step: $MR = C - RED$, element in a new recognizable matrix; delete every minimum knowledge quantity attributes from matrix M , add the attribute which including single element into the $CORE$, delete all the unit including $CORE$, calculating $W_{RED}(a_i)$, calculating, $RED = RED + CORE$, calculate the new recognizable matrix.

Step 7: returns RED . Returning all the satisfying reduction

After the continuous deleting unimportant attributes in the process of improved attributes reduction, unit element may be included in the recognizable matrix; its attribute number is 1, which means this attribute is indispensable in other attributes and irreplaceable in reduction besides some deleted attributes; if unit element is not included in the matrix, which means every deleted attributes can be expressed through other attributes; unit element is continuously added into $CORE$,all the attributes set in the $CORE$ and RED constitutes the final reduction result, which is the leftover set after the deletion.

Overall, the system expressed by the reduction result is as decisive as the former information system.

4 Experiment and Result Analysis

There are nearly 2,000 graph sequences in 6 expressions grasped from more than 200 people in the Cohn-Kanade's expression bank; every sequence begins from neutral

expression to the maximum state. If randomly selecting 40 humans' facial expressions as the experiment objects, the previous 10 people's expression sequence is supposed as the training sample, the rest 30 people's is supposed as the testing sample. There are 2,400 graphs in 10 photos from weak to intensive in every individual's every expression. The processing as following:

Training

Step 1: designate 32 feature points in every face in accordance with the geometric feature points method introduced in the 2.1, acquire 36 facial expressions; using the improved attributes reduction method introduced in 2.2 to reduce 36 features , acquire 14 features, each graph owns 14 features, each expression from each one owns 140 dimensions in feature vector matrix.

Step 2:equalizing 10 feature vector matrixes ($14 * 10$) to 1 average feature vector matrix ($14 * 10$), marked as $y_i (i = 1, 2, \dots, 6)$, representing 6 feature factor models respectively.

Testing

Select the rest 30 people, who hold 10 graphs individually; put these graphs sequences in the extraction and selection respectively, acquire 140 dimension of feature matrix x_i by using Manhattan range formula.

Calculate the Manhattan range formula of 6 expression model vectors categories, the testing expression belongs to the minimum expression model among the 6, dM as the distance desired, vector x_i as testing facial feature vector, vector $y_i (i = 1, 2, \dots, 6)$ as one of the vector models among 6 . Experimental results as shown in Fig. 2.

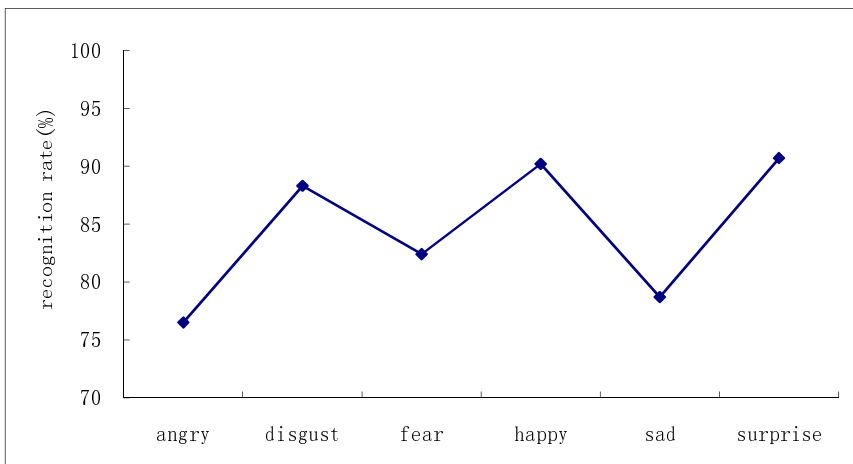


Fig. 2. Experimental results of facial expression recognition

The average recognition rate is 86.56% by using the method introduced in this paper (feature points + RS+ Manhattan distance), the rate is 84.63% by using features dot + COS. The method introduced in this paper has a higher level of recognition rate. Although it has the lower recognition rate comparing with feature points + Manhattan

distance, the amount for selection is only $1/3$, the time used is only $2/5$. It simplifies the calculation, reducing the occupancies rate on the computer resources, and greatly lowering the processing time, improving the efficiency.

5 Conclusions

The improved feature selection method based on a rough-set enable the selected feature subset to own the similar recognition ability as the former feature set through selection of the former set. Firstly extract the feature from the pretreating expression sequences, then make use of attributes reduction algorithm, finally recognize by using the Manhattan range, so greater efficiency has been made. Based on this, the next step is to integrate the whole feature extracted into the recognition system so as to better the facial recognition efficiency.

Acknowledgments. The authors would like to thank anonymous reviewers for their constructive comments on the paper. We are grateful to the department to provide facial expression database. We thank Jilin University that funded this work. This paper is supported by the Key Project of Science and Technology Development Plan for Jilin Province (Grant No. 20071152), funding by Jilin University “985 Project” Engineering Bionic Technology Innovation Platform.

References

1. Cootes, T.F., Edwards, G.J., Taylor, C.J.: Active appearance models. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 23, 681–685 (2001)
2. Cho, M.G.: A new gesture recognition algorithm and segmentation method of Korean scripts for gesture-allowed ink editor. *Information Sciences* 176(9), 1290–1303 (2006)
3. Cootes, T.F., Taylor, C.J., Cooper, D.H., Graham, J.: Active shape models – their training and application. *Computer Vision and Image Understanding* 61, 38–59 (1995)
4. Pantic, M., Rpthkrantz, L.J.M.: Expert system for automatic analysis of facial expressions. *Image and Vision Computing* 18(11), 881–905 (2000)
5. Jolliffe, I.T.: *Principal Component Analysis*. Springer, New-York (1986)
6. Krzanowski, W.J.: A Stopping Rule for structure-Preserving Variable Selection. *Statistics and Computing* 6, 51–56 (1996)
7. Ginneken, B.V., Frangi, A.F., Staal, J.J., Romeny, B.M., Viergever, M.A.: Active shape model segmentation with optimal features. *IEEE Transactions on Medical Imaging* 21(8), 924–933 (2002)
8. Khanum, A., Mufti, M., Javed, Y., et al.: Fuzzy case-based reasoning for facial expression recognition. *Fuzzy Sets and Systems* 160(2), 231–250 (2009)
9. Yan, S.C., Liu, C., Li, S.Z., Zhang, H.J., Shum, H.Y., Cheng, Q.S.: Face alignment using texture-constrained active shape models. *Image and Vision Computing* 12, 69–75 (2003)
10. Kim, J., Cetin, M., Willsky, A.S.: Nonparametric shape priors for active contour-based image, segmentation. *Signal Processing* 87(12), 3021–3044 (2007)
11. Larsen, R., Stegmann, M.B., Darkner, S., Forchhammer, S., Cootes, T.F., Ersboll, B.K.: Texture enhanced appearance models. *Computer Vision and Image Understanding* 106, 20–30 (2007)

12. Li, S.Z., Xue, Z., Teoh, E.K.: Bayesian shape model for facial feature extraction and recognition. *Pattern Recognition* 36, 2819–2833 (2004)
13. Shih, F.Y., Chuang, C.F.: Automatic extraction of head and face boundaries and facial features. *Information Sciences* 158, 117–130 (2004)
14. Wan, K.W., Lam, K.M., Ng, K.C.: An accurate active shape model for facial feature extraction. *Pattern Recognition Letters* 26, 2409–2423 (2005)
15. Wong, W.T., Shih, F.Y., Liu, J.: Shape-based image retrieval using support vector machines, Fourier descriptors and self-organizing maps. *Information Sciences* 177(8), 1878–1891 (2007)
16. Afzal, S., Sezgin, T.M., Gao, Y., Robinson, P.: Perception of emotional expressions in different representations using facial feature points. In: *Affective Computing and Intelligent Interaction and Workshops*, Amsterdam, Holland, pp. 1–6 (2009)
17. Mower, E., Mataric, M.J., Narayanan, S.S.: A framework for automatic human emotion classification using emotional profiles. *IEEE Trans. Audio Speech Language Process.* 19(5), 1057–1070 (2011)
18. Lee, C.-C., Mower, E., Busso, C., Lee, S., Narayanan, S.: Emotion recognition using a hierarchical binary decision tree approach. In: *Proc. Interspeech* (2009)

Concise Representations for State Spaces in Conformant Planning Tasks

Weisheng Li¹, Jiao Du¹, and Lifang Zhou²

¹ College of Computer Science and Technology, Chongqing University of Posts and Telecommunications, Chongqing 400065, China

² College of Software, Chongqing University of Posts and Telecommunications, Hongqing 400065, China

Liws@cqupt.edu.cn, Dujiao19880429@126.com,
Zhoulifang160966@sina.com

Abstract. Concise representations for state spaces in a conformant planning task based on the finite-domain representations are presented. The problems that literals in the finite-domain representations never used in a conformant planning task are analyzed. The paper proposed following methods: aborting void literal, translating the negative literal, and utilize domain transition graph to remove these literals. The experimental results show that the proposed methods can effectively reduce the size of state states.

Keywords: Conformant planning, State space, Literal.

1 Introduction

A conformant planning task is to decide whether there exists a linear sequence of actions that will achieve the goal from any initial state and any resolution of the non-determinism in the planning problem [1]. Smith and Weld designed the first conformant planning system CGP (Conformant Graphplan) in 1998 [2]. Hoffmann and Brafman developed the CFF (Conformant-FF), a domain independent planning system in 2004 [3]. The system extends the classical FF planner to treat initial state uncertainty expressed in the form of a CNF formula. In IPC-5 (the Fifth International Planning Competition), the IPPC (International Probabilistic Planning Competition) tracks included sub track for Probabilistic Planning, Non-Deterministic Planning and Conformant Planning [4, 5]. The conformant track as part of IPPC is hosted every other year from 2006 which boost the development of conformant planning.

The state spaces of a conformant planning task are a set of states representing all possible assignments of values to variables representing for the literals [6]. A conformant planning system consists of three modules: the translator, the knowledge compilation and the search engine. The first part is to transform the input files of a planning system into a concise representation for the world state, which is one of the important factors influencing the efficiency to solve the planning problem. In 2009, Helmert applied the FDR (Finite-Domain Representation) to represent the world state

in a PDDL (Planning Domain Definition Language) planning task [7]. FDR encodes the pair of mutually exclusive literals with numeric and propositional variables. Compared with the pure propositional representation, the FDR can represent the conformant planning task more concisely.

However, the state spaces for conformant planning based on the FDR for PPDDL planning tasks is still redundant. Firstly, some atoms which are contradict with the real world, such as atom $\text{on}(a, a)$ which means block a is on block a in the blocksworld problem of conformant track in IPC-5, are unnecessary to encode. Secondly, too many negative literals are encoded into the state variables. Finally, some encoded literals are not used to represent the state of the conformant planning task. The aim in this study is to remove the redundant state spaces in a conformant planning task based on FDR.

2 Conformant Planning

The standard domain definition language for conformant track is a subset of PPDDL added with *oneof* statements. The *oneof* statement has the followings form:

$$(\text{oneof } e_1, e_2, \dots, e_n) \quad (1)$$

where $e_k(1 \leq k \leq n)$ are probabilistic effects of the statement for uncertain action effects or:

$$(\text{oneof } s_1, s_2, \dots, s_n) \quad (2)$$

where $s_k(1 \leq k \leq n)$ are possible initial states of the statement for uncertain initial states. The semantics of *oneof* form is that when executing such form, one of e_i or s_i is chosen and applied to the plan.

A conformant planning model with FDR is given by a 5-tuple [8]:

$$\Pi = \langle V, S_0, S_G, A, O \rangle \quad (3)$$

where

V : finite sets of state variables. Each $v \in V$ is with the domain value D_v . The state variable is portioned into fluent and derived variable. Fluent is affected by operators $o \in O$. A derived variable is computed by evaluating axioms $a \in A$. The domain of a derived variable must contain the default value \perp . The undefined value means the value of this variable is unknown or does not matter.

S_0 : sets of specified all possible initial worlds. Each initial world is consisted of the conjunctions of state variable assignments over V .

S_G : sets of DNF formulas, called the goals of the task. Each goal is composed of the conjunction of state variables assignment over V .

A : sets of axioms over V . The axiom in the planning task is explained as a triple $\langle \text{cond}, v, d \rangle$, where cond is the condition or the body of the axiom, and pair $\langle v, d \rangle$ is the result or the head of the axiom. In the head of the axiom $\langle v, d \rangle$, v is a derived variable name an affect variable and $d \in D_v$ is the new value for v .

O : sets of operators with the form $\langle pre, post, prv \rangle$, where it is denoting the pre-, post and prevail-condition respectively. In the form $\langle pre, post, prv \rangle$, pre and prv is the precondition of the operator and $post$ is the effect of the operator. The variable in pre is affected variable while the variable in prv is unaffected variable. The effect post is defined as a tripe $\langle cond, v, d \rangle$, where $cond$ is a partial variable assignment, v is a fluent affected by the operator and d is a new value for V . For every operator $o = \langle pre, post, prv \rangle \in O$, it must be satisfied two restrictions shown as follows:

For all $v \in V$, if $pre[v] \neq \perp$, then $prv[v] \neq \perp$.

For all $v \in V$, $post[v] = \perp$ or $prv[v] = \perp$.

A valid plan to the conformant planning model is sets of applicable actions $\{a_0, a_1, a_2, \dots, a_n\}$. By building graph of the conformant planning task, finding a path including all of the actions $\{a_0, a_1, a_2, \dots, a_n\}$ in the graph from S_0 to S_G or prove that none of sets of actions exists.

The state space S for the planning task T is computed as

$$S(T) = \prod_{i=1}^n D_i \tag{4}$$

where n is the sum of the state variables in the planning task and D_i is the range of the state variable i .

3 Concise Representations for State Spaces

The finite-domain representation for planning tasks was successfully applied in the classical planning system such as FF [9], FastDownward [10], LAMA [11] and LAMA2011 [12]. In order to decrease the state space for conformant planning based on the FDR for PPDDL planning tasks, we propose following methods to solve the problem.

3.1 Void Literal

Void literal is that literal with different parameters convert to literal with the same arguments, i. e., atom on ($?x$ -block $?y$ -block) is grounded into on(b, b). However, void literal does not exist in the real world. Thus, abort the void literals can reduce the number of encoded literals.

3.2 Translating the Negative Literal

The negative literals are mainly distributed in the initial states and goal states in the finite-domain representations for a planning task. The negative literals in the initial states can be divided into two types: impliedly indicated by the opposite literals, and positive literals and the real negative literals. In the same way, the negative literals in the goal state can be divided into two types: impliedly indicated by the opposite literals if the goal is reachable, and need not encode the negative literal if the goal is unreachable.

In the domain definition language for the conformant track, (*oneof* s_1, s_2, \dots, s_n) indicates uncertain initial states and (*oneof* e_1, e_2, \dots, e_n) indicates uncertain action effects. After the conversion of *oneof* form, negative literals emerge in large numbers in a conformant planning task. Some experiments on domain *btuc*, *dispose*, *forest* and *uts* are shown in Fig.1.

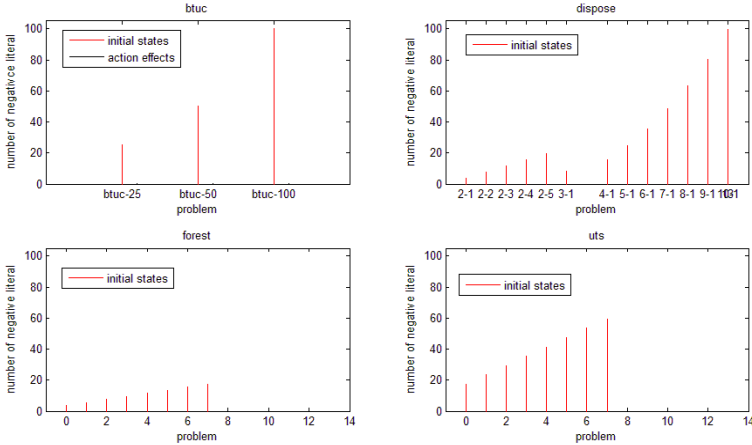


Fig. 1. Number of encoded negative literals in four domains: *btuc*, *dispose*, *forest*, and *uts*

In Fig. 1, the horizontal axis labeled the problems in IPC-domain and the vertical axis labeled the number of negative literals in the planning task. Negative literal come from the uncertain initial states and action effects in *dispose*, *forest* and *uts* and in the *btuc* domain. However, the number of negative literals derives from uncertain action effects is only one, the other negative literal are from uncertain initial states. Thus, we can reduce the negative literals derives from uncertain initial states to decrease the state spaces. Here, the *oneof* algorithm is shown as follows.

Algorithm 1. *oneof*:

Input: *oneof* (s_1, s_2, \dots, s_n)

Output: state_variable v

$i = 1$

set _{i} = { s_1, s_2, \dots, s_n }

state_variable = v_i

$D_i = \{j: s_i \mid 0 \leq i \leq n\}$

$i = i + 1$

simplify_state_variable($v_k \mid 1 \leq k \leq i$)

Return

Function simplify_stat_variable (v_k):

Invariants = {invariant₁, invariant₂, ..., invariant _{n} }

For $i = 1$ to n

$v_k = v_k - (\text{invariant}_i \cap v_k)$

Return v_k

3.3 Domain Transition Graph

After the FDR encodes the planning task, it contains all the information required to solve the task. But not all the encoded literals are useful for representing conformant planning task. Here, we utilize the domain translation graph to signify the relationship between the encoded literals in a planning task.

The domain translation graph is defined as:

$$G = \langle V, E \rangle \quad (5)$$

$$V = \{v_i \mid 1 \leq i \leq n, v_i \in D_i\} \quad (6)$$

$$E = \{e_i \mid 1 \leq i \leq n\} \quad (7)$$

where D_i is the range of the state variable i . Edge $e = \langle v_1, v_2 \rangle$, and v_1 is the precondition of the action, v_2 is the effect of the action; or v_1 is the condition of the action effect, v_2 is the effect of the action. The algorithm is illustrated as follows.

Algorithm 2. domain_transition_graph:

Input: The sets of world states S

Output: The domains of state variable var_domain

$S = \{s_i \mid 1 \leq i \leq n\}$, n is the number of all possible initial belief states

For $s \in S$

$s = \{\text{var}_j: \text{value}_j \mid 1 \leq j \leq m\}$, m is the number of variable in the task

 For var: value in s

 domain_graph = build_domain_transition_graph (var: value, task)

 reachable_domains_{var} = reachable_domains (var, domain_graph)

 End For

var_reachable_domains_j = \sum reachable_domains_i ($1 \leq j \leq m$)

var_domain = var_reachable_domains_{var}

End

Function build_domain_transiton_graph (var: value, task):

$v_0 = \text{value}$

 For action in task.action

$e = \langle \text{value}_1, \text{value}_2 \rangle$, value_1 is the affected variable value before executing action and value_2 is the affected variable value after executing action

 For axiom in task.axiom

$e = \langle \text{value}_1, \text{value}_2 \rangle$, $\text{value}_1 = -1$ and value_2 is the effect of the axiom

 End

 Function reachable_domains (var, domain_graph):

 For e_i in edge of domain_graph

 If var in e_i : ($e_i = \langle \text{var}, \text{var}_i \rangle$)

 Domain_domains + = $\{\text{var}_i\}$

 End If

 End For

 Return {domain_domains, var}

End

The process of building domain transition graph is to generate domains of the state variable for each initial state in all possible initial world states. Then, the domains of the state variable are the union set of domains from all of the initial states.

4 Experimental Results

Table 1 shows the state spaces for translator with 23 problems in five domains: blockworld, sortn, uts, coins and comm of conformant track in IPC-5. T_1 is the translator without the methods described in our paper and T_2 is the translator with our methods.

Table 1. Result of the proposed methods in problems of conformant track in IPC-5

Numbers		State variables		State space	
		T_1	T_2	T_1	T_2
problem					
block world	b2	5	5	$25*3^3$	$9*2^3$
	b4	9	9	7^4*3^5	5^4*2^5
	b5	11	11	8^5*3^6	6^5*2^6
sortn	s3	7	7	3^3*2^4	2^4
	s4	11	11	3^4*2^7	2^5
	p3	10	10	$4*3^9$	$3*2^8$
	p4	13	13	$5*3^{12}$	$4*2^{12}$
uts	p5	16	16	$6*3^{15}$	$5*2^{15}$
	p6	19	19	$7*3^{18}$	$6*2^{18}$
	p7	22	22	$8*3^{21}$	$7*2^{21}$
	p8	25	25	$9*3^{24}$	$8*2^{24}$
	p9	28	28	$10*3^{27}$	$9*2^{27}$
coins	p1	16	16	3^4*2^{12}	3^2*2^{14}
	p2	16	16	3^4*2^{12}	3^2*2^{14}
	p3	16	16	3^4*2^{12}	3^2*2^{14}
	p4	16	16	3^4*2^{12}	3^2*2^{14}
	p5	16	16	3^4*2^{12}	3^2*2^{14}
	p10	26	26	$5^4*3^2*2^{20}$	5^4*2^{21}
	p1	15	15	$6*3^4*2^{10}$	$5*2^8$
comm	p2	21	21	$7*3^6*2^{14}$	$6*2^{12}$
	p3	27	27	$8*3^8*2^{18}$	$7*2^{16}$
	p4	33	33	$9*3^{10}*2^{22}$	$8*2^{20}$
	p5	39	39	$10*3^{12}*2^{26}$	$9*2^{22}$
intro	problem of conformant track in IPC-5				

The state variables of T_1 , T_2 are the same because the two translators using the identical invariant synthesis algorithm proposed by Richter [13]. The state space of T_1 , T_2 is $\prod 3^k \times 2^{n-m-k}$, where n is the number of state variables, m is the number of invariant variables with domain values D_i , k is the number of variable with domain values $\{0: \text{positive literal}, 1: \text{negative literal}, 2: \text{none of those}\}$, $n-m-k$ is the number of variable with domain values $\{0: \text{positive literal}, 1: \text{none of those}\}$. The factors in state spaces of T_1 , T_2 are m , k , D_i and key influent factors in state space of T_1 , T_2 is k . The *oneof* algorithm proposed in our paper is to remove the negative literals in planning task. Thus, k is diminished. The *domain_transition_graph* algorithm proposed in our paper is to remove the unnecessary value of domain values. It is shown that the state space of T_2 is smaller than that of T_1 .

5 Conclusions

This paper presents concise representations for state spaces in a conformant planning task based on the finite-domain representations for PPDDL planning tasks. The problems that literals in the finite-domain representations never used in a conformant planning task are analyzed in detail. To solve this problem, the paper proposed following methods: aborting void literal, translating the negative literal, and utilize domain transition graph to decrease the state spaces. The experimental results show that the approach can reduce the size of state states effectively.

Acknowledgments. This research was supported in part by the National Natural Science Foundation of China (No. 61142011, No. 61100114).

References

1. Goldman, R.P., Boddy, M.S.: Expressive Planning and Explicit Knowledge. In: Proc. AIPS 1996 (1996)
2. Smith, D.E., Weld, D.S.: Conformant Graphplan. In: Proc. AAAI 1998 (1998)
3. Brafman, R.I., Hoffmann, J.: Conformant Planning via Heuristic Forward Search: A New Approach. *J. Artif. Intell. Res.* 14, 253–302 (2001)
4. Bonet, B., Givan, B.: Results of Conformant Track in the 5th International Planning Competition, <http://ldc.usb.vt/~bonet/ipc5/docs/results-conformant.pdf>
5. Hakan, L., Younes, S., Michael, L.L.: PPDDL1.0: An Extension to PDDL for Expressing Planning Domains with Probabilistic Effects. Technical report CMU-CS-04-167, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA (2004)
6. Sanner, S., Yoon, S.: IPPC Results Presentation. In: Proc. ICAPS 2011, pp. 16–17 (2011)
7. Helmert, M.: Concise finite-domain representations for PDDL planning tasks. *Artif. Intell.* 173, 503–535 (2009)
8. Li, W.S., Zhang, Z., Wang, W.X.: CPT-FDR: An Approach to Translating PPDDL Conformant Planning Tasks into Finite-Domain Representations. *Chinese J. Electro.* 21, 53–58 (2012)

9. Hoffmann, J., Nebel, B.: The FF Planning System: Fast Plan Generation through Heuristic Search. *J. Artif. Intell. Res.* 14, 253–302 (2001)
10. Helmert, M.: The Fast Downward Planning System. *J. Artif. Intell. Res.* 26, 191–246 (2006)
11. Richter, S., Westphal, M.: The LAMA Planner: Guiding Cost-Based Anytime Planning with Landmarks. *J. Artif. Intell. Res.* 39, 127–177 (2010)
12. Richter, S., Westphal, M., Helmert, M.: LAMA 2008 and 2011, <http://www.informatik.uni-freiburg.de/~srichter/papers/richter-et-al-ipc11.pdf>
13. Rintanen, J.: An Iterative Algorithm for Synthesizing Invariants. In: *Proc. 17th National Conference on Artificial Intelligence*. AAAI Press (2000)

On Fast Division Algorithm for Polynomials Using Newton Iteration

Zhengjun Cao* and Hanyue Cao

Department of Mathematics, Shanghai University, Shanghai, China
caozhj@shu.edu.cn

Abstract. The classical division algorithm for polynomials requires $O(n^2)$ operations for inputs of size n . Using reversal technique and Newton iteration, it can be improved to $O(M(n))$, where M is a multiplication time. But the method requires that the degree of the modulo, x^l , should be the power of 2. If l is not a power of 2 and $f(0) = 1$, Gathen and Gerhard suggest to compute the inverse, f^{-1} , modulo $x^{\lceil l/2^r \rceil}, x^{\lceil l/2^{r-1} \rceil}, \dots, x^{\lceil l/2 \rceil}, x^l$, separately. But they did not specify the iterative step. In this paper, we show that the original Newton iteration formula can be directly used to compute $f^{-1} \bmod x^l$ without any additional cost, when l is not a power of 2.

Keywords: Newton iteration, reversal technique, multiplication time.

1 Introduction

Let R be a ring (commutative, with 1). Let us define the length of a polynomial $f(X) \in R[X]$, denoted $\text{len}(f)$, to be the length of its coefficient vector; more precisely, we define

$$\text{len}(f) = \begin{cases} \deg(f) + 1 & \text{if } f \neq 0 \\ 1 & \text{if } f = 0 \end{cases}$$

Polynomials over a field form a Euclidean domain. This means that for all a, b with $b \neq 0$ there exist unique q, r such that $a = qb + r$ where $\deg r < \deg b$. The division problem is then to find q, r , given a, b . The classical division algorithm for polynomials requires $O(n^2)$ operations for inputs of size n . Concretely, we have the following results [12]:

Theorem 1. *Let a and b be arbitrary polynomials in $R[X]$.*

- (i) *We can compute $a \pm b$ with $O(\text{len}(a) + \text{len}(b))$ operations in R .*
- (ii) *We can compute $a \cdot b$ with $O(\text{len}(a)\text{len}(b))$ operations in R .*
- (iii) *If $b \neq 0$ and $lc(b)$ is a unit in R , we can compute $q, r \in R[X]$ such that $a = bq + r$ and $\deg(r) < \deg(b)$ with $O(\text{len}(b)\text{len}(q))$ operations in R .*

Fast algorithms for polynomials and integers are of great importance to computers algebra. We refer interested readers to [1-4, 6-9, 11]. In this paper, we concentrate on fast division algorithm for polynomials using Newton iteration. Using

* Corresponding author.

reversal technique and Newton iteration, it can be improved to $O(M(n))$, where M is a multiplication time. But the method requires that the degree of x^l should be the power of 2. If l is not a power of 2 and $f(0) = 1$, Gathen and Gerhard [5] suggest to compute the inverse, f^{-1} , modulo $x^{\lceil l/2^r \rceil}, x^{\lceil l/2^{r-1} \rceil}, \dots, x^{\lceil l/2 \rceil}, x^l$, separately. But they did not specify the iterative step. In this paper, we show that the original Newton iteration formula can be directly used to compute $f^{-1} \bmod x^l$ without any additional cost, when l is not a power of 2. We also correct an error in the cost analysis in Ref. [5].

2 Division Algorithm for Polynomials Using Newton Iteration

The description comes from Ref. [1].

Let R be a ring (commutative, with 1) and $a, b \in R[x]$ two polynomials of degree n and m , respectively. We assume that $m \leq n$ and that b is monic. We wish to find polynomials q and r in $R[x]$ satisfying $a = qb + r$ with $\text{degr } r < \text{degr } b$ (where, as usual, we assume that the zero polynomial has degree $-\infty$). Since b is monic, such q, r exist uniquely.

Substituting $1/x$ for the variable x and multiplying by x^n , we obtain

$$x^n a \left(\frac{1}{x} \right) = \left(x^{n-m} q \left(\frac{1}{x} \right) \right) \cdot \left(x^m b \left(\frac{1}{x} \right) \right) + x^{n-m+1} \left(x^{m-1} r \left(\frac{1}{x} \right) \right) \quad (1)$$

We define the *reversal* of a as $\text{rev}_k(a) = x^k a(1/x)$. When $k = n$, this is the polynomial with the coefficients of a reversed, that is, if $a = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$, then

$$\text{rev}(a) = \text{rev}_n(a) = a_0 x^n + a_1 x^{n-1} + \dots + a_{n-1} x + a_0$$

Equation (1) now reads

$$\text{rev}_n(a) = \text{rev}_{n-m}(q) \cdot \text{rev}_m(b) + x^{n-m+1} \text{rev}_{m-1}(r),$$

and therefore,

$$\text{rev}_n(a) \equiv \text{rev}_{n-m}(q) \cdot \text{rev}_m(b) \bmod x^{n-m+1}.$$

Notice that $\text{rev}_m(b)$ has constant coefficient 1 and thus is invertible modulo x^{n-m+1} . Hence we find

$$\text{rev}_{n-m}(q) \equiv \text{rev}_n(a) \cdot \text{rev}_m(b)^{-1} \bmod x^{n-m+1},$$

and obtain $q = \text{rev}_{n-m}(\text{rev}_{n-m}(q))$ and $r = a - qb$.

So now we have to solve the problem of finding, from a given $f \in R[x]$ and $l \in N$ with $f(0) = 1$, a $g \in R[x]$ satisfying $fg \equiv 1 \bmod x^l$. If l is a power of 2, then we can easily obtain the inversion by the following iteration step $g_{i+1} = 2g_i - fg_i^2$. In fact, if $fg_i \equiv 1 \bmod x^{2^i}$, then $x^{2^i} \mid 1 - fg_i, x^{2^{i+1}} \mid (1 - fg_i)^2$.

Hence, $x^{2^{i+1}} \mid 1 - f(2g_i - fg_i^2)$. Using the above iteration method, we have the following result:

Theorem 2. *Let R be a ring (commutative, with 1), $f, g_0, g_1, \dots, \in R[x]$, with $f(0) = 1, g_0 = 1$, and*

$$g_{i+1} \equiv 2g_i - fg_i^2 \pmod{x^{2^{i+1}}}$$

for all i . Then $fg_i \equiv 1 \pmod{x^{2^i}}$ for all $i \geq 0$.

By Theorem 1, we now obtain the following algorithm to compute the inverse of $f \pmod{x^l}$. We denote by \log the binary logarithm.

Algorithm 1: Inversion using Newton iteration

Input: $f \in R[x]$ with $f(0) = 1$, and $l \in N$.

Output: $g \in R[x]$ satisfying $fg \equiv 1 \pmod{x^l}$.

1. $g_0 \leftarrow 1, r \leftarrow \lceil \log l \rceil$
 2. **for** $i = 1, \dots, r$ **do**
 $g_i \leftarrow (2g_{i-1} - fg_{i-1}^2) \pmod{x^{2^i}}$
 3. **Return** g_r
-

From the algorithm 1, one can easily obtain the following.

Algorithm 2: Fast division with remainder

Input: $a, b \in R[x]$, where R is a ring (commutative, with 1) and $b \neq 0$ is monic.

Output: $q, r \in R[x]$ such that $a = qb + r$ and $\deg r < \deg b$.

1. **if** $\deg a < \deg b$ **then return** $q = 0$ and $r = a$
 2. $m \leftarrow \deg a - \deg b$
call Algorithm 1 to compute the inverse of $\text{rev}_{\deg b}(b) \in R[x]$ modulo x^{m+1}
 3. $q^* \leftarrow \text{rev}_{\deg a}(a) \cdot \text{rev}_{\deg b}(b)^{-1} \pmod{x^{m+1}}$
 4. **return** $q = \text{rev}_m(q^*)$ and $r = a - bq$
-

3 On the Form of l

The authors [5] stress that “if l is not a power of 2, then the above algorithm computes too many coefficients of the inverse.” They suggest to compute the inverse modulo

$$x^{\lceil l/2^r \rceil}, x^{\lceil l/2^{r-1} \rceil}, \dots, x^{\lceil l/2 \rceil}, x^l$$

For example, suppose $l = 11$, then one has to compute

$$x^{\lceil 11/2^4 \rceil} = x, x^{\lceil 11/2^3 \rceil} = x^2, x^{\lceil 11/2^2 \rceil} = x^3, x^{\lceil 11/2 \rceil} = x^6$$

In such case, one has to compute f^{-1} modulo x, x^2, x^3, x^6, x^{11} . It should be stressed that the authors did not specify the iterative step. More serious, the

sequence 1, 2, 3, 6, 11 does not form an addition chain [10]. Given a chain $\{a_i\}$ and f , we can define the following iterative step

$$g_{a_k} \equiv g_{a_i} + g_{a_j} - fg_{a_i}g_{a_j} \pmod{x^{a_k}}, \text{ if } a_k = a_i + a_j$$

In fact, the suggestion is somewhat misleading. If l is not a power of 2, the original algorithm 1 can be used to compute the inverse modulo x^l without any additional cost. It suffices to observe the following fact.

Fact 1. *If $0 < l \leq t$ and $x^t \mid 1 - fg$, then $x^l \mid 1 - fg$.*

The above fact is directly based on the divisibility characteristic. Based on the fact, we obtain the following algorithm.

Algorithm 3: Inversion using divisibility characteristic

Input: $f \in R[x]$ with $f(0) = 1$, and $l \in N$.
 Output: $g \in R[x]$ satisfying $fg \equiv 1 \pmod{x^l}$.

1. $g_0 \leftarrow 1, r \leftarrow \lceil \log l \rceil$
2. **for** $i = 1, \dots, r - 1$ **do**
 $g_i \leftarrow g_{i-1} \cdot (2 - f \cdot g_{i-1}) \pmod{x^{2^i}}$
3. $g_r \leftarrow g_{r-1} \cdot (2 - f \cdot g_{r-1}) \pmod{x^l}$
4. **Return** g_r

Correctness. It suffices to observe that $l \leq 2^r$ where $r = \lceil \log l \rceil$. Hence $x^l \mid x^{2^r}$. Since

$$x^{2^r} \mid 1 - f \cdot (2g_{r-1} - fg_{r-1}^2)$$

we have

$$x^l \mid 1 - f \cdot (2g_{r-1} - fg_{r-1}^2)$$

That means g_r is the inverse of f modulo x^l , too.

4 On the Cost Analysis

To make a sound cost analysis, we need the following definition of multiplication time and its properties.

Definition 1. *Let R be a ring (commutative, with 1). We call a function $M : N_{>0} \rightarrow R_{>0}$ a multiplication time for $R[x]$ if polynomials in $R[x]$ of degree less than n can be multiplied using at most $M(n)$ operations in R . Similarly, a function M as above is called a multiplication time for Z if two integers of length n can be multiplied using at most $M(n)$ word operations.*

For convenience, we will assume that the multiplication time satisfies

$$M(n)/n \geq M(m)/m \text{ if } n \geq m, \quad M(mn) \leq m^2M(n),$$

for all $n, m \in N_{>0}$. The first inequality yields the superlinearity properties

$$M(mn) \geq mM(n), \quad M(m+n) \geq M(n) + M(m), \quad \text{and } M(n) \geq n$$

for all $n, m \in N_{>0}$.

By the above definition and properties, the authors obtained the following result [5].

Theorem 3. *Algorithm 1 correctly computes the inverse of f modulo x^l . If $l = 2^r$ is a power of 2, then it uses at most $3M(l) + l \in O(M(l))$ arithmetic operations in R .*

Proof. In step 2, all powers of x up to 2^i can be dropped, and since

$$g_i \equiv g_{i-1}(2 - fg_i) \equiv g_{i-1} \pmod{x^{2^{i-1}}}, \tag{2}$$

also the powers of x less than 2^{i-1} . The cost for one iteration of step 2 is $M(2^{i-1})$ for the computation of g_{i-1}^2 , $M(2^i)$ for the product $fg_{i-1}^2 \pmod{x^{2^i}}$, and then the negative of the upper half of fg_{i-1}^2 modulo x^{2^i} is the upper half of g_i , taking 2^{i-1} operations. Thus we have $M(2^i) + M(2^{i-1}) + 2^{i-1} \leq \frac{3}{2}M(2^i) + 2^{i-1}$ in step 2, and the total running time is

$$\begin{aligned} \sum_{1 \leq i \leq r} \left(\frac{3}{2}M(2^i) + 2^{i-1} \right) &\leq \left(\frac{3}{2}M(2^r) + 2^{r-1} \right) \sum_{1 \leq i \leq r} 2^{i-r} \\ &< 3M(2^r) + 2^r = 3M(l) + l, \end{aligned} \tag{3}$$

where we have used $2M(n) \leq M(2n)$ for all $n \in N$.

There is a typo and an error in the above proof and theorem.

- In the above argument there is a typo (see Eq.(2)).
- The cost for one iteration of step 2 is $M(2^i)$ for the computation of g_{i-1}^2 instead of the original $M(2^{i-1})$, because it is computed under the module x^{2^i} , not $x^{2^{i-1}}$. Since the upper half of $f(g_{i-1}^2)$ modulo x^{2^i} is the same as g_i and the lower half of g_i is the same as g_{i-1} , the cost for the computation of $f(g_{i-1}^2)$ modulo x^{2^i} only needs $M(2^{i-1})$. Therefore, according to the original argument the bound should be

$$\begin{aligned} \sum_{1 \leq i \leq r} \left(\frac{3}{2}M(2^i) + 2^{i-1} \right) &\leq \left(\frac{3}{2}M(2^r) + 2^{r-1} \right) \sum_{1 \leq i \leq r} 2^{i-r} \\ &< 3M(2^r) + 2^r \leq 12M(l) + 2l \end{aligned} \tag{4}$$

The last estimation comes from $l \leq 2^r \leq 2l$.

Now, we make a formal cost analysis of algorithm 3.

Theorem 4. *Algorithm 3 correctly computes the inverse of f modulo x^l . It uses at most $5M(l) + l \in O(M(l))$ arithmetic operations in R .*

Proof. The cost for step 2 is $3M(2^{r-1}) + 2^{r-1}$ (see the above cost analysis). The cost for step 3 is bounded by $2M(l)$. Since $2^{r-1} \leq l \leq 2^r$, the total cost is $5M(l) + l$.

5 Conclusion

In this paper, we revisit the fast division algorithm using Newton iteration. We show that the original Newton iterative step can be still used for any arbitrary exponent l without the restriction that l should be the power of 2. We also make a formal cost analysis of the method. We think the new presentation is helpful to grasp the method entirely and deeply.

Acknowledgements. This work is supported by the National Natural Science Foundation of China (Project 11171205), and the Key Disciplines of Shanghai Municipality (S30104).

References

1. Bailey, D., Lee, K., Simon, H.: Using Strassen's Algorithm to Accelerate the Solution of Linear Systems. *J. of Supercomputing* 4(4), 357–371 (1991)
2. Burnikel, C., Ziegler, J.: Fast recursive division. Research report MPI-I-98-1-022, Max-Planck-Institut für Informatik, Saarbrücken, Germany (1998)
3. Cantor, D., Kaltofen, E.: On fast multiplication of polynomials over arbitrary algebras. *Acta Informatica* 28, 693–701 (1991)
4. Gathen, J.: Counting decomposable multivariate polynomials. *Appl. Algebra Eng. Commun. Comput.* 22(3), 165–185 (2011)
5. Gathen, J., Gerhard, J.: *Modern computer Algebra*, 3rd edn. Cambridge University Press (2003)
6. Gathen, J., Mignotte, M., Shparlinski, I.: Approximate polynomial GCD: Small degree and small height perturbations. *J. Symbolic Computation* 45(8), 879–886 (2010)
7. Gathen, J., Panario, D., Richmond, L.: Interval Partitions and Polynomial Factorization. *Algorithmica* 63(1-2), 363–397 (2012)
8. Jebelean, T.: Practical Integer Division with Karatsuba Complexity. In: *Proc. ISSAC 1997*, pp. 339–341. ACM Press (1997)
9. Johnson, S., Frigo, M.: A modified split-radix FFT with fewer arithmetic operations. *IEEE Trans. Signal Processing* 55(1), 111–119 (2007)
10. Knuth, D.: *The Art of Computer programming*, 3rd edn., vol. 2. Addison-Wesley (1997)
11. Lundy, T., Buskirk, J.: A new matrix approach to real FFTs and convolutions of length $2k$. *Computing* 80(1), 23–45 (2007)
12. Shoup, V.: *A Computational Introduction to Number Theory and Algebra*. Cambridge University Press (2005)

On the Security of an Improved Password Authentication Scheme Based on ECC

Ding Wang^{1,2}, Chun-guang Ma^{1,*}, Lan Shi¹, and Yu-heng Wang³

¹ College of Computer Science and Technology, Harbin Engineering University 145 Nantong Street, Harbin City 150001, China

² Automobile Management Institute of PLA, Bengbu City 233011, China

³ Golisano College of Computing and Information Sciences, Rochester Institute of Technology
102 Lomb Memorial Dr., Rochester, NY 14623, USA
wangdingg@mail.nankai.edu.cn

Abstract. The design of secure remote user authentication schemes for mobile applications is still an open and quite challenging problem, though many schemes have been published lately. Recently, Islam and Biswas pointed out that Lin and Hwang et al.'s password-based authentication scheme is vulnerable to various attacks, and then presented an improved scheme based on elliptic curve cryptography (ECC) to overcome the drawbacks. Based on heuristic security analysis, Islam and Biswas claimed that their scheme is secure and can withstand all related attacks. In this paper, however, we show that Islam and Biswas's scheme cannot achieve the claimed security goals and report its flaws: (1) It is vulnerable to offline password guessing attack, stolen verifier attack and denial of service (DoS) attack; (2) It fails to preserve user anonymity. The cryptanalysis demonstrates that the scheme under study is unfit for practical use.

Keywords: Authentication protocol, Elliptic curve cryptography, Cryptanalysis, Smart card, User anonymity.

1 Introduction

Since Lamport [1] introduced the first password-based authentication scheme in 1981, many password-based remote user authentication schemes [2–6] have been proposed, where a client remembers a password and the corresponding server holds the password or its verification data that are used to verify the client's knowledge of the password. These easy-to-remember passwords, called weak passwords, have low entropy and thus are potentially vulnerable to various sophisticated attacks, especially offline password guessing attack [7], which is the gravest threat a well-designed password authentication scheme must be able to thwart. A common feature among the published schemes is that computation efficiency and system security cannot be achieved at the same time. As the computation ability and battery capacity of mobile devices (e.g. PDAs, smart cards) are limited, the traditional public-key based remote authentication schemes are not suitable for mobile applications.

* Corresponding author.

Fortunately, it seems to see the dawn in recent two years, where several schemes based on ECC have been proposed to reduce computation cost while preserving security strength [8–12]. However, the reality of the situation is that this dilemma is only partially addressed and most of the ECC-based schemes were found severely flawed shortly after they were first put forward, so intensive further research is required. More recently, Islam and Biswas [13] proposed an advanced password authentication scheme based on ECC. The authors claimed that their scheme provides mutual authentication and is free from all known cryptographic attacks, such as replay attack, offline password guessing attack, insider attack and so on. Although their scheme is superior to the previous solutions for implementation on mobile devices, we find their scheme cannot achieve the claimed security: their scheme is vulnerable to the offline password guessing attack, the stolen verifier attack. Almost at the same time with us, He et al. [14] also have identified these defects in Islam-Biswas’s scheme. Hence, we went on to perform a further cryptanalysis on this protocol and observe that it is also prone to a denial of service (DoS) attack, and it transmits user’s identity in plain during the login request and thus user anonymity is not provided, while provision of user identity confidentiality is of great importance for a protocol in mobile environments [15].

The remainder of this paper is organized as follows: in Section 2, we review Islam-Biswas’s scheme. Section 3 describes the weaknesses of Islam-Biswas’s scheme. Section 4 concludes the paper.

2 Review of Islam-Biswas’s Scheme

In this section, we examine the password authentication scheme using smartcards proposed by Islam and Biswas [13] in 2011. Islam-Biswas’s scheme, summarized in Fig.1, consists of four phases: the registration phase, the authentication phase, the session key distribution phase and the password change phase. For ease of presentation, we employ some intuitive abbreviations and notations listed in Table 1.

Table 1. Notations

Symbol	Description
U_i	i^{th} user
S	remote server
ID_i	identity of user U_i
PW_i	password of user U_i
d_s	secret key of remote server S
G	base point of the elliptic curve group of order n such that $n \cdot G = O$
V_s	public key of remote server S , where $V_s = d_s \cdot G$
V_i	password-verifier of U_i , where $V_i = PW_i \cdot G$
K_x	secret key computed using $K = d_s \cdot V_i = PW_i \cdot V_s = (K_x, K_y)$
$E_{K_x}(\cdot)$	symmetric encryption with K_x
$H(\cdot)$	collision free one-way hash function
\oplus	the bitwise XOR operation
\parallel	the string concatenation operation
$A \Rightarrow B : M$	message M is transferred through a secure channel from A to B
$A \rightarrow B : M$	message M is transferred through a common channel from A to B

2.1 Registration Phase

Before the system begins, the server selects a large prime number p and two integer elements a and b , where $p > 2^{160}$ and $4a^3 + 27b^2 \pmod p \neq 0$. Then the server selects an elliptic curve equation E_p over finite field F_p : $y^2 = x^3 + ax + b \pmod p$. Let G be a base point of the elliptic curve with a prime order n and \mathcal{O} be a point at infinite, where $n \cdot G = \mathcal{O}$ and $n > 2^{160}$. The server chooses the private key d_s and computes the public key $V_s = PW_i \cdot G$. The registration phase involves the following operations:

Step R1. U_i chooses his identity ID_i and password PW_i , then computes $V_i = PW_i \cdot G$.

Step R2. $U_i \Rightarrow S: \{ID_i, V_i\}$.

Step R3. On receiving the registration message from U_i , the server S create an entry $(ID_i, V_i, \text{status-bit})$ in its database, where the *status-bit* indicates the status of the client, i.e., when the client is logged-in to the server the *status-bit* is set to one, otherwise it is set to zero.

2.2 Authentication Phase

When U_i wants to login to S , the following operations will be performed:

Step L1. U_i keys his identity ID_i and the password PW_i into the terminal. The client selects a random number r_i from $[1, n-1]$, computes $R_i = r_i \cdot V_s$ and $W_i = (r_i \cdot PW_i) \cdot G$. Then encrypts (ID_i, R_i, W_i) using a symmetric key K_x , where K_x is the x coordinate of $K = PW_i \cdot V_s = (K_x, K_y)$.

Step L2. $U_i \Rightarrow S: \{ID_i, E_{K_x}(ID_i \parallel R_i \parallel W_i)\}$.

Step L3. S computes the decryption key K_x by calculating $K = d_s \cdot V_i = (K_x, K_y)$ and then decrypts $E_{K_x}(ID_i \parallel R_i \parallel W_i)$ using K_x . Subsequently S compares decrypted ID_i with received ID_i , $\hat{e}(R_i, V_i)$ with $\hat{e}(W_i, V_s)$, respectively. If both conditions are satisfied, S selects a random number r_s and computes $W_s = r_s \cdot V_s = r_s \cdot d_s \cdot G$.

Step L4. $S \rightarrow U_i: \{W_i + W_s, H(W_s)\}$.

Step L5. U_i retrieves W_s by subtracting W_i from $W_i + W_s$. If the hashed result of retrieved W_s is equal to the received $H(W_s)$, then U_i performs the hash operation $H(W_i \parallel W_s)$ and sends it to the server.

Step L6. $U_i \rightarrow S: \{H(W_i \parallel W_s)\}$.

Step L7. The server S computes the hash value with its own copies of W_s and W_i and compares it with the received $H(W_i \parallel W_s)$, to accept or denied the login request. If the equality holds, the server grants the client's login request, otherwise rejects.

2.3 Session Key Distribution Phase and Password Change Phase

Since both the session key distribution phase and password change phase have little relevance with our discussion, they are omitted here.

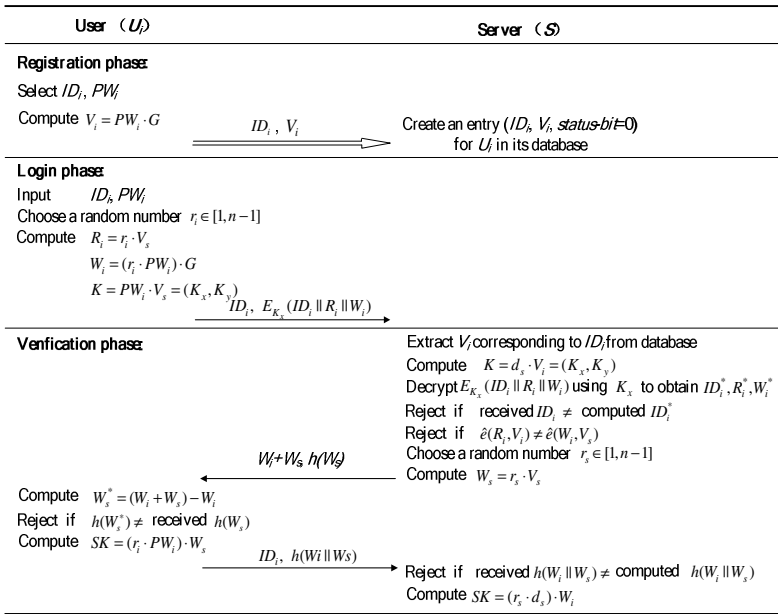


Fig. 1. Islam-Biswas’s remote user authentication scheme

3 Cryptanalysis of Islam-Biswas’s Scheme

With superior performance over other related schemes and a long list of arguments of security features that their scheme possesses presented, Islam-Biswas’s scheme seems desirable at first glance. However, their security arguments are still specific-attack-scenario-based and without some degree of rigorousness, and thus it is not fully convincing. We find that Islam-Biswas’s scheme still fails to serve its purposes and demonstrate its security flaws in the following.

3.1 Offline Password Guessing Attack

A remote user authentication scheme which is vulnerable to the offline password guessing attack must satisfy the following two conditions: (1) the user’s password is weak, and (2) there exists a piece of password-related information used as a comparison target for password guessing.

In Islam-Biswas’s scheme, a user is allowed to choose her own password at will during the registration and password change phases; the user usually tends to select a password, e.g., his birthday, which is easily remembered for his convenience. Hence, these easy-to-remember passwords, called weak passwords, have low entropy and thus are potentially vulnerable to offline password guessing attack.

Besides, user U_i ’s identity is transmitted in plaintext within the login request, it is not difficult for an adversary \mathcal{A} to identify the login request message sent by U_i .

Once the login request message $\{ID_i, E_{K_x}(ID_i \parallel R_i \parallel W_i)\}$ during any authentication process is intercepted by \mathcal{A} , an offline password guessing attack can be launched as follows:

- Step 1.** Guesses the value of PW_i to be PW_i^* from a dictionary space \mathcal{D} .
- Step 2.** Computes $K^* = PW_i^* \cdot V_s = (K_x^*, K_y^*)$, as V_s is the public key of server S .
- Step 3.** Decrypts the previously intercepted $E_{K_x}(ID_i \parallel R_i \parallel W_i)$ using K_x^* to obtain ID_i^* .
- Step 4.** Verifies the correctness of PW_i^* by checking if the computed ID_i^* is equal to the intercepted ID_i .
- Step 5.** Repeats Steps 1, 2, 3, and 4 of this procedure until the correct value of PW_i is found.

As the size of the password dictionary, i.e. $|\mathcal{D}|$, is very limited in practice, the above attack procedure can be completed in polynomial time. Moreover, the above attack we describe is very effective because it only requires the abilities of an eavesdropping attacker, and involves no expensive cryptographic operations.

After guessing the correct value of PW_i , \mathcal{A} can compute the valid symmetric key $K = PW_i \cdot V_s = (K_x, K_y)$. Then the attacker can impersonate U_i to send a valid login request message $\{ID_i, E_{K_x}(ID_i \parallel R_i \parallel W_i)\}$ to the service provider server S , since U_i 's identity ID_i can be intercepted from the channel and W_i can be fabricated with the correctly guessed PW_i . Upon receiving the fabricated login request, S will find no abnormality and responds with $\{W_i + W_s, H(W_s)\}$. Then \mathcal{A} can compute the valid W_s since she knows W_i . Hence the attacker \mathcal{A} can successfully masquerade as a legitimate user U_i to server S . On the other hand, the attacker may also impersonate the server S to U_i successfully in a similar way.

3.2 Stolen Verifier Attack

Let us consider the following scenarios. In case the verifier table in the database of the server S is leaked out or stolen by an adversary \mathcal{A} . With the obtained entry $(ID_i, V_i, status-bit)$ corresponding to U_i , she can guess out the password PW_i of U_i using the method as follows:

- Step 1.** Guesses the value of PW_i to be PW_i^* from a uniformly distributed dictionary.
- Step 2.** Computes $V_i^* = PW_i^* \cdot G$, as G is public.
- Step 3.** Verifies the correctness of PW_i^* by checking if the computed V_i^* is equal to the somehow obtained V_i .
- Step 4.** Repeats Steps 1, 2, and 3 of this procedure until the correct value of PW_i is found.

As the password dictionary size is very limited, the above attack procedure can be completed in polynomial time. Since the underlying assumption of the above attack introduced is much constrained, it is much less effective than the attack introduced in Section 3.1. However, it is still an insecure factor to be noticed.

3.3 Failure of Protecting the User's Anonymity

As violation concern of user privacy on e-commerce and industrial engineering applications is promptly raised among individuals, human right organizations and national governments, identity protection has become a very popular research topic in recent years. Many systems have been advanced, which implement different (and sometimes even contradictory) notions of what it means to be “anonymous. Instead of a single anonymity property, there are dozens of different flavors of anonymity, such as sender un-traceability, blender anonymity, sender k-anonymity and so on [16]. As for remote authentication schemes, user anonymity basically means initiator anonymity (i.e., sender anonymity), more precisely, it means the adversary could not have any knowledge of real identity of the initiator but may know whether two conversations originate from the same (unknown) entity. Comparatively, a more ideal anonymity property is initiator un-traceability (i.e., sender un-traceability), which means that the adversary can know neither who the initiator is nor whether two conversations originate from the same (unknown) initiator. A protocol with user anonymity prevents an adversary from acquiring sensitive personal information about an individual's preferences, lifestyles, social circle, shopping patterns, current location, etc. by analyzing the login information.

In Islam-Biswas's scheme, the user's identity ID is transmitted in plain, which may leak the identity of the logging user once the login messages were eavesdropped; the user's identity ID is static in all the login phases, which may facilitate the attacker to trace out the different login request messages belonging to the same user and to derive some information related to the user U_i . In a word, neither initiator anonymity nor initiator un-traceability can be preserved in their scheme.

3.4 Denial of Service Attack

Without any knowledge of the user private information like password or security parameters stored in smart card, an adversary \mathcal{A} can successfully launch a kind of denial of service attack, which is the so called “clogging attack” [17], in many non-DoS-resilient cryptography protocols. Let's see how this could happen with Islam-Biswas's scheme in place. The following is performed by the adversary \mathcal{A} :

Step 1. Sends the previously intercepted $\{ID_i, E_{K_x}(ID_i \parallel R_i \parallel W_i)\}$ to the server S .

Step 2. Ignores the reply from the server S .

The following is performed by the server:

Step 1'. On receiving the login request from U_i (actually \mathcal{A}), S computes the decryption key K_x by calculating $K = d_s \cdot V_i = (K_x, K_y)$ and then decrypts $E_{K_x}(ID_i \parallel R_i \parallel W_i)$ using K_x . Subsequently S compares decrypted ID_i with received ID_i , $\hat{e}(R_i, V_i)$ with $\hat{e}(W_i, V_s)$, respectively.

Step 2'. Selects a random number r_s and computes $W_s = r_s \cdot V_s = r_s \cdot d_s \cdot G$.

Step 3'. Sends out $\{W_i + W_s, H(W_s)\}$ and waits for the response from U_i (actually \mathcal{A}), which will never come.

Since ID_i and $E_{K_x}(ID_i \parallel R_i \parallel W_i)$ are valid, S will find no abnormality in Step 1' and then proceeds to Step 2'.

The point here is that, in the above attack, the adversary \mathcal{A} does not need to perform any special or expensive cryptographic operations but sending one message out. However, on the server side, in Step 1', S needs to perform one symmetric-key decryption and one bilinear pairing operation, which are computationally intensive. According to [18], the cost of one bilinear pairing operation is twenty times higher than that of one scale multiplication, and two times higher than that of one modulo exponentiation at the same security level. It should be noted that even DoS-resilient mechanisms (e.g. timeout or locking user account for a period of time after a predefined number of login failures) are introduced on server side, it may be not a real obstacle for attacker \mathcal{A} as it can initialize new sessions with different intercepted identities in an interleaving manner. Hence, \mathcal{A} can potentially performs the above attack procedure continuously, which will make the victimized server keeps computing the useless expensive operations rather than any real work. Thus \mathcal{A} clogs S with useless work and therefore S denies any legitimate user any service. If distributed DoS attacks are launched based on this strategy, the consequences will be more serious.

4 Conclusion

Smartcard-based password authentication technology has been widely deployed in various kinds of security-critical applications, and careful security considerations should be taken when designing such schemes. In this paper, we have shown that Islam-Biswas's scheme suffers from the offline password guessing attack, stolen-verifier attack and denial of service attack. In addition, their scheme fails to provide the property of user anonymity. In conclusion, although Islam-Biswas's scheme is very efficient and possesses many attractive features, it, in fact, does not provide all of the security properties that they claimed and only radical revisions of the protocol can possibly eliminate the identified flaws. Therefore, the scheme under study is not recommended for practical applications. In future work, we will propose an improvement over Islam-Biswas's scheme to overcome the identified drawbacks.

Acknowledgements. The authors would like to thank the anonymous reviewers for their valuable comments and constructive suggestions. This research was supported by the National Natural Science Foundation of China (NSFC) under Grants No. 61170241 and No. 61073042.

References

1. Lamport, L.: Password authentication with insecure communication. *Communications of the ACM* 24(11), 770–772 (1981)
2. Liao, I.E., Lee, C.C., Hwang, M.S.: A password authentication scheme over insecure networks. *Journal of Computer and System Sciences* 72(4), 727–740 (2006)

3. Song, R.: Advanced smart card based password authentication protocol. *Computer Standards & Interfaces* 32(5), 321–325 (2010)
4. Yeh, K.H., Su, C., Lo, N.W., Li, Y., Hung, Y.X.: Two robust remote user authentication protocols using smart cards. *Journal of Systems and Software* 83(12), 2556–2565 (2010)
5. Ma, C.-G., Wang, D., Zhang, Q.-M.: Cryptanalysis and Improvement of Sood et al.'s Dynamic ID-Based Authentication Scheme. In: Ramanujam, R., Ramaswamy, S. (eds.) ICDCIT 2012. LNCS, vol. 7154, pp. 141–152. Springer, Heidelberg (2012)
6. Ma, C.-G., Wang, D., Zhao, P., Wang, Y.-H.: A New Dynamic ID-Based Remote User Authentication Scheme with Forward Secrecy. In: Wang, H., Zou, L., Huang, G., He, J., Pang, C., Zhang, H.L., Zhao, D., Yi, Z. (eds.) APWeb 2012 Workshops. LNCS, vol. 7234, pp. 199–211. Springer, Heidelberg (2012)
7. Klein, D.V.: Foiling the cracker: A survey of, and improvements to, password security. In: *Proceedings of the 2nd USENIX Security Workshop*, pp. 5–14 (1990)
8. Wang, R.C., Juang, W.S., Lei, C.L.: Robust authentication and key agreement scheme preserving the privacy of secret key. *Computer Communications* 34(3), 274–280 (2011)
9. Wu, S.H., Zhu, Y.F., Pu, Q.: Robust smart-cards-based user authentication scheme with user anonymity. *Security and Communication Networks* 5(2), 236–248 (2012)
10. Wei, J., Hu, X., Liu, W.: An improved authentication scheme for telecare medicine information systems. *Journal of Medical Systems* (2011), doi:10.1007/s10916-012-9835-1
11. Pu, Q., Wang, J., Zhao, R.: Strong authentication scheme for telecare medicine information systems. *Journal of Medical Systems* (2011), doi:10.1007/s10916-011-9735-9
12. He, D.B., Chen, J.H., Zhang, R.: A more secure authentication scheme for telecare medicine information systems. *Journal of Medical Systems* (2011), doi:10.1007/s10916-011-9658-5
13. Islam, S.H., Biswas, G.P.: Design of improved password authentication and update scheme based on elliptic curve cryptography. *Mathematical and Computer Modelling* (2011), doi:10.1016/j.mcm.2011.07.001
14. He, D.B.: Comments on a password authentication and update scheme based on elliptic curve cryptography. *Cryptology ePrint Archive, Report 2011/411* (2011), <http://eprintH.iacr.org/2011/411.pdf>
15. Wan, Z., Zhu, B., Deng, R.H., Bao, F., Ananda, A.L.: Dos-resistant access control protocol with identity confidentiality for wireless networks. In: *2005 IEEE Wireless Communications and Networking Conference (WCNC)*, vol. 3, pp. 1521–1526. IEEE Press, New York (2005)
16. Hughes, D., Shmatikov, V.: Information hiding, anonymity and privacy: a modular approach. *Journal of Computer Security* 12(1), 3–36 (2004)
17. Roy, S., Das, A.K., Li, Y.: Cryptanalysis and security enhancement of an advanced authentication scheme using smart cards, and a key agreement scheme for two-party communication. In: *2011 IEEE 30th International Performance Computing and Communications Conference (IPCCC)*, pp. 1–7. IEEE Press, New York (2011)
18. Cao, X., Kou, W., Du, X.: A pairing-free identity-based authenticated key agreement protocol with minimal message exchanges. *Information Sciences* 180(15), 2895–2903 (2010)

Strategies to Develop Personal Knowledge Management Ability Based on M-Learning

Lin Hu

Jilin University of Finance and Economics, Changchun, China
huhu315@126.com

Abstract. With m-learning being more and more popular among college students, there exist some problems in m-learning. How to achieve the best learning effect in m-learning has become a hot issue among researchers and educators. PKM (personal knowledge management) draws our attention because under the circumstances of m-learning, there are no teachers to instruct students' study. Then students' personal knowledge management ability is quite important. How to develop their PKM ability? In this paper, the author puts forward several strategies to develop students' PKM ability.

Keywords: strategy, mobile learning, personal knowledge management, personal knowledge management ability.

1 Introduction

Mobile learning (m-learning) refers to study anytime, anywhere by means of mobile digital devices [1]. The mobile devices must effectively present the study content as well as can provide mutual communications between teachers and students. With rapid development of mobile technology, mobile learning becomes possible for everyone. In mobile learning contexts, because of the lack of teachers' supervision and guidance, the awareness of self-management ability is vital for m-learning effect.

2 Problems Existing in M-Learning

As students put more emphasis on English learning and mobile devices have higher and higher penetration, m-learning is more and more popular among students, especially among college students, who have relatively lower study pressure, relative flexible study time and are adept at internet technology. So there appears an embarrassing phenomenon in colleges. College students don't listen to their teachers whereas after class they will study by themselves. In general, m-learning of college students lacks of planning. Although they have more time spent in m-learning, the effects are not obvious. They are clearly aware of the assistance function of mobile devices to their study but they don't know the features of m-learning and they can't scientifically and effectively arrange time and content. Some students have lower ability to search for useful m-learning resources on the net. So they can't guarantee

high study efficiency. During the study process, students aren't aware of the importance of participatory interactive and cooperative study so that m-learning devices can't be well-used to exert its study assistance role. What's more, students and teachers lack of interaction under the circumstances of m-learning. Teachers can't grasp the studies of the students. So the immediate interactive function of the mobile devices can't be fully made use of.

3 Knowledge Management Ability and Its Function on College English M-Learning

Considering knowledge management, some say it came from an utterance of Peter Druck, whereas some say it came from Enovation Consultant Firm in Massachusetts at the beginning of the nineties. The initial application of it is in the field of enterprise management. Personal knowledge management is the outcome which is applied at the individual level. It can be understood as a kind of thought, a kind of device to make personal improvement in mobile context. Before discussing personal knowledge management, it is quite necessary to get to know relative information about knowledge.

The definition of personal knowledge management was put forward by Prof. Paul Dorsey of the University of Michigan, USA. He thought that personal knowledge management should be regarded as a set of skills and methods of solving problems not only on the logic level but on the actual level. To be specific, it should include the skills of information search, evaluation, organization, analysis and information collaboration etc. Prof. Paul Dorsey put forward seven skills (as is shown in Fig. 1) of personal knowledge management which include: obtaining information, evaluating information, organizing information, analyzing information, expressing information, ensuring information safely and collaborating information. These seven skills are quite important in the course of personal knowledge management.

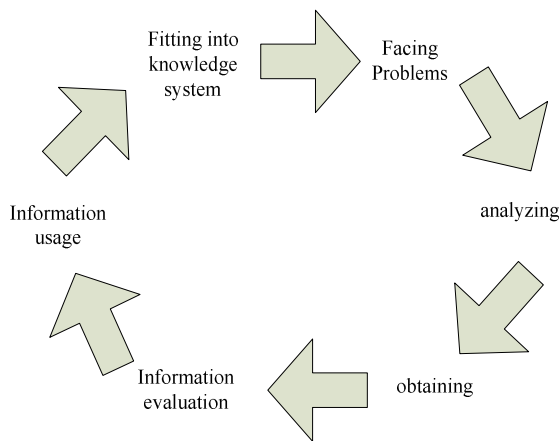


Fig. 1. Personal Knowledge Application Model

We could understand the definition of personal knowledge management from three aspects. At first, we can manage the knowledge we have got. Then, we could learn new knowledge, learn from the experience of others to make up the gap of knowledge and thought. At last, we can make use of the knowledge that we have grasped plus others' essence of thought to realize the transmission from implicit knowledge to explicit knowledge as well as to stimulate knowledge innovation.

The future society puts forward challenges for personal development. As knowledge is updated faster and faster, traditional education mode hasn't adapted to the needs of the era, which requires everyone to perform lifelong education which is also the target of m-learning. Because m-learning has its specific features, personal knowledge management ability has great role in it. Knowledge management asks people to employ modern technological methods to accumulate, transform, share, create and renew knowledge. With the precondition of insisting on study, on one hand, everyone should try his/her best to systemize knowledge. On the other hand, because everyone has his/her own knowledge structure, we should consciously communicate with others, sharing and creating knowledge, to get the best learning efficiency.

4 Fostering Management Ability in M-Learning

Before you begin to format your paper, first write and save the content as a separate text file. Keep your text and graphic files separate until after the text has been formatted and styled. Do not use hard tabs, and limit use of hard returns to only one return at the end of a paragraph. Do not add any kind of pagination anywhere in the paper. Do not number text heads-the template will do that for you.

Finally, complete content and organizational editing before formatting. Please take note of the following items when proofreading spelling and grammar:

4.1 Correcting Study Attitude

Students should have correct understanding to m-learning [5]. Students should be clear about the meaning of m-learning: arousing their potential; promoting their development; caring about their own needs; pursuing distinguishing development; cultivating study ability; realizing sustainable development. The spoon-fed education in our country has wiped out students' innate passion for knowledge. Students have been the containers for knowledge. Although knowledge is right, students are oppressive. New era will bring students the right of choice. The popularity of mobile devices will bring the best opportunity for students to experience m-learning.

4.2 Fostering the Awareness of Management

Fostering the awareness of management is the precondition of achieving the best mobile learning effect [6]. Teachers may explain the importance of self-management to the students through lectures, group discussions etc. While cultivating the

awareness of management, teachers should pay special attention to the following. Firstly, students should be clear about the target of study. Teachers should tell students the specific needs and levels they should reach during teaching. Secondly, students should have the awareness of grasping the learning materials. Teachers have the responsibility to cultivate students to analyze learning materials to be aware of its structure, difficulty level in order to reasonably distribute time and attention. Thirdly, students should have the awareness of using self-management strategies. Fourthly, students should be clear about their own cognitive styles to explore learning methods suitable for them.

4.3 Helping Students Select M-Learning Resources

Abundant learning resources are the basis of m-learning but the information resources on the net are numerous and uneven in quality. Some students often can't resist the temptation to click or download some irrelevant links with their study. They are often lost in the front of numerous links. Most students are aware of the assistance role of m-learning but when looking up a number of learning materials, they often make a hasty inspection and make a choice at random. Some students even lose interest because of the messy information and abandon looking up. Teachers should take the responsibility for helping students select what materials are good for their study.

4.4 Helping Students Establish Learning Target and Work Out Study Plan

Establishing learning target can inspire students to study. Under the exam-oriented education system, freshmen are not clear about the features of college English study and the importance of m-learning in college. So they are blind to study. After freshmen enter the college, teachers must instruct them in order to help them adapt to college study. What's more, students should have correct self-positioning. Teachers may grasp the existing problems in English study through interview and questionnaire etc. Teachers also may help students find their gap through admission test or level test etc. Under the instruction of teachers, students work out long-term study plan and short-term study plan according to their own conditions. The plans are involved in the study of the textbooks and the study of extended knowledge.

4.5 Creating Actual M-Learning Tasks to Inspire Students' Motivation

Learning motivation is an inner process or inner mental state that inspires learning activities, maintains inspired learning activities and leads to a fixed learning target [7]. Through arousing learning motivation, students may generate intense interest in m-learning so that they will strengthen their initiative of self-supervision to elevate their level of meta-cognitive supervision. The tasks they undertake in college English classes should have close connection with their everyday lives and also should be challenging. Only such tasks can arouse students' interest and determination to settle them. For example, when studying the text "Writing Three Thank You Letters",

students should be reminded of thinking about their beloved people. Teachers should ask them to write letters to their beloved people to thank them. Through this familiar and interesting task, students may grasp and understand the knowledge in this text. In the above task, students can form groups to discuss, cooperate, communicate, and report what they have found to cultivate team spirit as well as exercise the ability to settle actual problems. Teamwork contributes to improve supervision ability of meta-cognition.

4.6 Inspiring Students' Self-management Experience through Mobile Devices

Reflection is a very important form to carry out self-supervision. Through reflection, self-management experience can be stimulated [8]. Students will gradually better their self-management ability and gradually form better learning habits. During college English teaching, teachers require students not only to know what to do but also to continuously make self-reminding and self-reflection.

Because there are many ways to solve problems, students achieve their full potential while exploring the solutions to the problems. Students should give more thought about teachers' solutions, partners' solutions and their own solutions. We actively stimulate students' initiative of self-doubting and self-reflecting on which students organize and adjust their studies. Teachers may direct students to keep self-reflection diaries which students may upload to a designated email through their mobile devices. Teachers may judge students' achievements through these self-reflection diaries.

4.7 Elevating Students' Supervision Ability

During the process of m-learning, students should often check out their performance of phrasal targets, their studies and study fruits. Supervision is to check out the performance of the tasks, to see if they have made profitable use of time, to analyze and evaluate their cognitive activities in time [9]. In the stage of selecting and understanding supervision strategies, students should carry out supervision to the strategies. Meanwhile, students should continuously evaluate their learning situation. Teachers may provide instruction of supervision strategies before and after study activities. So teachers should construct a context to offer their support for students to activate the strategies and methods. Teachers should provide more examples on the basis of teaching content, clearly explaining the scopes of application and conditions of application of the strategies and providing enough strategy practice. Teachers may urge students to keep learning diaries or weekly diaries to make students reflect their learning process. Teachers also may instruct students to form study group made up of four or five students. Each member makes his/her study plan public and the other members can supervise his/her plan implementation. Each group may hold a discussion each week to check out the performance of the plans as well as to share information.

Here, we strongly recommend teachers to urge students to use “self-questioning sheet” to supervise their studies. Through this self-questioning sheet, students will gradually be familiar with and accept this method. As time goes by, students may properly revise their plans and strategies according to the self-questioning sheet. Finally, students can design the sheet by themselves rather than by teachers. The transition from other-directed to self-directed is achieved [10]. In the following table 1, the three stages of learning process are showed in the self-questioning sheet.

Table 1. Three Stages of Self-questioning Sheet

Before learning	During the learning process	After learning
What are the tasks fulfilled in the class?	Could you finish your tasks according to your own plans?	What have I learned?
What are the components of the tasks? What skills, methods and resources do you need to grasp to finish the tasks?	Are your procedures correct? Could the methods mentioned in the textbook be changed? Are there any simpler methods?	Are the learning effects the same with my previous expectations?
Have you ever met similar tasks? Are there any reference functions to the tasks?	What if I meet some difficulties when finishing the tasks?	Am I satisfied with the effects? Are there any revisions?
		While looking back on the whole process, what should I do next time?

4.8 Helping Students Make Correct Self-evaluation

Self-evaluation refers to the reflection of study. It is a very important meta-cognitive strategy because it can help students make self-improvement, erect self-confidence and regulate learning methods to make better improvement in their academic records. So after a class, it is quite helpful to carry out a survey about students’ self-evaluation which can help cultivate the awareness of self-management and promote students’ m-learning practice. According to the features and problems existing in m-learning, students may check out their performance of daily study plan or phrasal study plan. Usually when learning attitude is active, learning notes are kept clear and the content is comprehensive. Each m-learning activity will bring new learning content. Summarizing study experience, students can find out suitable learning methods and learning strategies. Through writing about learning experience, students can check out the encountered problems in order to correct them in time. Through the overall understanding of recent study, students may adjust their learning targets and plans.

5 Summary

From above, management ability has instructive and practical meaning for m-learning. M-learning in college under the instruction of self-management is workable and effective. This theory will promote education reform.

References

1. Mobile Learning Anytime Everytime [DB/OL],
<http://www.linezine.com/2.1/features/cqmmwiyp.html>
2. O'Malley, J.M., et al.: Learning strategy applications with students if English as a second language. *TSSOL Quarterly* 19, 128–136 (1985b)
3. Wang, H.: Automatic Learning among College Students. *Social Sciences Review* (5) (2007)
4. Flavell: Metacognition and cognitive monitoring: A new area of psychology inquiry. In: Nelson, T.O. (ed.) *Metacognition: Core Readings*, pp. 3–8. Allyn and Bacon, Boston (1979)
5. Brown, T.H.: Beyond constructivism: Exploring future learning paradigms. *Education Today* 2, 1–11 (2005)
6. Li, Y., Li, R.: Cultivating Meta-cognitive Ability to Teach Students How to Learn. *China Education Journal* (1), 32–35 (1999)
7. Wang, Y., Zhang, L.: How to Cultivate Students' Meta-cognitive Ability. *China Audio-Visual Education Journal* (8), 21–23 (2000)
8. Li, D.: Problems on Teach Students How to Learn. *Journal of Northwest Normal University (Social Sciences)* (1), 3–9 (1994)
9. Dickinson, L.: *Self-instruction in Language Learning*. Pergamen, Oxford (1981)
10. Pelletier, C.M.: *Successful Teaching Strategies: An Effective Instruction of Teaching Practice*. China Light Industry Press, Beijing (2002)

Improved Detection Algorithm for MIMO Wireless Communication System Based on Chase Detector

Li Liu¹, Jinkuan Wang¹, Dongmei Yan¹, and Fulai Liu²

¹ School of Information Science & Engineering, Northeastern University, 110819, Liaoning, China

² Engineering Optimization and Smart Antenna Institute, Northeastern University at Qinhuangdao, 066004, Hebei, China
liuliqhd@126.com

Abstract. In order to get better trade-off between detection performance and complexity in multiple-input multiple-output (MIMO) wireless communication system, an improved detection algorithm based on Chase detector was presented here. Original Chase detector was combined of a list detector followed by parallel banks of V-BLAST sub-detectors. There was error propagation in this algorithm and the complexity was much high. The order of symbol detection was becoming critical because of serial detection in Chase detector. So an improved detection algorithm for reducing error propagation was proposed. Sorted QR decomposition and partial Maximum likelihood detection were performed firstly to create candidate list in order to decrease error in the first step. Parallel detections were used in sub detectors with improved QRD-M to improve the bit error performance with lower complexity further more. The proposed algorithm could obtain proper trade-off between complexity and performance, and simulation experiment results show its validity.

Keywords: MIMO detection, Chase algorithm, error propagation, sorted QR decomposition, parallel detection.

1 Introduction

Multiple-input Multiple-output (MIMO) systems capacity can be increased enormous without additional the bandwidth or transmitted power in rich scattering channel [1],[2]. At the transmitter, serial data streams are converted to parallel, with each data symbol transmitted by different antenna. In order to detect transmitted symbols at receiver, symbol detection algorithm for MIMO systems have attracted much interests in recent years, and many detection techniques has been proposed such as Bell-Labs layered space-time detection (V-BLAST)[3], SD[4],[5], LR[6],[7], QRD-M[8],[9], tree search[10],[11], Chase[12], etc. Maximum likelihood (ML) detection algorithm is the optimum detection algorithm at the performance of bit-error rate (BER), but the computational complexity growing exponentially with the order of modulation and the number of transmit-antennas. Other detections have lower computation

complexity comparing with ML detection; meanwhile, the performance is also lower than ML detection. So lots of efforts have been put into the search of detection algorithms achieving ML or near-ML performance with lower complexity.

The general Chase detector for MIMO detection is a combination of a list detector and parallel banks of sub-detectors. By changing the list length, Chase detector could be regarded as the unified framework of existing detectors including ML, line detection and V-BLAST. From the detection mechanism of the Chase detector, it can be noticed that Chase detection is successive interference cancellation (SIC) detection algorithm in essence. When error symbol is produced in the first detection stage, it will lead to error spread in later sub-detectors. The importance of first detected symbol correctly should be treated carefully.

In order to suppress the error propagation and reduce the complexity in Chase detector, sorted QR decomposition algorithm according to the rule of SNR is performed firstly and partial ML detection is performed then to create candidate list without calculation of pseudo-inverses. This work can reduce the errors in the first detection step. Then parallel detection is selected as sub detectors with improved QRD-M to improve the bit error performance with lower complexity further more. The proposed algorithm can obtained proper trade off between complexity and performance, and simulation experiment results show its validity.

2 Algorithm Description

In this part, the system model and original Chase detector are presented firstly; proposed detection algorithm and simulation are given later.

2.1 System Description

In a MIMO system with N_t transmit antennas and N_r receiver antennas ($N_t \leq N_r$). The received signal complex vector $\tilde{\mathbf{y}}(t)$ can be represented as

$$\tilde{\mathbf{y}}(t) = \tilde{\mathbf{H}}(t)\tilde{\mathbf{x}}(t) + \tilde{\mathbf{n}}(t) \quad (1)$$

Where $\tilde{\mathbf{x}}(t)$ is the $N_t \times 1$ transmitted signal, with $\tilde{\mathbf{y}}(t)$ the $N_r \times 1$ received signal, and $\tilde{\mathbf{n}}(t)$ is the noise symbols. The element of $\tilde{\mathbf{H}}(t)$ represents complex channel gains between transmitter and receive antennas at the discrete time t .

The complex MIMO system model in Eq.(1) can be equivalent to a real model as follows

$$\begin{bmatrix} \text{Re}(\tilde{\mathbf{y}}) \\ \text{Im}(\tilde{\mathbf{y}}) \end{bmatrix} = \begin{bmatrix} \text{Re}(\tilde{\mathbf{H}}) & -\text{Im}(\tilde{\mathbf{H}}) \\ \text{Im}(\tilde{\mathbf{H}}) & \text{Re}(\tilde{\mathbf{H}}) \end{bmatrix} \begin{bmatrix} \text{Re}(\tilde{\mathbf{x}}) \\ \text{Im}(\tilde{\mathbf{x}}) \end{bmatrix} + \begin{bmatrix} \text{Re}(\tilde{\mathbf{n}}) \\ \text{Im}(\tilde{\mathbf{n}}) \end{bmatrix} \quad (2)$$

Definitions the new real dimensions $N_T = 2N_t$, $N_R = 2N_r$. The real symbol alphabet is now Ω , e.g., in the case of 16-QAM, $\Omega = \{-1, -3, 1, 3\}$, the equivalent real model

$$\bar{y}(t) = \bar{H}(t)x(t) + \bar{n}(t) \tag{3}$$

ML detection of the transmitted signal can be formulated as finding

$$\hat{x}_{ML} = \arg \min_{x \in \Omega^{N_T}} \{ \|\bar{y} - \bar{H}x\|_2^2 \} \tag{4}$$

When performing MLD, an exhaustive searching over the whole real alphabet has to be done, the computational complexity grows exponentially with the increasing of the number of transmit antenna and the constellation size.

2.2 Original Chase Detection Framework.

The Chase detector defines a framework for existing MIMO detection algorithms. Block diagram of Chase detector is shown in Fig.1.

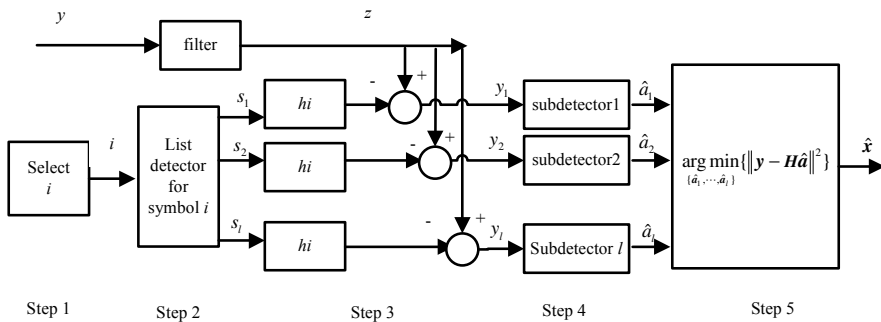


Fig. 1. Block diagram of the Chase detector

Five steps are contained in Chase detector, these steps are outlined in the framework:

Step1. Identify which symbol to be detected first according to the channel matrix H .

Step2. Generate a sorted list $L = \{s_1, \dots, s_l\}$ of candidate values for i^{th} symbol, defined as the l elements of the alphabet nearest to z_i , where $z = Gy$ is the output of linear filter.

Step3. Generate a set of l residual vectors $\{y_1, \dots, y_l\}$ by canceling the contribution to \mathbf{y} from the i^{th} symbol, assuming each candidate from the list is correct in turn, $y_j = \mathbf{y} - h_i s_j$

Step4. Assign each of $\{y_1, \dots, y_l\}$ to its own independent sub-detector to detect the remaining $N_t - 1$ symbols (all but the i^{th} symbol).

Step5. Choose the candidate \hat{a} which has the minimum mean square error as the final hard decision

$$\hat{x} = \arg \min_{\{\hat{a}_1, \dots, \hat{a}_l\}} \{\|y - H\hat{a}\|^2\} \quad (5)$$

3 Improved Chase Detection Algorithm

In Chase detector, which symbol to be detected first is critical to all the overall performance. Maximum post SNR ordering methods is proved to be the optimal detection order in V-BLAST with iterative pseudo-inverses calculating. For reducing error propagation, symbols should be detected with the optimal order. In order to select the first detected symbol correctly with lower complexity, an improved Chase detection algorithm is proposed here, the diagram is shown in Fig.2. The proposed detector is different to the original Chase detector in step1, step2 and step4.

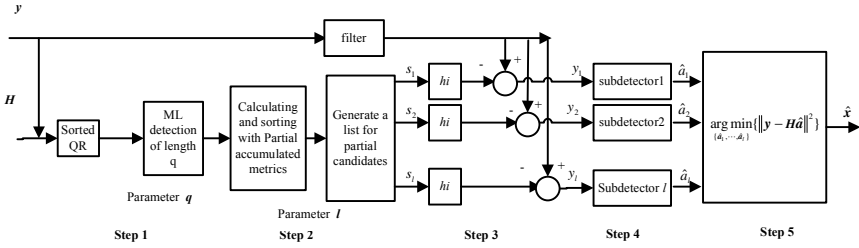


Fig. 2. Block diagram of modified Chase detectors

3.1 Sorted QR Decomposition with SNR Rules

The SNR of i^{th} layer can be described

$$SNR_i = \frac{E[|x_i|^2]}{E\{\bar{n}_i\}^2 [R^{-1}(R^{-1})^H]_{i,i}} = \frac{1}{\sigma^2 \|(R)^{-1}_{i,:}\|_2} \sim \|R_{i,:}\|_2 \quad (6)$$

where, $[\cdot]_{i,i}$ is the i^{th} diagonal component of matrix $[\cdot]$, $\|R_{i,:}\|_2$ is the 2-norm of the i^{th} row of R .

So the optimized ordered rule of maximum SNR is equivalent to the order of row 2-norm of R . In order to reduce the error produced in the first detection step, the detection should be started with the maximum SNR. Since the detection is starting from the last layers, so we can get the order of $\|R_{i,:}\|_2^2$ from minimum to maximum. With perfect knowledge of channel, sorted QR decomposition based on Householder transformation is performed firstly.

Because of the unitary of Householder transformation, it orders by columns 2-norm only once. The sorting of the columns is from small to big as the SNR orders. Based on the order rule, the column with the smallest 2-norm is performed QR decomposition in the first to set the signal with smallest SNR at the upper layers. The

estimated signal \hat{x}_1 has the minimal SNR, the detection of the left layers are influenced relatively smaller, the estimation with \hat{x}_{N_t} has the maximum SNR, the error produce in the first detected layer would be smaller, so the error propagation can be reduced.

3.2 Creation of Candidate List

In order to reduce the error produced in the first detection step, a parameter p is set in step 1 for deciding the numbers of the first detected layers. Instead of selecting only one layer, p layers are selected and partial ML detection is performed for the selected layers.

After the p layers (their layer's numbers are $N_t, N_t - 1, \dots, N_t - p + 1$, respectively) are performed ML search, the partial accumulated metrics are calculated and sorted as

$$PAM(\lambda) = \sum_{i=N_t-p+1}^{N_t} \left| y_i - \sum_{j=i}^{N_t} R_{i,j} \hat{x}_j \right|^2 \quad (7)$$

where $\lambda = 1, \dots, q^p$. Ω is the set of modulation constellation, $|\Omega| = q$ is the cardinality of modulation constellation. Without loss of generality that the symbol with lower index has smaller metrics, $PAM(1) \leq PAM(2) \leq \dots \leq PAM(\lambda)$. Now we get q^p available partial candidate paths. Because of the property of the exhaustive search algorithm, symbols with high SNR are detected at this stage to guarantee the performance with little error. Then these sequences are sorted with the partial accumulated metrics.

To reduce error propagation further, in step2, another parameter l is set for defining candidate list length. The list of partial candidates is generated in which l candidate signal sequences with smaller partial accumulated metrics are selected in step2. More candidates are selected to look for more possible solution.

3.3 Sub Detectors with Parallel QRD-M

In step4 QRD-M algorithm is selected as the sub detector to search the left layers. Difference to original QRD-M, here they are performed in parallel separately with the partial candidate sequence in the list created in step2. There is ordered operation only in each sub detector, and there is no sorted operation between sub detectors which can reduce sorting so the calculating complexity reduced obviously.

With each output from step4, the candidate signal sequence with the smallest accumulated metrics is decision, quantization, and re-arranged as the order as the transmitted signal.

$$\hat{x} = \arg \min_x (\|y - R\hat{x}\|^2) = \arg \min_x \left(\sum_{i=1}^{N_t} \left| y_i - \sum_{j=i}^{N_t} R_{i,j} \hat{x}_j \right|^2 \right) \quad (8)$$

The maximum likelihood detection of the transmitted signal can be formulated as finding

$$\hat{x} = \arg \min_{x \in \{\Omega\}^{N_T}} \{\|y - Hx\|^2\} \quad (9)$$

3.4 Simulations

Channel is assumed to be Rayleigh flat-faded and no correlation between sub-channels. And the receiver knows the knowledge of channel perfectly. Detection performance of proposed Chase detector compares with VBLAST and ML detection. Improved Chase algorithm with parameters as $p = 2, l = 2, M = 2$ (line 1) and $p = 2, l = 4, M = 4$ (line 2). MIMO system is consisted with $N_T = 4, N_R = 4$ and 16-QAM modulation. The BER curves are shown in Fig.3.

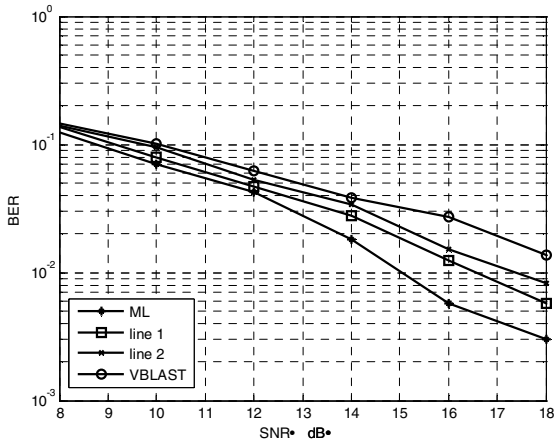


Fig. 3. BER comparison of detection performance

From the comparison results, it can be found that the performance of the proposed algorithm outperform V-BLAST as p and l increases. The performance is near that of ML detection with lower calculation complexity. Even only increasing the list length l , the BER performance can be improved significantly.

4 Conclusions

In order to reduce the error producing at the first detection layers, sorted QR algorithm was used here. The detection order was sorted as maximum SNR without calculating pseudo-inverses of channel. An improved Chase detector was proposed, which different to original Chase detector by setting two adjustable parameters. Error

propagation is decreased in early detection layers. Trade-off of the complexity and performance would be obtained properly by modifying the number of maximum likelihood detection layers and the list length. Simulation results show that the presented algorithm is superior to VBLAST algorithm and can achieve performance of MLD with lower computational complexity.

Acknowledgments. This paper has been supported by the National Natural Science Foundation of China under Grant No. 60904035 and 61004052, and by Directive Plan of Science Research from the Bureau of Education of Hebei Province, China, under Grant No. Z2009105.

References

1. Foschini, G.J., Gans, M.: On limits of wireless communications in a fading environment when using multiple antennas. *Wireless Personal Communications* 6(3), 311–335 (1998)
2. Telatar, E.: Capacity of multi-antenna Gaussian channel. *Europ. Trans. Telecommun.* 10, 585–595 (1999)
3. Foschini, G.J.: Layered space-time architecture for wireless communication in a fading environment when using multi-element antennas. *Bell Labs. Tech. J.* 1, 41–59 (1996)
4. Zheng, C.W., Chu, X.Z., McAllister, J., Woods, R.: Real-Valued Fixed-Complexity Sphere Decoder for High Dimensional QAM-MIMO Systems. *IEEE Transactions on Signal Processing* 59, 4493–4499 (2011)
5. Seethaler, D., Bolcskei, H.: Performance and Complexity Analysis of Infinity-Norm Sphere-Decoding. *IEEE Transactions on Information Theory* 56, 1085–1105 (2010)
6. Zhang, W., Qiao, S.Z., Wei, Y.M.: A Diagonal Lattice Reduction Algorithm for MIMO Detection. *IEEE Signal Processing Letters* 19, 311–314 (2012)
7. Chen, C.E., Sheen, W.H.: A New Lattice Reduction Algorithm for LR-Aided MIMO Linear Detection. *IEEE Transactions on Wireless Communications* 10, 2417–2422 (2011)
8. Li, W., Cheng, S.X., Wang, H.F.: An Improved QRD-M Algorithm in MIMO Communications. In: *Global Telecommunications Conference*, pp. 4380–4384 (2007)
9. Kim, B.S., Choi, K.: A Very Low Complexity QRD-M Algorithm Based on Limited Tree Search for MIMO Systems. In: *Vehicular Technology Conference*, pp. 1246–1250 (2008)
10. Jia, Y.G., Andrieu, C., Piechocki, R.J., Sandell, M.: Depth-First and Breadth-First Search Based Multilevel SGA Algorithms for Near Optimal Symbol Detection in MIMO Systems. *IEEE Trans. on Wireless Communications* 7, 1052–1061 (2008)
11. Kang, H.G., Song, I., Oh, J., Lee, J., Yoon, S.: Breadth-First Signal Decoder: A Novel Maximum-Likelihood Scheme for Multi-Input–Multi-Output Systems. *IEEE Trans. on Vehicular Technology* 57, 1576–1584 (2008)
12. Waters, D.W., Barry, J.R.: The Chase Family of Detection Algorithm for Multiple-Input Multiple-Output Channel. *IEEE Trans. on Signal Processing* 56, 739–747 (2008)

Improving Recommendation Performance through Ontology-Based Semantic Similarity

Mingxin Gan^{1,*}, Xue Dou¹, and Rui Jiang^{2,*}

¹ School of Economics and Management, University of Science and Technology Beijing, Beijing, 100083, China

² Bioinformatics Division, TNLIST/Department of Automation, Tsinghua University, Beijing 100084, China

ganmx@ustb.edu.cn, ruijiang@tsinghua.edu.cn

Abstract. Making personalized recommendation according to preferences of users is of great importance in recommender systems. Currently most book recommender systems take advantage of relational databases for the representation of knowledge and depend on historical data for the calculation of relationships between books. This scheme, though having been widely used in existing methods based on the collaborative filtering strategy, overlooks intrinsic semantic relationships between books. To overcome this limitation, we propose a novel approach called COSEY (COLlaborative filtering based on item SEMantic similarity) to achieve personalized recommendation of books. We derive semantic similarities between books based on semantic similarities between concepts in an ontology that describes categories of books using our previously proposed method DOPCA, and we incorporate such similarities between books into the item-based collaborative filtering strategy to achieve personalized recommendation. We validate the proposed COSEY approach through comprehensive experiments and show the superior performance of this approach over existing methods in the recommendation of books.

Keywords: recommender systems, ontology, semantic similarity, degrees of overlap in paths, depth of the lowest common ancestor node.

1 Introduction

The amount of resources that spread over the world wide web has been growing exponentially in the past decade, due to the explosive development of the internet technology and the accompanying informatization that pervades the world. However, our capability of capturing the right information does not grow in a compatible speed [1]. To fill out this gap, recommender systems have been developed for helping people to achieve a quick and accurate pinpointing of information in the vast ocean of resources [2]. For example, a collaborative filtering system infers interests of a user by utilizing preferences of other users [3-4]. A content-based approach relies on characteristics of resources to make recommendations [5-6]. A knowledge-based system

* Corresponding authors.

suggests resources based on the inference about the need of users [7]. Moreover, these methods may also be integrated to achieve a more accurate recommendation [8].

However, these methods suffer from such common problems as the introduction of new users and new resources (cold start), the modeling of preferences of users, the extraction of suitable features of resources, and the intrinsic sparseness of relations between users and resources. It has been pointed out that a high quality model of user interests can greatly promote the performance of a recommender system that uses the collaborative filtering scheme, and such a model should have good expansibility. Nevertheless, traditional forms of knowledge description generally lack the capability of portraying interests of users with veracity and scalability.

On the other hand, successful applications of semantic similarity have been found in not only such fields in computer science as word sense disambiguation and web service discovery, but also such fields in life sciences as the derivation of functional similarity between genes [9]. Particularly, it has been shown that the incorporation of a domain ontology can effectively improve the performance of a recommender system through the inherent semantic description mechanism and the structure supporting logical reasoning of the ontology. With the introduction of a domain ontology, items can be categorized into a unified concept hierarchy, and similarities between items can be calculated through semantic relationships of concepts [11]. However, most existing methods can only obtain semantic similarity between concepts with single inheritance relation rather than multi-inheritances relation. Hence, it is not trivial to develop an effective method to calculate semantic similarity based on an ontology.

An ontology is typically represented as a directed acyclic graph (DAG), in which nodes correspond to concepts, and edges denote relationships between the concepts. There have been three categories of methods for evaluating relatedness of concepts: methods based on the structure of an ontology [12-14], methods relying on information content of concepts [15], and methods utilizing multiple properties of an ontology in a hybrid manner. We have previously proposed a method called DOPCA that relies on the structure of an ontology to calculate semantic similarity between concepts. This method combines two similarity measures, the degrees of overlap in paths (DOP) and the depth of the lowest common ancestor node (DLCA), and uses their weighted summation to quantify the relatedness of concepts. This method is capable of overcoming the limitation of existing methods that overlook the existence of multiple lowest common ancestor nodes, and is flexible when applied to ontologies of different domains.

With these understandings, we pursue the goal of applying our DOPCA method to calculate semantic similarities between books, and then use such a similarity measure to enhance the item-based collaborative filtering algorithm. More specifically, we propose a novel approach called COSEY (Collaborative Filtering based on Item Semantic Similarity) to make prediction of book ratings for users. We calculate semantic similarities between concepts in a book ontology using DOPCA, and we derive semantic similarity between books relying on associations between books and concepts. We perform comprehensive validation experiments on BookCrossing data and show the superior of our approach over existing methods.

2 Ontology Based Concept Semantic Similarity

We have previously proposed a method called DOPCA that relies on the structure of an ontology to calculate semantic similarity. This method uses the weighted summation of the degrees of overlap in paths (DOP) and the depth of the lowest common ancestor node (DLCA) to quantify the relatedness of terms in an ontology. Here we briefly introduce this method as follows.

Assumptions and definitions. An ontology has a directed acyclic graph (DAG) structure that is represented as $G = (V, E)$, where V is a set of vertices denoting terms in the ontology and E a set of edges denoting semantic relationships between the terms. In the case that an ontology has two or more types of semantic relationships, we distinguish the set of edges into two or more sub-sets, one for each type of relationship. For example, the gene ontology (GO) defines two relationships, “is-a” and “part-of”, between concepts, we then represent GO as $G = (V, E_i, E_p)$, where E_i is the set of “is-a” relationships, and E_p the set of “part-of” relationships.

Given a term A in an ontology G , we define $GA = (VA, EA)$ as a subgraph that includes A and all its ancestors. Therefore VA is the set of terms including A and its ancestors, and EA is the set of edges connecting the terms in VA . Given two terms A and B , we define the union subgraph of their corresponding subgraphs GA and GB as a subgraph $GA \cup B = (VA \cup B, EA \cup B)$, where $VA \cup B = VA \cup VB$ and $EA \cup B = EA \cup EB$. Similarly, given two terms A and B , we define their intersection subgraph as $GA \cap B = (VA \cap B, EA \cap B)$, where $VA \cap B = VA \cap VB$ and $EA \cap B = EA \cap EB$.

Following the literature [15], we assume the properties for a semantic similarity measure that relies on an ontology. With these assumptions and definitions, we propose the following degrees of overlap in paths (DOP) and the depth of the lowest common ancestor (DLCA) similarity measures.

Degrees of overlap in paths. Lin calculated semantic similarity between two terms A and B as the ratio between the amount of information needed to state their commonality and the information required to fully describe them, Considering the DAG structure, we simplify the definition proposed by Lin and propose the following degrees of overlap in paths (DOP) between the terms A and B as

$$sim_{DOP}(A, B) = \frac{w_v \times |V_{A \cap B}| + w_e \times |E_{A \cap B}|}{w_v \times |V_{A \cup B}| + w_e \times |E_{A \cup B}|} \quad (1)$$

where $|A|$ stands for the cardinality of A , w_v the weight of the vertices, and w_e the weight of the edges. Obviously, this formula is equivalent to

$$sim_{DOP}(A, B) = \frac{|V_{A \cap B}| + r_e \times |E_{A \cap B}|}{|V_{A \cup B}| + r_e \times |E_{A \cup B}|} \quad (2)$$

with $r_e = w_e / w_v$ being the ratio of the weights (relative weight). In the case that the ontology contains two or more types of semantic relationships, we can easily extend the above formula by incorporating multiple ratios for the corresponding multiple weights. For example, for the gene ontology, which has an “is-a” relationship and a “part-of” relationship, the formula should be

$$sim_{DOP}(A, B) = \frac{|V_{A \cap B}| + r_{is-a} \times |E_{A \cap B}^{is-a}| + r_{part-of} \times |E_{A \cap B}^{part-of}|}{|V_{A \cup B}| + r_{is-a} \times |E_{A \cup B}^{is-a}| + r_{part-of} \times |E_{A \cup B}^{part-of}|} \quad (3)$$

where r_{is-a} is the relative weight for is-a edges, and $r_{part-of}$ is the relative weight for part-of edges. Note that the assignment of the relative weights is critical in the above formulae. In particular, in the case that terms A and B have more than one LCA node, the paths from the root to these LCA nodes will all be calculated in our method, such that the problem of multiple LCA nodes can be solved.

Depth of the lowest common ancestor node. According to the property of semantic similarity, that is, the semantic similarity increases when the LCA node of two terms A and B becomes deeper in the ontology graph, the LCA node plays an important role in the calculation of the semantic similarity. Generally, if two terms share an ancestor that is deep in the ontology structure, their semantic similarity should be larger than those whose common ancestor locates shallow in the ontology structure. This is reasonable because the terms in deeper locations typically represent more concrete concepts and thus contribute more to the semantic similarity.

With this understanding, we adopt the exponential function that is used by Zhang et al [14] to calculate the semantic similarity that is represented by the depth of the LCA node, as

$$sim_{DLCA}(A, B) = \exp(-\lambda / D_{LCA}) \quad (4)$$

where λ is a free parameter between 0 and 1. This function ensures higher semantic similarity for terms whose LCA node locates deeper in the ontology.

The DOPCA model. With the above degrees of overlap in path (DOP) method and the depth of the lowest common ancestor node (DLCA) approach, we are able to calculate two semantic similarity measures. We then propose to combine these two quantities to obtain a single semantic similarity measure, as

$$sim_{DLCA}(A, B) = w_{DOP} Sim_{DOP} + w_{DLCA} Sim_{DLCA} \quad (5)$$

In this formula, Sim_{DOP} represents the semantic similarity of two terms that is calculated from not only the commonality of the two terms but also the locations of the two terms in the entire ontology graph, addressing the problem of multiple LCA nodes that most existing methods have ignored. On the other hand, Sim_{DLCA} represents our general understanding of semantic similarity and makes sure that the properties of a similarity measure are satisfied.

3 Collaborative Filtering Based on Item Semantic Similarity

Based on a constructed book ontology and the DOPCA method, we proposed a new method named COSEY (Collaborative Filtering based on Item Semantic Similarity) to make prediction of book ratings for users in a recommender system.

We first define symbols to describe our approach as follows. Let u be the number of users and b the number of books in our recommender system. Let $u = (u_1, \dots, u_u)$ be all users and $b = (b_1, \dots, b_b)$ all books. Let $r_u = (r_{u1}, \dots, r_{ub})$ be the ratings that the

user u gives to all books. Let $t_b = (t_{b1}, \dots, t_{bc_b})$ be the set of all concepts that are associated with the book b .

We then summarize the proposed method as follows. First, we calculate the semantic similarity between every two terms t and s in the book ontology by directly applying the above DOPCA method, as

$$Sim(t, s) = w_{DOP}Sim_{DOP}(t, s) + w_{DLCA}Sim_{DLCA}(t, s) \quad (6)$$

where w_{DOP} determines the contribution of the DOP measure and w_{DLCA} determines the contribution of the DLCA approach, respectively. We require that $w_{DOP} + w_{DLCA} = 1$. Second, on the basis of semantic similarities between concepts in the ontology, we calculate semantic similarity between a concept t and a set of concepts S as

$$Sim(t, S) = \max_{s \in S} (Sim(t, s)) \quad (7)$$

that is, the maximum similarity between the concept t and every concept s in the set S . Third, we calculate semantic similarity between a set of concept T and a set of concepts S as

$$Sim(T, S) = \frac{\sum_{t \in T} Sim(t, S) \sum_{s \in S} Sim(T, s)}{|T| + |S|} \quad (8)$$

that is, the average of similarity between every concept in a set and the other set of concepts. Since a book is just a set of concepts, the semantic similarity between two books b and c can be simply calculated as the semantic similarity of two sets of concept B and C , where B is the set of concepts corresponding to book b , and C is the set of concepts corresponding to book c . Fourth, we calculate the rating that a user u gives to a book b as

$$\hat{r}_{ub} = \frac{\sum_{b' \in b} Sim(b, b') \times r_{ub'}}{\sum_{b' \in b} Sim(b, b')} \quad (9)$$

that is simply the weighted arithmetic average of all known rating that the user u gives to books, with the weight being the normalized similarity of all other books to the book b . Finally, we sort the rating of all books in descending order and recommend the top ranking books to the user.

4 Experimental Evaluation

Construction of book ontology. The construction of a complete ontology for books is too ambitious to be a feasible work, due to the high degree of specificity that leads to a very large number of concepts. Therefore, we focus only on employing an ontology that can be easily extracted from some of the available classification standards for books, instead of building an ontology covering all possible types of book items. For this purpose, we chose Amazon.com as the main input for constructing the domain

ontology for books. More specifically, we have defined new properties and classes to accommodate some missing features. Along with the multiple hierarchies of classes that serve to categorize the books and their attributes, the ontology contains labeled properties joining each item to its attributes.

Validation sets. We choose some of the selected book ontology to verify our method. Our data set consists of a subset of the book ratings that were collected by Ziegler in a 4-week crawl from the Book-Crossing community. For each distinct ISBN, we mined Amazon.com’s Book Taxonomy and collected the category, title, URL, and editorial reviews. Only the explicit ratings, expressed on a scale from 3-10, are taken into account in our experiments. The data set was converted into a user-item matrix belonging to those users with 20 or more ratings. For evaluation, we used 5-fold cross-validation. For each fold, 80% of the book ratings were included in the training set, which was utilized to compute similarities among users. The remaining 20% of the ratings were included in the test set, which was used for predicting ratings.

Experimental metrics. To measure prediction accuracy, we rely on a commonly used metric Mean Absolute Error (MAE), which measures the average absolute deviation between a predicted rating and the user’s true rating:

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |q_{ui} - r_{yi}| \quad (10)$$

where N is the total number of ratings over all users, q_{ui} is the predicted rating for user u on item i , and r_{ui} is the actual rating.

Results and analysis. The proposed COSEY method relies on DOPCA to calculate semantic similarity of concepts in the constructed ontology. To show the effectiveness of DOPCA, we also implement the method of Wang for calculating semantic similarity of concepts. We denote the method of COSEY with DOPCA as COSEY-DOPCA) and that use COSEY with Wang’s method as COSEY-Wang.

Table 1. MAE Results of COSEY-DOPCA vs. COSEY-Wang

Sample size	MAE (COSEY-Wang)	MAE(COSEY-DOPCA)
50	0.682	0.617
60	0.651	0.601
70	0.654	0.572
80	0.633	0.589
90	0.621	0.536
100	0.590	0.480

We perform a more comprehensive evaluation of the proposed method by calculating the criterion of MAE for a number of books selected at random from the 5-fold cross-validation. In detail, for each method, we computed a predicted rating for each item that was rated by that user’s neighbors, excluding the items that were rated by the user. We then generate a recommendation list for each specific user by sorting the books in descending order with respect to the predicted rating for each item, and calculate the MAE across the predicted ratings created by each method along with the increasing number of sample size from 50 to 100. As shown in Table 1, the MAE

value for COSEY-Wang is 0.682, and the value for COSEY-DOPCA is 0.617, on the sample size of 50. We also notice that on whatever size of the sample, the MAE value of COSEY-Wang is higher than that of COSEY-DOPCA. To achieve the trend of MAE result more clearly, we draw a trend line in Figure 1. Obviously, to each method, along with the increasing of the sample size, the MAE value is decreasing. It means the MAE value is effected by the amount of the sample size. Moreover, the fact that the MAE value of COSEY-DOPCA is lower than that of COSEY-Wang following the trend indicates that ratings predicted through COSEY method based on DOPCA is more accurate than that based on Wang.

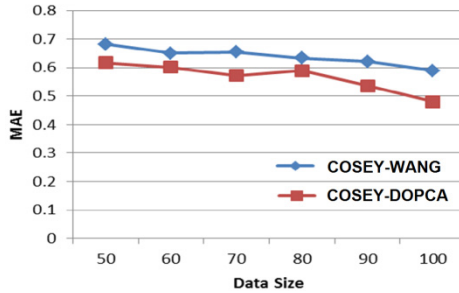


Fig. 1. MAE Result for COSEY

5 Conclusions and Discussion

In this paper, we propose a novel approach called COSEY to achieve personalized recommendation for books. A hallmark of COSEY is the incorporation of an ontology to the description of books, and thus relationships between books can be quantified through the calculation of similarities between concepts in the ontology. Obviously, this scheme overcomes the typical drawback of existing method of item-based collaborative filtering, which utilizes historical to calculate similarities between items (books) and intrinsically has the cold-start issue.

The second main characteristic of our method is the use of a method called DOPCA to quantify similarities of concepts in an ontology. DOPCA can successfully solve the problem of describing multi-inheritances of concepts semantic in a domain ontology through overlap degrees of paths in an ontology. Hence, items belong to more than one class can be considered under the domain ontology based on DOPCA. This is largely supporting to solve the problem of similarity calculating between item sets which belong to multi-classes.

We have demonstrated the effectiveness of the proposed COSEY framework and the DOPCA method through comprehensive experiments. We have implemented two version of COSEY, one based on DOPCA and another based on an existing method for calculating semantic similarity of concepts in an ontology. We have shown that COSEY with DOPCA can achieve higher prediction accuracy for ratings of books.

At present, the application of ontology has been promoted in some recommender systems such as digital library. However, since researches on ontology based concept

semantic similarity are still young while the mainstream recommendation methods remains collaborative filtering approaches based on users' similarity, it still remains a long way to go to realize real semantic-driven recommendation.

Acknowledgments. This work was partly supported by the National Natural Science Foundation of China under Grants No. 71101010, 61175002 and 71172169, the Fundamental Research Funds for the Central Universities under Grant No. FRF-BR-11-019A, and TNLIST Cross Discipline Foundation.

References

1. Adomavicius, G., Tuzhilin, A.: Toward the Next Generation of Recommender Systems: A Survey of the State-of-the-Art and Possible Extensions. *IEEE Transactions on Knowledge and Data Engineering* 17(6), 734–747 (2005)
2. Resnick, P., Varian, H.R.: Recommender Systems. *Communications of the ACM* 40(3), 56–58 (1997)
3. Liu, R., Jia, C., Zhou, T., Sun, D., et al.: Personal Recommendation via Modified Collaborative Filtering. *Physica A* 338, 462–468 (2009)
4. Chen, Y., Cheng, L.: A Novel Collaborative Filtering Approach for Recommending Ranked Items. *Expert Systems with Applications* 34(4), 2396–2405 (2008)
5. Belkin, N., Croft, B.: Information Filtering and Information Retrieval. *Communications of the ACM* 35(12), 29–37 (1992)
6. Balabanovic, M., Shoham, Y.: Fab: Content based Collaborative Recommendation. *Communications of the ACM* 40(3), 66–72 (1997)
7. Prasad, B.: A Knowledge-based Product Recommendation System for e-Commerce. *International Journal of Intelligent Information and Database Systems* 1(1), 18–36 (2007)
8. Burke, R.: Hybrid Recommender Systems: Survey and Experiments. *User Modeling and User-Adapted Interaction* 12(4), 331–370 (2002)
9. Pesquita, C., Faria, D., Falcao, A.O., et al.: Semantic Similarity in Biomedical Ontologies. *PLoS Computational Biology* 5(7) (2009)
10. Blanco-Fernandez, Y., Lopez-Nores, M., Pazos-Arias, J.J., Garcia-Duque, J.: An Improvement for Semantics-based Recommender Systems Grounded on Attaching Temporal Information to Ontologies and User Profiles. *Engineering Applications of Artificial Intelligence* 24(8), 1385–1397 (2011)
11. Deshpande, M., Karypis, G.: Item-based Top-N Recommendation Algorithms. *ACM Transactions on Information Systems* 22(1), 143–177 (2004)
12. Rada, R., Mili, H., Bicknell, E., Blettner, M., et al.: Development and Application of a Metric on Semantic Nets. *IEEE Transactions on Systems, Man, and Cybernetics* 19(1), 17–30 (1989)
13. Wang, J., Du, Z., Payattakool, R., Yu, P.S., et al.: A New Method to Measure the Semantic Similarity of GO Terms. *Bioinformatics* 23(10), 1274–1281 (2007)
14. Zhang, S., Shang, X., Wang, M., et al.: A New Measure Based on Gene Ontology for semantic similarity of Genes. In: *WASE International Conference on Information Engineering*, pp. 85–88. IEEE Press, Los Alamitos (2010)
15. Resnik, P.: Using Information Content to Evaluate Semantic Similarity in a Taxonomy. In: *14th International Joint Conference on Artificial Intelligence* (1995)

Formal Modeling and Model Checking Analysis of the Wishbone System-on-Chip Bus Protocol

Ricai Luo and Hua Tan

Dep. of Computer & Information Science, Hechi University, YiZhou, GuangXi, China
luoricai@163.com, 164666827@qq.com

Abstract. The Wishbone System-on-Chip bus protocol, which is developed by the Silicore Corporation, in connection with its characteristics and complexity, it is verified by the model checking approach. Firstly, the communication model of IP cores is created, and a FSM modeling approach for the model is proposed. Secondly, the non-starvation and fairness properties are specified using the computation tree logic. Finally these properties are verified against the model with the help of the model checking tool SMV. The result shows that there is a bus starvation scenario which will be caused by the unfairness of arbiter. This research demonstrates that there are some flaws with the specification of Wishbone System-on-Chip bus protocol. It also reflects that the arbitration mechanism is prone to flaw and therefore the formal modeling and verification is necessary.

Keywords: Wishbone Bus Protocol, System-On-chip, Finite State Machine, Model Checking, Computation Tree Logic.

1 Introduction

With the increasing scale of integrated circuit design, and the increasing complexity of functions, system-on-chip (SOC) design emerges as a new generation of IC design technology to meet this market demand. In the context of this design technology, on-chip system bus and its protocol has become the key to application.

Nowadays, many companies have introduced and developed a number of on-chip bus standards, such as ARM Company's AMBA [1] bus, IBM Company's Core Connect[2] bus, completely open and completely free OCP [3] bus and Wishbone [4] bus. It is a matter of concern whether these buses and their protocols have possible or potential defects.

The checking of on-chip bus protocol has become the concern of the researchers. Abhik Roychoudhury et al have verified AMBA bus protocol widely popular in industry [5], it is found through checking that potential bus starvation situation exists due to the incomplete specifications of the protocol. Pankaj Chauhan et al found two potential problems in the protocol specifications during their checking of the PCI bus [6]. In reference [7], IBM Company's Core Connect bus is verified, and defects are found as the results of that checking; in reference [8], the compliance testing research on on-chip bus of SOC is carried out. These studies reflect that the on-chip bus protocol has some possible defects to a certain extent. Though, in practical

applications, it does not seem to have any problem, this potential possibility of emerging problem may not be excluded. It may cause huge economic losses due to this potential error, such as the floating-point arithmetic error in the Pentium chip produced by Intel in 1994, although the probability of occurrence of this error is one out of hundreds of millions, the error has caused huge economic losses of \$ 475 million to Intel Corporation. For bus protocol, its good or poor design does not affect the efficiency of data transmission between the IP cores and the reuse flexibility of the IP core, but also affect the stability and reliability of the entire SOC chip.

Wishbone on-chip bus is also an on-chip bus widely used in SOC design based on IP reuse. Based on the characteristics and complexity of Wishbone bus, this paper introduces a modeling method, and carries out the confirmatory analysis through using model checking technology.

2 Introduction of WISHBONE Bus Protocol

Wishbone bus is developed by Silicore Company, and now the maintenance of it is performed by Open-Cores organization. Since it is completely free, and it is an on-chip bus developed according to source code, it has strong market competitiveness and influence compared to other buses and has good prospects. It is also a bus architecture currently widely used in SOC design based on IP reuse, and its bus system architecture is shown in Fig. 1.

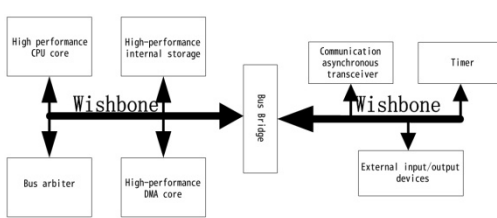


Fig. 1. Wishbone bus system architecture

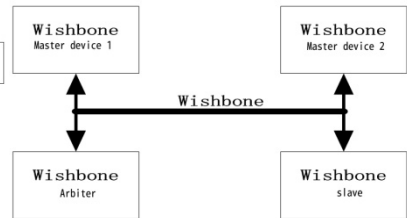


Fig. 2. Shared bus mode

The difference between Wishbone bus and other buses is that it supports four different connection modes: point to point, data flow, shared bus and cross interconnection. Interconnection structure has variability, which greatly increases the flexibility of the IP core interconnection. Where, the shared bus interconnection model is a typical connection mode, and compared to other connection modes, its implementation is simpler, and it uses less logic cells and allocation of resources, applicable for sharing a bus for multiple IP cores to simultaneously connect multiple master-slave equipments in order to increase throughput. The system block diagram is shown in Fig. 2.

According to the definition of the Wishbone bus protocol, the IP core module connected to the bus distinguishes the signal by Master / Slave interface, the module with Master interface is the initiator of the bus operation, and the module with Slave interface is the responder of Master interface module.

Wishbone bus has three unique transaction transmission modes: single-byte reading, single-byte writing and single-byte block transferring. The single-byte writing operation timing diagram in Wishbone bus protocol specifications is shown in Fig. 3.

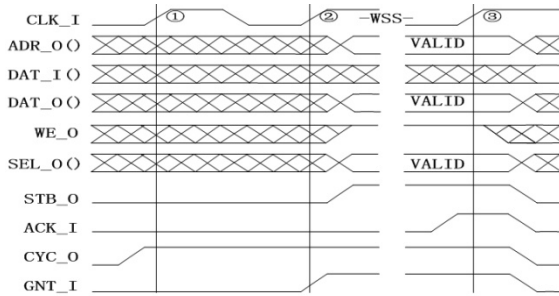


Fig. 3. Single-byte writing operation timing diagram

The CLK_I in the Figure is clock signal, when the Master interface module performs the single-byte writing operation to Slave interface module, CYC_O signal is sent from Master port when CLK_I rises for triggering for the first time, and is set as 1 in order to initialize the transaction process, while it sends request to access bus operation to the bus arbitration; if bus arbiter detects the CYC_O signal set as 1, and send GNT_I signal and sets it as 1, it indicates that Master interface module has the right to access to bus operation, and then Master port sends STB_O signal to the targeted Slave port and sets it as 1 to show to Slave port that it is in an active state, also sends a signal for requesting to write data, and is ready to write the transaction bus cycle; if Slavs port detects STB_O and STB_O signal is set as 1, Slavs port will make the action of writing data, after writing data is completed, it will response to ACK_I signal and set the signal as 1, otherwise, it will set it as 0; if the processing time is delayed, ACK_I signal will not be set as 1 until writing data is completed; Master port will detect ACK_I to determine whether Slave port has successfully completed the operation of writing data, if completed, it will set CYC_O and STB_O signals as 0 to end the transaction, and the bus arbiter also sets GNT_I signal as 0. The single-byte writing data operation is completed.

In Shared bus interconnection mode, multiple master and slave devices share a bus, and arbitration mechanism is added, where, the arbiter determines the right to operate of Master to the bus. In the single-byte writing operation, the latch signal LOCK and address type signal TGA and other signals are included in the main unit interface module except for the Master interface signal shown in Fig. 4 . The checking analysis for Wishbone bus protocol in this paper is based on a shared bus interconnect model.

3 Formal Modeling of Protocol

Formal modeling analysis has been applied to the modeling of complex system. For a system, it should be determined which formal method will be used for modeling according to its characteristics to develop an effective formal model.

This section analyzes the modeling of bus protocol by three steps according to its characteristics and complexity of Wishbone bus protocol.

A. IP interaction model of the protocol

According to the specifications of Wishbone bus protocol, this section establishes IP interaction model interconnected in a manner of a shared bus as shown in Fig. 4.

This model consists of two master devices, one slave device and one arbitration device, and all these devices are connected in a manner of a shared bus. Based on this model, when Master1 will access the slave, it must send bus operation request to the arbiter, if Arbiter detects that no other master device occupies the bus, Master1 is allowed to operate the bus, if it is detected that other master device is occupying the bus, the request of Master1 will be rejected. Only when Master1 gets the right to operate the bus, it can access to Slave. Similarly, same process is required if Master2 will access Slave.

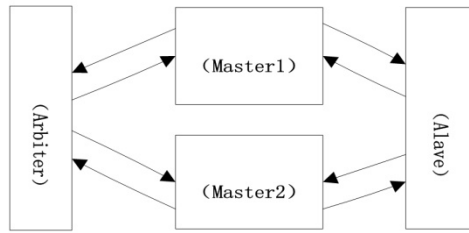


Fig. 4. IP interaction model of the protocol

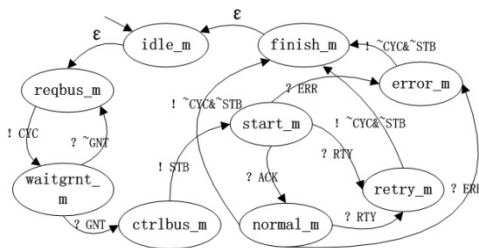


Fig. 5. FSM state transition diagram of Master

B. Main body modeling of the protocol

The specifications of Wishbone bus protocol are described based on the sequence diagram and the way of natural language. Based on its characteristics and complexity, in order to facilitate the analysis for Wishbone bus protocol, this paper uses formal method of finite state machine to specify its specifications. On the basis of analysis of the required bus signals and status, this paper describes the finite state machine used for the behaviors of different subjects involved in the implementation of the protocol according to the IP interaction model, the“~” marker shown in the finite state transition diagram indicates low level signal and other markers indicate high level signal. The finite state machine transition processes for Master 1 and Master 2 are the same, so they are described in the same FSM state transition diagram as shown in Fig. 5.

Before the starting of data transmission, Master is in the idle_m state, which is the initial state. When the Master is ready to access to the Slave, the Master will automatically skip to the reqbus_m state, and send the CYC signal of accessing to the bus to enter into waitgrnt_m state; if the Master receives ~ GNT signal, it will return back to reqbus_m state; if receiving GNT signal, it will enter into ctrlbus_m state; if the Master getting the right to operate the bus accesses the Slave, it will send STB signal, and the Master enter into start_m state, which indicates the starting of data transmission; after receiving the response signal of ACK from the Slave, if the Master enters into normal_m state, it indicates that the normal data transmission may be started; if receiving the response signal of RTY from the Slave , the Master will enter into retry_m state; after receiving a signal of ERR , If the Master will enter into error_m state; when the Master is in normal_m state, if receiving the ERR signal, the Master will enter into the error state; if receiving the the RTY signal, the Master will enter into retry_m state. When it is in normal_m state, retry_m, state and error_m state, the main unit of Master will finish the bus operation after sending~CYC and ~ STB signals, and finally get back to idle_m state automatically.

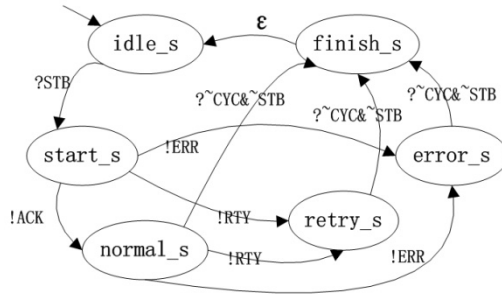


Fig. 6. FSM state transition diagram of slave

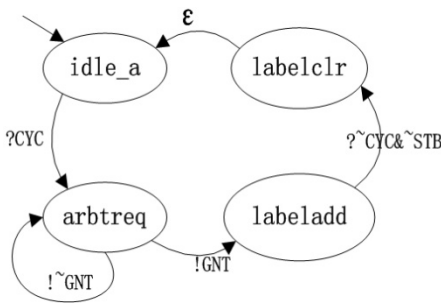


Fig. 7. FSM state transition diagram of Arbiter

Before the starting of data transmission, Master is in the idle_m state, which is the initial state. When the Master is ready to access to the Slave, the Master will automatically skip to the reqbus_m state, and send the CYC signal of accessing to the bus to enter into waitgrnt_m state; if the Master receives ~ GNT signal, it will return

back to reqbus_m state; if receiving GNT signal, it will enter into ctrlbus_m state; if the Master getting the right to operate the bus accesses the Slave, it will send STB signal, and the Master enter into start_m state, which indicates the starting of data transmission; after receiving the response signal of ACK from the Slave, if the Master enters into normal_m state, it indicates that the normal data transmission may be started; if receiving the response signal of RTY from the Slave, the Master will enter into retry_m state; after receiving a signal of ERR, If the Master will enter into error_m state; when the Master is in normal_m state, if receiving the ERR signal, the Master will enter into the error state; if receiving the the RTY signal, the Master will enter into retry_m state. When it is in normal_m state, retry_m, state and error_m state, the main unit of Master will finish the bus operation after sending \sim CYC and \sim STB signals, and finally get back to idle_m state automatically.

The state transition diagram of the Slave is shown in Fig. 6. From the starting, the Slave is in idle_s state; when receiving the STB signal from the Master, the Slave will enter into the start_s state to prepare the data transmission; if sending the ACK signal, the Slave will enter into normal_s state, which indicates the normal cycle of the bus; if sending the RTY signal, the Slave will enter into retry_s state, which requires the retransmission of the bus cycle; if sending the ERR signal, the Slave will enter into error_s state, which indicates that error happens to the bus cycle. When the Slave is in the normal_s state, if sending the ERR signal, it will enter into error_s state, which indicates that error happens to the operation, if sending the RTY signal, it will enter into retry_s state. The Slave will finish the bus operation after sending \sim CYC and \sim STB signals, and finally get back to idle_s state automatically.

The state transition diagram of the bus Arbiter is shown in Fig. 7. The bus Arbiter will enter into arbtreq state when detecting the request signal of CYC from the Master in the process of query; if the bus has been already occupied by other Masters, the bus Arbiter will send \sim GNT signal, which indicates that the operation is not allowed, and the bus Arbiter will remain at arbtreq state. If there is no other Master occupying the bus, the bus Arbiter will send GNT signal, and then enter into labeladd state to label and record the Master. The Arbiter will clear the label after receiving \sim CYC and \sim STB signals, and then enter into labelclr state to release the bus lock protection, and finally get back to idle_a automatically.

C. CTL description of system property

When specifying the formal specifications of the system property, generally, the computation tree logic [9] is used for formal description. In this description, the behavioral properties to be verified in the system model are expressed in CTL formulas. The CTL formula consists of atomic propositions (describing the basic elements of a state), logic connectors (\wedge, \vee, \neg) and the modal operator. The modal operator is divided into two parts: one part is the path quantifier, including A (Always, all future paths) and E (Exists, at least a path); another part is the modal operator, including G (Global, all states at present and in future), F (Future, a state at present or in future), X (next-time, the next state) and U (until, until a state).

The generation rules of CTL formulas are as follows:

Atomic proposition is the CTL formula;

If p, q are CTL formulas, $(\neg p), (p \vee q), (p \wedge q), (AG p), (AF p), (EG p), (EX p), (EF p), (E(p U q)), (A(p U q))$ are CTL formulas.

It should be noted that the operators A, E, and G, F, X, U must appear in pairs; otherwise, they are not CTL formulas. And each CTL formula is true or false in a given state.

For bus protocol, the properties required to be complied with are the starvation-free property, fairness and exclusiveness. The CTL description of the three properties of Wishbone bus protocol is as follows.

(1) No-starvation property: it is the first property of the bus protocol, which is also called activity. When the Master sends a request to the bus, the Master will eventually be allowed. It should be expressed by CTL formula: $\text{no_starve: AG (CYC}_m \rightarrow \text{AF GNT}_m)$; where, m represents any Master, when a Master of m sends request to the bus, it will eventually get the bus request permission after sending the request signal of CYC.

(2) Fairness: the property required for Arbiter by the protocol. It indicates that after the Master sends the bus request to the Arbiter, the Arbiter believes that it is not marked, and the request will eventually be allowed. It should be expressed by CTL formula: $\text{Fair: AG ((CYC}_m \ \& \ \sim \text{Label}_m) \rightarrow \text{AF GNT}_m)$;

where, m represents any Master, when the Master sends the request signal of CYC, and the signal has not been marked by the Arbiter, the Arbiter expresses trust to it, and the request will eventually be allowed.

(3) Mutual exclusion: in the system with multiple Masters, the same resource can not simultaneously be accessed. It ensures that no error occurs when multiple devices in the same system access to the same resource. It should be expressed by CTL formula: $\text{mutex: AG } (\sim \text{STB}_{m1} \vee \sim \text{STB}_{m2})$;

where, $M1$ and $m2$ represent two different Masters, the formula indicates that the two Masters can not access to the Slave.

4 Model Checking Analysis of the Protocol

Model checking is a formal checking method with high degree of automation, originally proposed by Clarke and Emerson. In formal checking, model checking achieves remarkable performances, and is successfully applied to modeling analysis of network security protocols and e-commerce protocols. Model checking process consists of three phases: (1) system modeling: To abstract the design to be verified as the finite state system model (such as the finite state machine, Kripke structure, etc.); (2) nature specifications: to use the temporal logic (such as the linear temporal logic of LTL, computation tree logic of CTL, etc.) to describe behavioral property of the design to be verified; (3) property checking: to convert the finite state system model to the input language for the corresponding tool, and it will be input into with the behavioral property expressed by tense logical, and the model checking tools (such as SPIN, SMV, etc.) are used for automatic analysis.

SMV is the formal checking tool based on the technology of Symbolic Model Checking [10]; its typical characteristics is to use binary decision diagram BDD for the data structure [11], which can effectively mitigate the state space explosion problem existing in the model test method largely to promote the practical application

of the model test method. When using SMV to perform the verification analysis, firstly, the system should be modeled as a finite state machine and converted into the SMV input language, and then the formal description of system properties should be carried out in the CTL formula, finally these two parts will be input to the SMV system for running. If the finite state system satisfies the properties described by the CTL, the output is TRUE, otherwise the output is FALSE, and the corresponding counter-example will be illustrated. The schematic diagram is shown in Fig. 8.

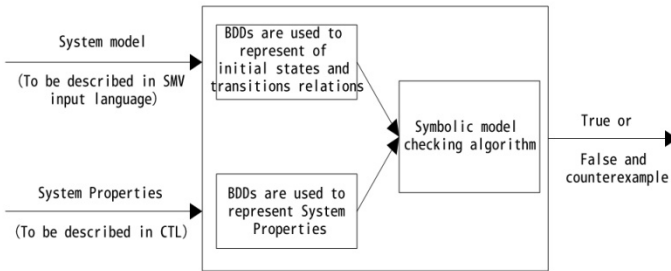


Fig. 8. Schematic diagram of SMV

```

Name | Layer
-----|-----
G 7. (top level)

Source | Trace | Log

File Show

module main()
{
  m1:master(a.grant1,s.resp);
  m2:master(a.grant2,s.resp);
  s:slave(m1.srb,m2.srb,m1.cyc,m2.cyc,a.master);
  a:arbiter(m1.cyc,m2.cyc,m1.srb,m2.srb,m1.state,m2.state);
  --no_starve1:assert G(m1.cyc -> F a.grant1);
  no_starve2:assert G(m2.cyc -> F a.grant2);
  --fair1:assert G(m1.cyc & ~a.label[1]) -> F a.grant1;
  --fair2:assert G(m2.cyc & ~a.label[2]) -> F a.grant2;
  --mutex: assert f(m1.srb | ~m2.srb);
  --using fair2 prove no_starve2;
  --assume fair2;
}

module master(grant,resp)
{
  srb:boolean;
  cyc:boolean;
  state:idle_m, reqbus_m, waitgrant_m, ctrlbus_m, start_m, normal_m, retry_m, error_m, finish_m;
  init(srb)=0;
  next(srb)=mase(
    ~grant | state=normal_m | state=retry_m | state=error_m | state=finish_m | state=idle_m);
  next(srb) & (state=ctrlbus_m | state=start_m) |;
}
    
```

Fig. 9. Part of SMV procedures when checking

Property	Result	Time
no_starve2	false	Mon Nov 24 15:12:25 滑 觀 开 演 模 块 解 释 2008

Source	Trace	Log
File	Edit	Run View
.....	1	2
a.cyc	0	1
a.grant1	1	1
a.grant2	0	0
a.grn	1	1
a.label[1]	0	0
a.label[2]	0	0
a.master	1	1
a.state	idle_a	idle_a
a.srb	0	0
m1.cyc	1	0
m1.state	idle_m	idle_m
m1.srb	0	0
m2.cyc	1	0
m2.state	idle_m	idle_m
m2.srb	0	0
s.cyc	0	0
s.prev.master	0	1

Fig. 10. Checking results of the property of no_starve2

The later developed versions of McMillan and Cadence Berkeley Labs for the SMV tool which is used for checking Wishbone bus protocol in this paper are called Cadence SMV [12]. to describe SMV input language used for finite state machine for modeling in part 2.2, and input it with system properties described by CTL formulas in part 2.3 into the SMV model checking tools for running. It is found through checking that the Master1 can meet the three properties of the system, but for the Master2, the system has no-starvation property and fairness, but can meet the mutual exclusion. Fig. 9 shows part of the SMV procedures submitted to the checking tool. Fig. 10 shows the result that the value is false after starvation property (i.e. property of no_starve2) is checked, and in the second half of the figure, a track leading to the result is given.

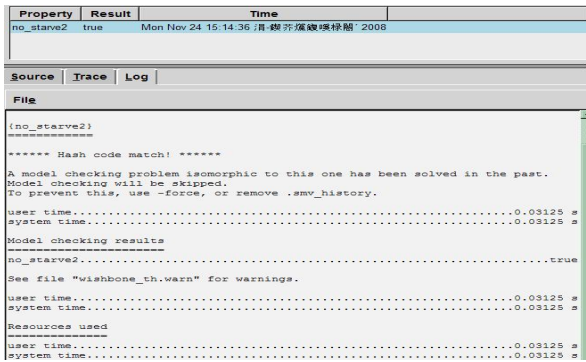


Fig. 11. Checking results of no-starvation property based on the assumption of meeting fairness

It can be seen from the track reflected by the checking tool that the bus request sent by Master2 to the Arbiter has not been allowed, which indicates the bus starvation, while the request sent by the Master1 is always allowed. After checking the fairness, it is found that the fairness (fair2) of Master2 to the bus protocol can not be complied with, and track obtained from the checking is the same with that from no compliance with no_starve2. In this paper, the further research and analysis is carried out, it is found that the no-starvation property of the bus protocol on the Master2 may be caused by the fairness of the Master2.

Based on the above situations, in this paper, the fairness assumption of the Arbiter is added when performing the modeling, and in the case that assumption conditions are met, the re-checking of no-starvation property (i.e. property no_starve2) is performed, and the result that no starvation property is true is obtained as shown in Fig. 11.

Based on the checking results, the starvation of Wishbone bus in a shared bus interconnection is caused by unfairness of the bus Arbiter. In terms of the integrity of the protocol specifications, it is the defect of Wishbone on-chip bus protocol. The study results also reflect that in the shared bus on-chip system, the good or poor design of the fairness of the arbitration mechanism has a significant impact on the system.

5 Conclusions

This paper analyzes the formal modeling of Wishbone bus protocol in a shared bus interconnection mode, and checks it using SMV model checking tool. The checking results show that the Wishbone on-chip bus protocol has defects. This study also shows that the fairness of the bus arbitration mechanism exerts significant influence on the performance of the system composed of multiple master and slave units in a shared bus interconnection mode. It will be our next research how to improve the shortcomings of the protocol.

Acknowledgments. Supported by the Key Program of Hechi University(Grant No. 2011YBZ-N002).

References

1. ARM. Advanced microcontroller bus architecture specification [S/OL] (1999), http://www.arm.com/armtech/AMBA_spec
2. IBM. 32-bit processor local bus architecture specifications [S/OL], Version 2.9, <http://www.ibm.com/chips/products/coreconnect/>
3. Open-Core Protocol Int. Partnership Association Inc. Open-core protocol specification [S/OL], Release 1.0 (2001), <http://www.ocpip.org>
4. WISHBONE, Revision B.3 Specification [S/OL], <http://www.opencores.org/-projects.cgi/web/wishbone/wishbone>
5. Roychoudhury, A., Mitra, T., Karri, S.R.: Using Formal Techniques to Debug the AMBA System-on-Chip Bus Protocol. In: The Design, Automation, and Test Europe Conference, Munich, Germany, pp. 828–833 (March 2003)
6. Chauhan, P., Clarke, E.M., Lu, Y., Wang, D.: Verifying IP-Core based System-On-Chip Designs. In: The IEEE International AS IC/SOC Conference, pp. 27–31 (September 1999)
7. Goel, A., Lee, W.R.: Formal verification of an IBM CoreConnect processor local bus arbiter core. In: DAC 2000, pp. 196–200 (2000)
8. Lin, H.-M., Yen, C.-C., Shih, C.-H., Jou, J.-Y.: On compliance test of on-chip bus for SOC. In: ASP-DAC 2004, pp. 328–333 (2004)
9. Clarke, E.M., Emerson, E.A., Sistla, A.P.: Automatic verification of finite-state concurrent systems using temporal logic specifications. *ACM Transactions on Programming Languages and Systems* 8(2), 244–263 (1986)
10. McMillan, K.L.: Symbolic Model Checking: An Approach to the State Explosion Problem. Carnegie-Mellon University publication CMU-CS-92-131 (May 1992)
11. Bryant, R.E.: Graph-based algorithms for Boolean function manipulation. *IEEE Transactions on Computers* C-35(8) (1986)
12. Cadence Berkeley Laboratories, California, USA. The SMV Model Checker (1999), <http://www-cad.eecs.berkeley.edu/kenmcmil/smv/>

A Time Synchronization Method for Wireless Sensor Networks

Chao Zou and Yueming Lu

Key Laboratory of Trustworthy Distributed Computing and Service (BUPT), Ministry of Education, Beijing University of Posts and Telecommunications, Beijing, 100876, China
{zouchaogood, ymlu}@bupt.edu.cn

Abstract. As consistent time scale is essential to facilitate group operations and improve network performance, time synchronization is regarded as a critical piece of infrastructure for distributed network measurement and control systems, especially for wireless sensor networks (WSNs). However, existing time synchronization algorithms for WSNs either do not provide enough scalability to achieve compatibility with other sync protocols, or do not fully take into account the characteristics of WSNs. This paper proposes a time synchronization method (TSM) to achieve precise time synchronization and reach the frequency drift compensation in WSNs at the same time. Evaluations show that TSM synchronizes wireless sensor nodes precisely with a magnitude of microsecond. Moreover, it has a good performance of stability and energy efficiency.

Keywords: Time synchronization, WSNs, Precision time synchronization, IEEE 1588, frequency drift compensation.

1 Time Synchronization and Challenges

As a complex of micro-electro-mechanism system (MEMS), wireless communication, transducer and digital electronic technology [1, 2], wireless sensor networks (WSNs) are seen as a new access technique for data acquiring and processing. With the rapid technical developments in these areas, WSNs are applied in more and more applications, such as intelligent home, medical diagnostics, environmental monitoring, battle-field surveillance, and so on [3, 4].

As a crucial piece of infrastructure for communication networks, especially for wireless sensor networks, time synchronization, which aims at realizing a common time scale in the whole networks, is essential for a number of group system operations and applications: time division multiplexing, coordination controlling, nodes localization, data consistency and fusion, energy-efficient scheduling and power management, etc. Consequently, inconsistent time scales among nodes in WSNs lead to ineffective information communication and group operations.

Many methods to synchronize WSNs have proposed in recent years. We can classify them into two categories. The first kind is using global positioning system (GPS). However, this method is impractical for large scale WSNs since having a GPS receiver on every node is a costly proposition (in term of money, power and size) and

the GPS signal is inaccessible in indoor or other harsh environment where cannot receive satellite signals. The second kind of time synchronization solutions is designed specifically for WSNs, for example, reference broadcast synchronization (RBS), timing-sync protocol for sensor networks (TPSN), flooding time synchronization protocol (FTSP) and tree-based synchronization algorithms [3], etc. Nonetheless, those synchronization protocols are incompatible with each other, especially for conventional internet network, and result in network partition. Therefore, a practical and standardization protocol for time synchronization in WSNs, as well as distributed applications, is necessary [5].

Employing the mature wired synchronization protocols to synchronize nodes in WSNs is an excellent option, but they cannot be directly used in WSNs since wired protocols do not take WSNs features (in terms of low bit rate, low energy consumption, limited storage and computing resource) into account, such as IEEE 1588 which is the precision time protocol (PTP) for wired network, especially for Ethernet. So, different situations between traditional wired network and wireless sensor networks make it valuable and significant to extend IEEE 1588 time synchronization protocol to WSNs.

The objective of this paper is to propose a time synchronization method (TSM) which enables precision time synchronization on WSNs nodes over Internet. In order to achieve this goal, TSM designed a system prototype which includes a PTP gateway connected with internet and wireless sensor nodes to realize the convergence between the conventional internet and WSN. At the same time, it does not only achieve time synchronization, but also reach frequency drift compensation supporting more precision synchronization.

2 IEEE 1588 Time Synchronization Algorithm and IEEE 802.15.4

In this section, the basic algorithm and communication protocol employed in TSM are introduced briefly as follows.

In the TSM, we exploited the basic IEEE 1588 time synchronization algorithm. Similar to other time synchronization protocols over packet, IEEE 1588 defines four types of timestamp message to realize synchronization: SYNC, FOLLOW_UP, DELAY_REQ and DELAY_RESP. The details and some modification version of IEEE 1588 are presented in [5, 6, 7, 8].

However, as a precision time synchronization standard for wireless network in the first place, IEEE 1588 cannot directly be used in WSNs without necessary improvement based on WSNs features.

Meanwhile, TSM also uses IEEE 802.15.4 as the communication protocol between nodes in WSNs and 6LoWPAN as the exchange protocol between Ethernet and WSNs. IEEE standard 802.15.4 [9] defines a communication and interconnection standard of devices via radio communication whose key targets are low energy consumption, low cost and low bit rate. This standard supports an over-the-air data rate of 250 kb/s in theory. In addition, two key advantages of employing IEEE 802.15.4 are accessing to a hardware synchronization signal and the integration of a

microcontroller with the transceiver in the same chip. Consequently, IEEE 802.15.4 is well suited for the communication of WSNs. Besides, we also exploit 6LoWPAN, which is the acronym of IPv6 over Low power Wireless Personal Area Networks, in our scheme because it access IPv6 Internet-compatible datagrams over IEEE 802.15.4 and can interoperate with other IEEE 802.15.4 protocol, such as zigbee [10].

In addition, in order to connect the IEEE 802.15.4 wireless networks and wired networks like Ethernet, A special gateway is required. Meanwhile, the gateway also functions as a wireless router. The method proposed engages the 6LoWPAN wireless protocol to realize communication with IPv6 based on IEEE 802.15.4. In our scheme, both wireless sensor nodes and the associated gateway have a RF transceiver unit with IEEE 802.15.4 capability, a microcontroller unit and a time synchronization processing unit.

3 TSM

The goal of TSM is to extend the IEEE 1588 core concepts to WSNs. However, it cannot be achieved without enhancements according to features of WSNs. Three primary challenges must be conquered before applying the wired standard to wireless:

Different resource requirements: wired standard wired IEEE 1588 asks for a relatively large amount of processing power, computing capability and storage space which is inadvisable for battery energy and limited resource wireless nodes.

Different frame length: the length of the standard wired IEEE 1588 synchronization message (166 bytes) is too long for IEEE 802.15.4 standard (128 bytes limited), such as zigbee or other related wireless protocol. It must be fragmented or shortened. In addition, a specific gateway is needed at the same time.

Different data rate: the data rate of wireless (250 Kb/s for IEEE 802.15.4 at frequency of 2.4 GHz in the ideal case) [9] is far lower than wired (generally tens to hundreds Mb/s). Consequently, at least a relatively large gateway buffer is required.

To solve the mentioned problems, the TSM is presented in two steps in following subsections: First, TSM exploits the core theory of IEEE 1588 algorithm to realize time synchronization. Second, TSM introduces a frequency drift compensation mechanism to enhance the stability of the slave time reference.

3.1 Time Synchronization

As shown in fig. 1, the first step in TSM is to realize time synchronization. The proposed sync procedure adheres to the core of 1588 only with some modifications for wireless sensor network. First, it abbreviates the sync message by gateways in accordance with 6LoWPAN protocol. Because messages are fragmented two or more parts, more power is required and more network congestion happens comparing with our scheme. Second, FOLLOW_UP messages are no longer necessary even not be used in the proposed vision, except for special cases that demand the full format version. Most importantly, it reduces energy consumption and network congestion caused by the overhead of time synchronization without degrading sync performance.

Last but not the least, there must be a gateway to ensure the data compatibility between the wired network (Ethernet) and the wireless network (6LoWPAN).

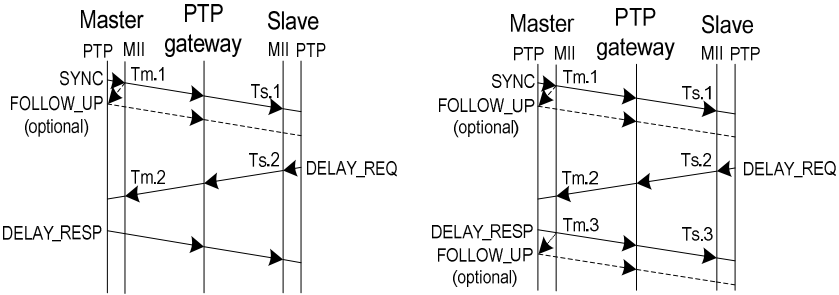


Fig. 1. Step 1: time synchronization for wireless sensor nodes **Fig. 2.** Step 2: frequency drift compensation for wireless sensor nodes

Note that the time stamping point is critical for the accuracy of time synchronization. Known from existing works [11, 12], the time stamping point should be as close to the physical layer as possible to obtain a minimum jitter of timestamp value. To attain a highly level of accurate time synchronization, the most desirable stamping place is the Media Independent Interface (MII) between Media Access Control (MAC) layer and physical (PHY) layer.

The synchronization packet consists of four parts: preamble, start of frame delimiter (SFD), frame length and MAC protocol data. Time stamps in TSM are taken at MII when the last bit of SFD is detected by sensor nodes with the supporting of the hardware. Then, nodes process the value of the local clock counter in accordance with the synchronization protocol.

3.2 Frequency Drift Compensation

Affected by many factors, such as oscillator aging, temperature and humidity, the slave clocks frequency will change over time, which is known as frequency drift. Consequently, slave nodes clock depart from the master clock over time and need time synchronization periodically. However, frequent time synchronization leads to high power consumption and increase network congestion. In order to keep frequency drift between master and slave as low as possible, it is necessary to introduce a feasible and effective drift compensation mechanism in WSNs. Furthermore, the compensation mechanism also improves the synchronization accuracy.

Although discussed in some papers [5, 13], this problem is not solved perfectly up to now. In these solutions, a high demand of computing and storage capacity is oppressive for WSNs. The scheme proposed in this paper presents a drift compensation mechanism which removes frequency drift without the requirement of extra message exchange or storage space.

Based on the first part of procedure depicted in previous subsection, the second step of TSM is to realize frequency drift compensation as shown in fig. 2. In our scheme, the value of the DELAY_RESP message is recorded in master time

reference. At the same time, the master transmits DELAY_RESP message with the value of the timestamp $Tm.2$ and $Tm.3$. The slave node measures and stores the receipt time of DELAY_RESP. Then, slave nodes calculate the frequency drift ratio ρ with equation (1):

$$\rho = \frac{Tm.3 - Tm.1}{Ts.3 - Ts.1} \quad (1)$$

In slave clocks, the adjustment procedure for drift compensation follows three steps:

Step 1: the DELAY_RESP transmission timestamp $Ts.2$ is corrected to remove frequency drift. The corrected value $T^*s.2$ in the slave time reference is obtained from equation (2):

$$T^*s.2 = Ts.1 + \rho(Ts.2 - Ts.1) \quad (2)$$

Step 2: the accurate value of offset, i.e. $offset^*$, between master and slave is evaluated based on conventional concept of IEEE 1588 as follow:

$$offset^* = \frac{(Ts.1 - Tm.1) + (T^*s.2 - Tm.2)}{2} \quad (3)$$

Step 3: the slave node completes time synchronization by exploiting the corrected value of $offset^*$, the synchronized time reference in slave clock $Ts.m$ is obtained as in equation (4):

$$T.s_m = Ts - offset^* \quad (4)$$

After execution of the above steps, the slave clock realizes time synchronizing with the master clock. In a synchronization cycle in terms of the interval between two consecutive synchronization procedures, the timestamp of a generic event timestamp $Ts.e$ can be removed frequency drift and the related corrected value $T^*s.e$ equation (5):

$$T^*s.e = T^*s.2_m - \rho(Ts.e - T^*s.2_m) \quad (5)$$

Where $T^*s.2_m$ is the corrected value of $T^*s.2$ in the synchronized slave time reference based on equation (4).

Due to the frequency drift, the slave time drifts apart from the master time after time synchronization. Consequently, periodical time synchronization is demanded. However, frequent message exchanges in WSNs increase energy consumption and the probability of network congestion. For better energy efficiency and decreasing message exchanges caused by time synchronization, the interval of synchronization should be as long as possible. It is obvious that the fixed short default interval of synchronization does not meet the requirements of WSNs. The sync interval in TSM is not fixed, but be variable with the accuracy requirement of different applications. The synchronization interval is evaluated as equation (6):

$$interval_{sync} = \frac{error_{max}}{|1 - \rho|} \quad (6)$$

Where $error_{max}$ is the maximum time error permitted by the specific application.

4 Performance Evaluation

The performance of TSM is evaluated by simulations. Our simulating network consists of one master clock sever, an Ethernet node (switch or router), a special gateway (designed in the proposed protocol) and a wireless sensor. Note that the wireless sensors communicate with the gateway on the air.

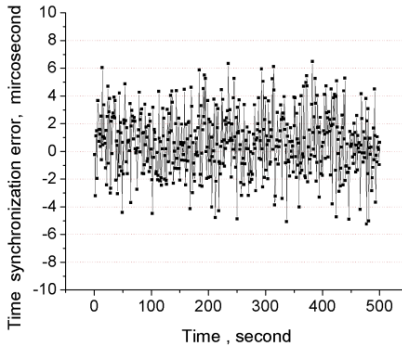


Fig. 3. Time synchronization error using TSM

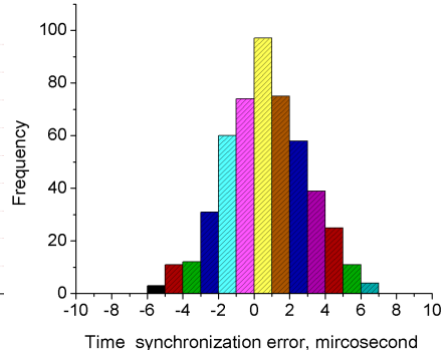


Fig. 4. Distribution of time synchronization error using TSM

The first evaluation is to evaluate the accuracy of TSM. The synchronization process is triggered at a frequency of per second, and time stamps are collected from both master and slave nodes to calculated offset in slave clock once per second. The experiment was repeated 500 times. Fig. 3 depicts the values of synchronization error using TSM, i.e. the remaining offset in the slave node from the master sever after time synchronization, and fig. 4 depicts the distribution of the result. The two graphs imply that the closer the offset is to zero and the more concentrated of the distribution, the more precise time synchronization, and the zero value indicates perfect time synchronizing. A magnitude of sub-millisecond can meet the precision requirement of a large number of applications in WSNs. The maximum absolute value of the error using TSM is about 5 microseconds. Known from the above evaluation, TSM can satisfy the strict precision requirement of most WSNs applications.

To evaluate the stability and energy consumption of the proposed algorithm, we also use the conventional IEEE 1588 to synchronize the wireless sensor node in our evaluation, which does not have frequency drift compensation mechanism and the sync messages are separated two parts due to the message length constraint at a fixed sync interval in WSNs.

Fig. 5 plots the offset in slave from the master sever using different synchronization method as a function of time after synchronizing procedure. The offset of the conventional IEEE 1588 drifts apart from the master clock reference gradually after the synchronization procedure, because the frequency of the clock crystal is not exactly equal to that of master's caused by frequency drift effect. On the contrary, the offset employing TSM changes little over time by exploiting frequency drift compensation.

In a wireless sensor system, all nodes are normally is a sleep mode and only transmit messages when necessary. Table 1 plots the time synchronizing interval as a

function of the time precision requirement of applications. In conventional IEEE 1588 protocol, the sync interval is fixed (the default value is 2s), but the interval of TSM is variable and be adjusted with the precision requirement of applications. The power consumption, which directly determines the battery even the whole node lifetime in WSNs, increases obviously with frequency sync updating. Furthermore, frequent message exchange in WSN increase the probability of network blocking and message retransmission which make significant effect on other communication in whole network.

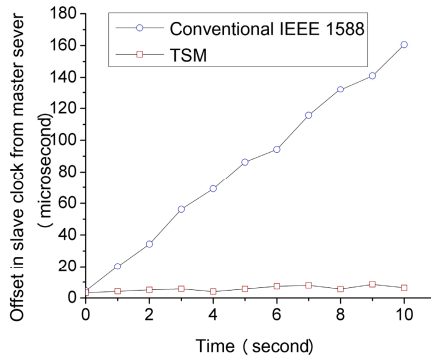


Fig. 2. Offset of the slave clock using different synchronization methods as a function of time after synchronizing procedure

Table 1. Time synchronization interval vs. the precision requirement of application

Precision(μ s)	30.82	50	100	200	500
Interval (second)					
Conventional IEEE 1588	2.0	2.0	2.0	2.0	2.0
TSM	2.0	3.2	6.5	13.1	32.5

5 Conclusion and Future Work

Time synchronization is a common requirement for most of applications in distributed devices networks. Especially, consistent time scales among wireless sensor nodes is essential to facilitate group operations and improve network performance in WSNs. For providing a practical and compatible time synchronization standard for WSNs, this paper proposes TSM to synchronize distributed sensor nodes in WSNs. The performance evaluation manifests that TSM provides a microsecond magnitude of global time synchronization for WSNs with low energy consumption and good stability. Our future work is to propose more sophisticated time synchronization based on IEEE 1588 v2 and to apply them in more practical applications.

Acknowledgements. This research was supported in part by National 863 Program (No. 2011AA01A205) and National 973 Program (No. 2011CB302702).

References

1. Akyildiz, I.F., Su, W., et al.: Wireless Sensor Networks: A Survey 38(4), 393–422 (2002)
2. Noh, K., Serpedin, E., Qaraqe, K.: A New Approach for Time Synchronization in Wireless Networks: Pairwise Broadcast Synchronization. *IEEE Transactions on Communications* 7(9), 3318–3322 (2008)
3. Lasassmeh, S.M., Conrad, J.M.: Time Synchronization in Wireless Sensor Networks: A Survey. In: *IEEE SoutheastCon 2010*, pp. 242–245 (March 2010)
4. Ping, S.: Delay Measurement Time Synchronization for Wireless. Intel Research, IRB-TR-03 (2003)
5. Cho, H., Jung, J., et al.: Precision Time Synchronization Using IEEE 1588 for Wireless Sensor Networks. In: *International Conference on Computational Science and Engineering*, pp. 579–586 (October 2009)
6. Akyildiz, I.F., Su, W., Sankarasubramaniam, Y., Cayirci, E.: Wireless Sensor Networks: A Survey. *Computer Networks* 38(4), 393–422 (2002)
7. Lv, S., Lu, Y., Ji, F.: An Enhanced IEEE 1588 Time Synchronization for Asymmetric Communication Link in PTN. *IEEE Communications Letters* 4(8), 764–767 (2010)
8. Du, Z., Lu, Y., Ji, F.: An Enhanced End-to-End Transparent Clock Mechanism with a Fixed Delay Ratio. *IEEE Communications Letters* 15(8), 872–874 (2011)
9. IEEE Computer society, IEEE Std. 802.15.4 TM-2003 (2003)
10. Wobschall, D., Ma, Y.: Synchronization of wireless sensor networks using a modified IEEE 1588 protocol. In: *2010 International IEEE Symposium on Precision Clock Synchronization for Measurement Control and Communication (ISPCS)*, pp. 67–70 (October 2010)
11. Cooklev, T., Eidson, J.C., Pakdaman, A.: An Implementation of IEEE 1588 Over IEEE 802.11b for Synchronization of Wireless Local Area Network Nodes. *IEEE Transactions on Instrumentation and Measurement* 56(5), 1632–1639 (2007)
12. Rosselot, D.: Simple, Accurate Time Synchronization in an Ethernet Physical Layer Device. In: *2007 International IEEE Symposium on Precision Clock Synchronization (ISPCS) for Measurement, Control and Communication, Vienna, Austria*, pp. 123–127 (October 2007)
13. Ferrari, P., Flammini, A., et al.: IEEE 1588-based Synchronization System for a Displacement Sensor Network. In: *Instrumentation and Measurement Technology Conference, Sorrento*, vol. 47(2), pp. 254–260 (April 2006)

Real-Valued Negative Selection Algorithm with Variable-Sized Self Radius

Jinquan Zeng^{1,2}, Weiwen Tang², Caiming Liu³, Jianbin Hu⁴, and Lingxi Peng⁵

¹ School of Computer Science & Engineering, University of Electronic Science and Technology of China, 610054 Chengdu, China

² Sichuan Communication Research Planning & Designing Co., Ltd, 610041 Chengdu, China

³ Laboratory of Intelligent Information Processing and Application, Leshan Normal University, 614004 Leshan, China

⁴ School of Electronics & Information, Nantong University, 226019 Nantong, China

⁵ Department of Computer and Education Software, Guangzhou University, 510006 Guangzhou, China

zengjq@uestc.edu.cn

Abstract. Negative selection algorithm (NSA) generates the detectors based on the self space. Due to the drawbacks of the current representation of the self space in NSAs, the generated detectors cannot enough cover the non-self space and at the same time, cover some of the self space. In order to overcome the drawbacks, a new scheme of the representation of the self space is introduced with variable-sized self radius, which is called VSRNSA. Using the variable-sized self radius to represent the self space, we can generate the more quality detectors. The algorithm is tested using the well-known real world datasets; preliminary results show that the new approach enhances NSAs in increasing detection rates and decrease false alarm rates, and without increase in complexity.

Keywords: Artificial Immune Systems, Negative Selection Algorithm, Anomaly Detection.

1 Introduction

Biological immune systems (BIS) have many characteristics such as uniqueness, autonomous, recognition of foreigners, distributed detection, and noise tolerance [1]. Inspired by BISs, Artificial Immune Systems (AIS) have become one of the relatively new areas of soft computing [2-5] and AISs generally include clonal selection based algorithms, negative selection based algorithms and artificial immune network models [7-9].

One of the major algorithms developed within AISs is the NSA, proposed by Forrest et al. [10]. The NSA can only use self samples to train detectors for classifying unseen data as self or non-self and its typical applications include anomaly detection, fault detection, especially, network security. Early works in NSAs used the problem in binary representation [10]. However, many applications are natural to be described in real-valued space and cannot be processed by NSAs in binary

representation [11]. Recently, more concerns focused on real-valued negative selection algorithm (RNS) [12-13]. The algorithms use a real-valued representation of the self/non-self space and can speed up the detector generation process [14]. Another important variation among RNSs, is V-detector [15], which uses variable-sized detector and terminates training stage when enough coverage is achieved.

In most of real applications, we can not get all of self samples, and so we have to use some of self samples to train detectors. In order to decrease false alarm rate, the self radius is adopted and large self radius means that more unseen data is divided into self space, resulting in low detection rate and vice versa. However, almost in all NSAs, the self radius is in constant size, these methods cannot build an appropriate self profile of the system. This paper tries to address the issue and a new scheme of the representation of the self space is introduced with variable-sized self radius, which is called VSRNSA. It is applied to perform anomaly detection for well-known real world datasets; preliminary results show that the new approach enhances NSAs in increasing detection rates and decrease false alarm rates, and without increase in complexity. It also is a general approach that can be applied to different anomaly detection problems.

2 Real-Valued Negative Selection Algorithm with Variable-Sized Self Radius

2.1 Constructing the Self Profile of the System

In BIS, negative selection is a mechanism to protect body against self-reactive lymphocytes and self-reactive T-cells are eliminated by a controlled death. As a result, only self-tolerant T-cells survive the negative selection process and are allowed to leave the thymus. Similarly, the negative selection algorithm generates detector set by eliminating any detector candidate that match elements from a group of self samples. In real applications, it is difficult to get all of self samples, such as benign files in every computer. In order to train detectors, we have to use only some of self sample to build the profile of the system that reflects the normal behavior. In RNSs, the self-radius is introduced to allow other elements to be considered as self elements which lie close to the self-center and represents the allowed variability of the self samples, illustrated in Fig. 1. The self radius of self sample specifies the capability of its generalization (the elements within the self radius of the self sample is considered as self elements). The bigger the self radius is, the more generalization the self sample is. Fig.1 also illustrates that the const-sized self radius cannot construct an appropriate profile of the system. The self radius is too small, the self space cannot be covered enough and high false positive rate is occurred. The self radius is too large, part of self elements cover the non-self space and result in false negative errors. Fig.2 illustrates that variable self radius can appropriately cover self space and build the profile of the system. In conventional NSAs, the self radius is in constant size, so the appropriate profile of the system cannot be built. While our proposed approach adopts variable-sized self radius and the appropriate profile of the system can be built, the appropriate profile of the system can increase the true positive rate and decrease the false positive rate.

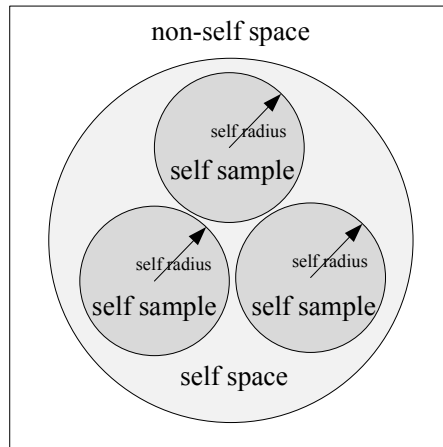


Fig. 1. Constructing the self profile using constant-sized self radius in 2-dimensional space

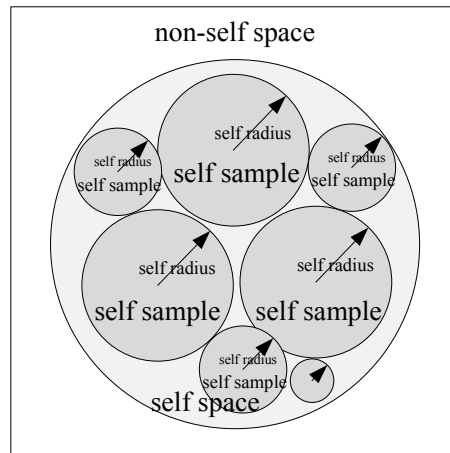


Fig. 2. Constructing the self profile using variable-sized self radius in 2-dimensional space

2.2 The Proposed Algorithm

A NSA consists of two phases, training and detecting phase. In training phase, the detectors are generated randomly and those that match any self samples using Euclidean distance matching rule are eliminated. Comparing with the version of constant-sized self radius, the most important differences in VSRNSA lie in steps 2 through 8 in Fig.3. We count the distance among the set of self samples and assign a variable-sized self radius based on the total distance of every sample to other self samples. Now that we let each self sample in the training set has its own self radius in addition to the distance to other self samples. Big distance means that the self sample is far from other self sample, and so the number of self elements near the self sample

Algorithm 1: VSRNSA($S, m, n, T_{\max}, c_0, MSC$)

Input: S =the set of self samples, n =the self radius coefficient, T_{\max} =maximum number of detector, c_0 =estimated coverage, MSC =maximum self coverage

Output: D = the Set of generated Detectors

```

1 Begin
2   for each  $s \in S$  do
3     total=0
4     for each  $p \in S$  do
5       total=total +dist(s,p) //dist(s,p), Euclidean distance between s and p
6     endfor
7     s.r= total/(|S|-1)/n // s.r, the self radius of the sample s
8   endfor
9    $D \leftarrow \emptyset$ 
10  repeat
11   $t \leftarrow 0$ 
12   $T \leftarrow 0$ 
13   $r \leftarrow \infty$ 
14   $x \leftarrow$  random point from  $[0, 1]^n$ 
15  for each  $d \in D$  do
16    if  $\text{dist}(d, x) \leq d.r$  then // d.r is the detection radius of the detector d
17       $t \leftarrow t+1$ 
18      if  $t \geq 1/(1-c_0)$  then
19        return  $D$ 
20      endif
21    goto 12
22  endif
23 endfor
24 for each  $s \in S$  do
25    $l \leftarrow \text{dist}(s, x)$ 
26   if  $l-s.r \leq r$  then
27      $r \leftarrow l-s.r$ 
28   endif
29   if  $r > s.r$  then
30      $x.r=l- s.r$ 
31      $D \leftarrow D \cup \{x\}$ 
32   else
32      $T \leftarrow T+1$ 
33     if  $T > 1/(1-MS C)$  then
34       exit
35     endif
36   endif
37 until  $|D| = T_{\max}$ 
38 end

```

Fig. 3. Real-Valued Negative Selection Algorithm with Variable Self Radius (VSRNSA)

is little and low self radius is assigned to the self sample. If the distance is little, it shows that the number of the self elements near the self samples is high and so the big self sample is assigned to the self sample. Based the constructing self profile using the variable-sized self radius, the detecting radius of detectors and the end of the algorithm is similar with constant-sized self radius [15]. The detection radius of the detectors is decided by the closest self sample and the end of the algorithm is decided by estimated coverage c_0 , maximum number of detectors T_{max} and maximum self coverage MSC [15].

3 Experiments and Results

To study the property and possible advantages of VSRNSA, we performed the experiments with a classical dataset used extensively in the pattern recognition literature, the Fisher's Iris data set, which includes three different classes of flowers: setosa, virginica and vericolor. In the dataset, each element is described by four attributes and each class is different from the others. To measure the distance on the same scale, the data is normalized first, depicted by equation (1).

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \quad (1)$$

Fig. 4. shows that the self radius coefficient affects the detection rate and false alarm rate. It is easy to see the effect of the self radius coefficient on the results. The larger self radius coefficient would result in the smaller self radius and then high detection rate but high false alarm rate are attained, thus suitable when we need high detection rate. On the other hand, the smaller self radius coefficient would result in the larger self radius and low high false alarm rate but low detection rate, thus suitable when we need low false alarm rate.

In order to determine the performance and possible advantages of VSRNSA, we compared the results obtained using the constant-sized self radius, namely V-detector. V-detector combines real-valued negative selection algorithm and variable-sized detectors and enhances the negative selection algorithm in efficiency. Nevertheless, the self radius is constant sized in V-detector and cannot build an appropriate profile of the system. Table 1 illustrates the comparison using the Fisher's Iris data set. Different experiments were performed in each case using one of the three classes as the normal and the other two as abnormal. The training data was either partially or completely composed by elements of the normal class. The self radius coefficient value was 30 and the results shown were the average of 100 different runs for each method.

A good detection system should have high detection rate and low false alarm rate. Our proposed approach, VSRNSA, has higher detection rate but lower false alarm rate, e.g. when the training data is the 50% elements of the setosa, the detection rate and false alarm rate of V-detector are 99.91% and 2.36%, but the detection rate of VSRNSA increases to 99.94% and the false alarm rate of VSRNSA decreases to 2.2%. Table 1 also shows that VSRNSA has another advantage, which is VSRNSA needs smaller number of detectors and has a good coverage, e.g. when the training

data is the total elements of the verginica, V-detector needs 245.63 detectors and the detection rate is 86.59%, but VSRNSA only needs 169.38 detectors and the detection rate increases to 90.19%.

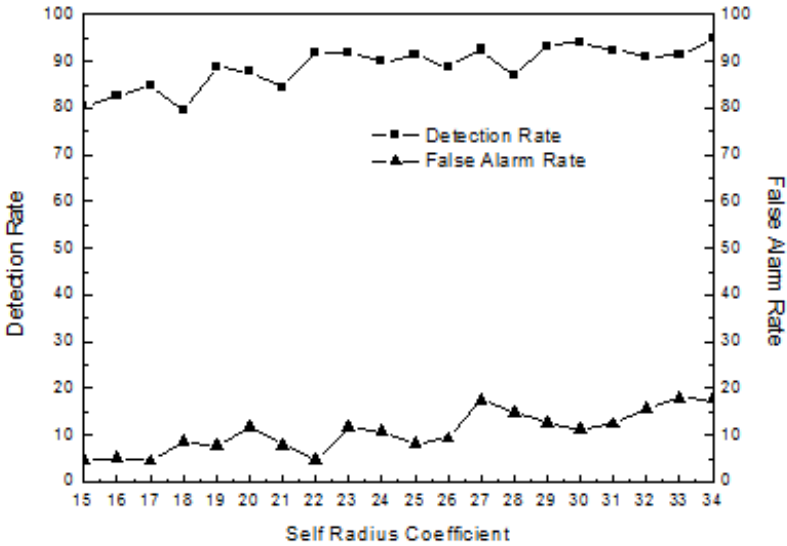


Fig. 4. The coefficient of self radius affect on the detection rate and false alarm rate

Table 1. Comparison between VSRNSA and V-detector using Fisher’s Iris dataset

Training Data	Algorithm	Detection Rate	False Alarm Rate	Number Detectors	of
Setosa	V-detector	99.91%	2.36%	16.02	
50%	VSRNSA	99.94%	2.20%	14.49	
Setosa	V-detector	99.96%	0.00%	22.51	
100%	VSRNSA	99.97%	0.00%	20.39	
versicolor	V-detector	90.86%	4.80%	113.49	
50%	VSRNSA	90.39%	4.52%	106.66	
versicolor	V-detector	86.94%	0.00%	156.14	
100%	VSRNSA	86.45%	0.00%	148.41	
verginica	V-detector	89.92%	17.56%	137.66	
50%	VSRNSA	92.87%	22.88%	92.36	
verginica	V-detector	86.59%	0.00%	245.63	
100%	VSRNSA	90.19%	0.00%	169.38	

4 Conclusion

The paper presented a real-valued negative selection algorithm with variable-sized self radius, VSRNSA, which is called VSRNSA. Using the variable-sized self radius to represent the self space, we can construct the appropriate profile of the system, and then

generate the more quality detectors, which can increase the true positive rate and decrease the false positive rate. The experiment results show that VSRNSA is an efficient NSA and offers the characteristics of high detection rate and low false alarm rate.

Acknowledgments. This work is supported by special technology development fund for research institutes of the Ministry of Science and Technology of China (2009EG126226, 2010EG126236, and 2011EG126038), China Postdoctoral Science Foundation (20100480074), Supported by the Fundamental Research Funds for the Central Universities (ZYGX2011J069) and NSFC (61103249, 61100150).

References

1. de Castro, L., Zuben, F.: Artificial Immune Systems: Part I – Basic Theory and Applications. TR – DCA 01/99 (1999)
2. Zhang, F.B., Yue, X., Wang, D.W., Xi, L.: A Principal Components Weighted Real-valued Negative Selection Algorithm. *International Journal of Digital Content Technology and its Applications* 5(6), 313–324 (2011)
3. Zhang, F.Y., Qi, D.Y.: Run-time malware detection based on positive selection. *Journal in Computer Virology* 7(4), 267–277 (2011)
4. Greensmith, J., Aickelin, U., Cayzer, S.: Introducing Dendritic Cells as a Novel Immune-Inspired Algorithm for Anomaly Detection. In: Jacob, C., Pilat, M.L., Bentley, P.J., Timmis, J.I. (eds.) ICARIS 2005. LNCS, vol. 3627, pp. 153–167. Springer, Heidelberg (2005)
5. Manzoor, S., Shafiq, M.Z., Tabish, S.M., Farooq, M.: A Sense of ‘Danger’ for Windows Processes. In: Andrews, P.S., Timmis, J., Owens, N.D.L., Aickelin, U., Hart, E., Hone, A., Tyrrell, A.M. (eds.) ICARIS 2009. LNCS, vol. 5666, pp. 220–233. Springer, Heidelberg (2009)
6. Dervovic, D., Zuniga-Pflucker, J.C.: Positive selection of T cells, an in vitro view. *Semin. Immunol.* 22(5), 276–286 (2010)
7. de Castro, L., Zuben, F.: Learning and Optimization Using the Clonal Selection Principle. *IEEE Transactions on Evolutionary Computation* 6(3), 239–251 (2002)
8. Berna, H.U., Sadan, K.K.: A review of clonal selection algorithm and its applications. *Artificial Intelligence Review* 36(2), 117–138 (2011)
9. Al-Enezi, J.R., Abbod, M.F., Alsharhan, S.: Artificial Immune Systems-models, algorithms and applications. *International Journal of Research and Reviews in Applied Sciences* 3(2), 118–131 (2010)
10. Forrest, S., Perelson, A.S., Allen, L., Cherukuri, R.: Self-nonsel self discrimination in a computer. In: *Proceedings of the 1994 IEEE Symposium on Research in Security and Privacy*. IEEE Computer Society Press (1994)
11. Dai, H., Yang, Y., Li, C.: Distance Maintaining Compact Quantum Crossover Based Clonal Selection Algorithm. *JCIT* 5(10), 56–65 (2010)
12. Stibor, T., Timmis, J., Eckert, C.: A Comparative Study of Real-Valued Negative Selection to Statistical Anomaly Detection Techniques. In: Jacob, C., Pilat, M.L., Bentley, P.J., Timmis, J. (eds.) ICARIS 2005. LNCS, vol. 3627, pp. 262–275. Springer, Heidelberg (2005)
13. Zhou, J., Dasgupta, D.: Applicability issues of the real-valued negative selection algorithms. In: *Proceedings of Genetic and Evolutionary Computation Conference* (2006)
14. Gonzalez, F.A., Dasgupta, D.: Anomaly detection using real-valued negative selection. *Genetic Programming and Evolvable Machines* 4, 383–403 (2003)
15. Ji, Z., Dasgupta, D.: Real-Valued Negative Selection Algorithm with Variable-Sized Detectors. In: Deb, K., Tari, Z. (eds.) GECCO 2004, Part I. LNCS, vol. 3102, pp. 287–298. Springer, Heidelberg (2004)

Scale Effect on Soil Attribute Prediction in a Complex Landscape Region

Zhenfu Wu, Yanfeng Zhao^{*}, Li Qi, and Jie Chen

School of Water Conservancy and Environment, Zhengzhou University, Zhengzhou 450001,
China

{wfgjt1988, wish1005}@163.com, {yfzha, jchen}@zzu.edu.cn

Abstract. Total 283 soil samples were collected in 1220 km² area of Dengfeng county, Henan province, according to a nested sample strategy. Several scenarios were designed to research the scale effect on mapping soil organic material (O.M) with regression kriging interpolation. It was found that the trend of soil O.M on the elevation factor was macroscopical and that could be fitted optimally using only large scale data. If small scale data were added in simulation the precision of trend equation would decrease. However small scale data was contributive to the prediction of the residue in regression kriging, which could reveal not only spatial variability of residue in small scale but also enhanced the spatial structure in large scale and improved effectively the prediction. Therefore, the optimal way of soil O.M in regression kriging was that extracting the trend using only large scale data and simulating the residual using data of the both scales.

Keywords: Soil attribute, Scale effect, The trend, The residual.

1 Introduction

Spatial variability of soil attributes comes of multiple spatial scales[1]. The relationship between property of soil and environmental factor is affected strongly by spatial scale, same is true among different soil attributes [2]. Research under a single scale has own limitations and one-sidedness to some extent. Spatial variation pattern of soil attribute at different scale reveals scale effect on soil attribute[3-7]. Prediction accuracy can be improved effectively and spatial variability at different scale can be exploited fully researching soil prediction model in the scale point of view, that provides a reliable basis for studies on precision fertilization and soil evolution[8-10]. Li D Y (2008) indicated soil organic carbon in topsoil was different among town-scale, county-scale and city-scale[11]. Spatial heterogeneity of soil nutrient and soil character at different spatial scales, for example, soil pH, soil organic matter, and soil trace element, has been described (Zhao H X et al., 2005; Zhao J et al., 2006) [12-14]. Researches above showed that soil attributes of soil organic carbon, heavy metal, pH, organic matter, trace elements, were different significantly and displayed variety of spatial variation law, in the multi-scale. Meanwhile, multi-scale nested kriging

^{*} Corresponding author.

interpolation of soil attributes based on various spatial law improved mapping accuracy effectively, comparing with ordinary kriging under a single scale(Yu Q et al., 2007; Huo X N et al., 2009) [15-16].

For a exact description of spatial variability in small-scale as much as possible, researchers generally nested sample selectively on the basis of large-scale sampling to understand details of spatial variation in a fewer scale. However, how to synthesize datasets both in large-scale and small-scale in mapping process implies the scale effect, and the mapping result depends on modeling method. The present study were more limited to compare differences of soil attributes among multiple scales and to interpolation nested simply based multi-scale datas. Nevertheless, the impact of scale effect on prediction of soil attributes, etc, need some more further researches and discussions. In this paper, we taked soil organic matter of a complex landscape area in Dengfeng county as an example and focused on the impact of scale effect on soil prediction model.

2 Materials and Methods

2.1 Study Area

The study area, Dengfeng county(112°43'E ~ 113°48'E, 34°16'N ~ 34°35'N), is located in the southern region of the Mount Songshan in Henan province. The total area of Dengfeng is 1220km². Elevations of the county are usually less than 1510m and yet more than 80m. The topography of Dengfeng is characterized by middle mountains, low mountain, hills and river valley. Dengfeng has a the north temperate monsoon climate with distinct seasons, mild climate and sufficient sunlight, however, local climate is completely different due to complex terrain. Soil parent material in Dengfeng county are residual/slope deposit, proluvium and loess. The main soil types in the county are brown soil, cinnamon soil and aquic soil.

2.2 Soil Sampling and Chemical Analysis

92 topsoil samples were collected with 2km sampling interval as large-scale data. 45,73,52 and 31 samples were collected as small-scale data with 300m, 500m, 200m and 300m respectively in four landscape region which divided from the study area according to topography, parent material, soil types and other factors: Quartz Sandstone Zone in middle mountain, Diluvium Zone in low mountain, Purple Soil Zone and Proluvial fan in mountain front. 10 samples were shared in large-scale and small-scale, then a total of 283 topsoil samples were collected finally(Fig. 1). All sample sites were recorded using a hand-held global position system (GPS) and related information such as land use history, vegetation, and micro morphological characteristics were also recorded in detail.

The collected soil samples were air-dried and then ground to 60 meshes for chemical analysis. All of the soil samples were analyzed for soil organic matter. Soil organic matter determination was conducted with external heating potassium dichromate volumetric method, and specific operating processes refered to reference.

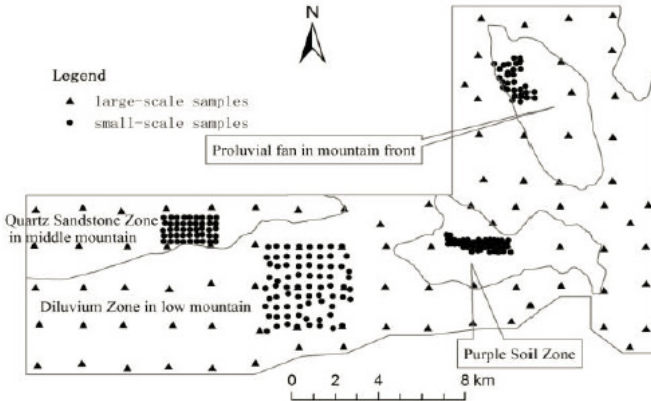


Fig. 1. Distribution of soil samples

2.3 Kriging Interpolation and Multi-scale Nested Semivariance Model

In a complex terrain region, the trend is generally eliminated using the trend surface equation by universal kriging or regression kriging, and then the residual is estimated with ordinary kriging. Universal kriging and regression kriging both divide the regionalized variables into trend value (deterministic component) and normally distributed residual value (stochastic component), the trend of former is just related to the space coordinates, yet the latter is related to other covariates.

Variation functions are the main tool for analyzing spatial structure which reflects and describes the spatial structure characteristics of regionalized variables. In general, regionalized variable contain different changes of scales, its spatial structure can not be express with just a simple model, two or more models is needed to indicate spatial structure information of different distance, that named nested structure:

$$\gamma(h) = \gamma_0(h) + \gamma_1(h) + \dots + \gamma_n(h) = \sum_{i=0}^n \gamma_i(h) \quad (1)$$

Where $\gamma_0(h)$ is the nugget variance (including measured error and spatial variability which can not be characterized within minimum sampling interval) of a nested structure and $\gamma_1(h)$ represent components of different scales.

3 Results and Discussion

3.1 Statistics

The classical statistic character of soil organic matter in multi-scale (Table 1) showed the coefficient of variation of soil organic matter in small-scale various landscapes was from 16.81% to 36.15% and belonged to a moderate degree of variation,

according to the classification of the coefficient of variation: <10% for the weak variation, between 10% and 100% for moderate variability and >100% for the strong variability. Among four landscape regions, proluvial fan in mountain front had a highest average content of soil organic matter of 25.26 g kg⁻¹, it probably mainly due to the soil parent material of alluvium which soil developed on was more fertile generally and higher soil organic matter, in this region.

Table 1. Statistic characters of soil organic matter in different scales

Scales	Landscape regions	Space /m	Min /g kg ⁻¹	Max /g kg ⁻¹	Mean /g kg ⁻¹	Std.D eviation	Var.coef ficient/ %
Small scale	Purple Soil Zone	200	9.76	30.84	15.72	4.71	29.92
	Quartz Sandstone Zone in middle mountain	300	6.84	23.61	15.31	3.46	22.61
	Proluvial fan in mountain front	300	18.60	33.35	25.26	4.25	16.81
	Diluvium Zone in low mountain	500	7.23	35.20	18.70	5.74	30.69
Large scale		2000	5.26	35.60	19.56	7.23	36.96

In large-scale, the average level of soil organic matter was 19.56 g kg⁻¹ within the range of the mean of four landscape regions; coefficient of variation(36.96%) was greater than any landscape region, as well as standard deviation, it showed soil organic matter of large-scale samples which covered entire study area had a greater spatial variability than small-scale one.

3.2 Scale Effect of the Trend Model

The correlation between elevation and Soil organic matter of 92 large-scale samples was significant, not only that, the correlation was significant also when computing small-scale samples alone, or calculating small-scale samples together with large-scale samples (table 2). The significant correlation provided a basis for using regression kriging method to predict soil organic matter and proved that the trend could be simulated with three possibilities. The correlation coefficient of elevation and “large-scale” data, “small-scale” data and “large-scale + small-scale” data was -0.461, -0.163 and -0.278, showed the trend of soil organic matter depended on elevation factor was a macro one that could be better characterized with 2km spacing large-scale samples. Adding into small-scale samples could not enhance a trend fitting, though increasing the number of samples.

Table 2. Correlation between soil organic matter and elevation

	Large scale	Small scale	Large scale + small scale
elevation	-0.461(**)	-0.163(**)	-0.278(**)

** P < 0.01 Correlation is significant at the 0.01 level.

Trend equations of three dataset (“large-scale”, “small-scale” and “large scale + small scale”) depended on elevation factor expressed respectively as Y_1, Y_2 and Y_3 (where X was elevation):

$$Y_1 = -0.0465 * X + 36.485 \tag{2}$$

$$Y_2 = -0.0197 * X + 25.500 \tag{3}$$

$$Y_3 = -0.0331 * X + 30.944 \tag{4}$$

3.3 Scale Effect of Spatial Variation of Residual Value

Residual value, that represents random content of corresponding scale, is equal to measured values minus trend values. In the following, paper takes the residual value after removing Y_1 as a example for exploring the scale effect of spatial variation of residual. Semivariograms, with 2km lag size, of “large-scale” samples residual and “large-scale + small-scale” samples residual (fig. 2-1, fig. 2-2), showed the spatial structure of 92 “large-scale” samples residual was a pure nugget effect and was completely random, could not drawn spatial variability information within 2km distance; while the semivariogram of “large-scale + small-scale” samples residual reached a sill 43.00 $g^2 kg^{-2}$ at a distance of around 9600m, nugget coefficient was about 51.3%. And the nested features in multi-scale be showed clearly by the residual semivariogram (fig. 2-3) of “large-scale + small- scale” samples taking the minimum sampling interval 200m as lag distance. Obviously, the spatial variability of residual within 2km lag distance would be enhanced when adding “small-scale” data on the basis of “small-scale” data, then spatial autocorrelation analysis for data exploration with 2km lag distance became possibly.

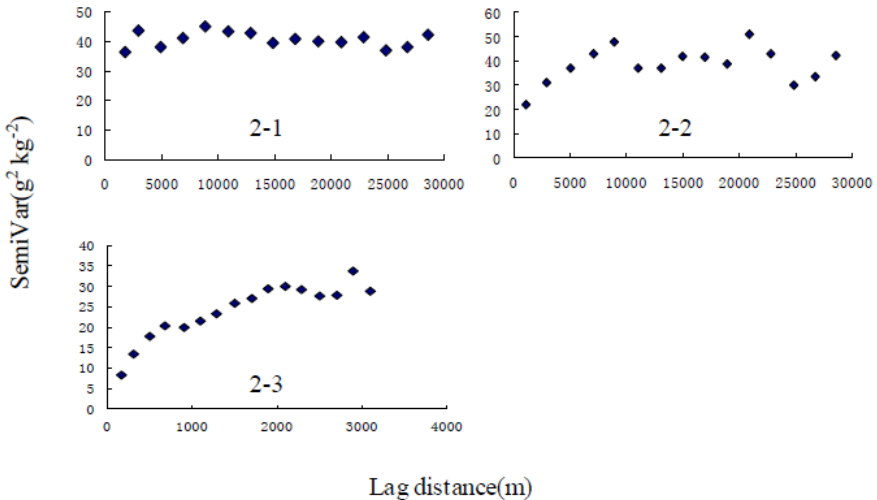


Fig. 2. Semi-Variogram of Residual after removing Y_1

Fitting residual semivariance function based multi-scale nested character was necessary due to the residual nested phenomenon of “large-scale + small-scale” samples. In the multi-scale nested structure, the lag distance in large scale was 2km, the sampling interval in large scale, furthermore, the one in small scale was 200m, the minimum sampling interval in small scale (Table 3). The nugget coefficient, in small scale, of “large-scale +small-scale” samples was 15.8%, the range in small scale was 950m and was 9649m in large scale. The nugget variance of nested semivariance was $3.00 \text{ g}^2 \text{ kg}^{-2}$, in addition, and semivariances in different scale satisfied spherical model.

Table 3. Parameters of the Residual Semivariance function of soil organic matter

Residual	Scales	Model	Range /m	C_0 / $\text{g}^2 \text{ kg}^{-2}$	C / $\text{g}^2 \text{ kg}^{-2}$	$C_0/(C_0+C)$ /%
Large scale	Large scale	Spherical model	9810	38.78	2.87	93.1
Large scale + Small scale	Small scale	Spherical model	950	3.00	16.00	15.8
	Large scale	Spherical model	9649	3.00	22.70	11.7

3.4 Influences of Changing Model on Predicting Outcomes

Predicting Outcomes. It was prediction map of soil organic matter that trend values estimated by the trend surface equation of $Y_1 \sim Y_3$ plused the corresponding residual forecast value.

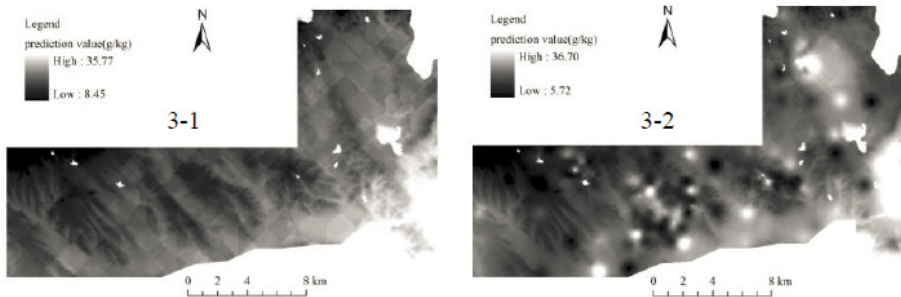


Fig. 3. Prediction map of soil organic matter with Y_1

Under the condition of Y_1 , when residual values were predicted with two datasets of “large-scale” samples and “large-scale + small-scale” samples respectively, soil organic matter prediction value was in the range of $8.45 \text{ g kg}^{-1} \sim 35.77 \text{ g kg}^{-1}$ and $5.72 \text{ g kg}^{-1} \sim 36.70 \text{ g kg}^{-1}$ (Fig. 3-1, Fig. 3-2); under Y_2 , was in the range of $11.50 \text{ g kg}^{-1} \sim 29.89 \text{ g kg}^{-1}$ and $6.46 \text{ g kg}^{-1} \sim 33.74 \text{ g kg}^{-1}$ (fig not shown); under Y_3 , was in the range of $10.12 \text{ g kg}^{-1} \sim 32.75 \text{ g kg}^{-1}$ and $6.20 \text{ g kg}^{-1} \sim 33.79 \text{ g kg}^{-1}$ (fig not shown). The minimum value and maximum value of measured values were 5.26 g kg^{-1} and 35.6 g kg^{-1} , the result could be

reached evidently comparing with prediction value ranges: compressed smooth effect of predictions was obvious when estimating residual value with large-scale samples only, no matter which trend surface equation was taken.

In order to identify the impact of different combinations of trend equations and residual simulations on results of soil predicting map, the prediction accuracy of estimating residual values with two datasets of “large-scale” samples and “large-scale + small-scale” samples respectively, under conditions of Y_1 , Y_2 and Y_3 , was compared (Table 4 ~ Table 6).

Table 4. Prediction accuracy of Y_1 + residual simulation with different dataset

Scales	Landscape regions	residual			
		Large scale		Large scale	scale+small
		RMSE	R	RMSE	R
		/g kg ⁻¹		/g kg ⁻¹	
Large scale		6.06	0.563**	0.88	0.997**
	Purple Soil Zone	5.82	-0.195	1.16	0.976**
	Quartz Sandstone Zone in middle mountain	4.32	0.121	1.27	0.985**
Small scale	Proluvial fan in mountain front	7.26	0.187	0.92	0.985**
	Diluvium Zone in low mountain	5.35	0.420	1.12	0.992**

Note: RMSE means the root mean square error, R means the correlation coefficient, it is same in the following.

* $P < 0.05$ Correlation is significant at the 0.05 level.

** $P < 0.01$ Correlation is significant at the 0.01 level.

Table 5. Prediction accuracy of Y_2 + residual simulation with different dataset

Scales	Landscape regions	residual			
		Large scale		Large scale	scale+small
		RMSE	R	RMSE	R
		/g kg ⁻¹		/g kg ⁻¹	
Large scale		8.31	0.429**	2.34	0.904**
	Purple Soil Zone	8.09	-0.320	2.51	0.889**
	Quartz Sandstone Zone in middle mountain	6.34	0.060	2.70	0.869**
Small scale	Proluvial fan in mountain front	9.22	0.139	2.22	0.902**
	Diluvium Zone in low mountain	7.51	0.217	2.59	0.885**

Influence of Changing Model. Table 4 ~ Table 6 showed that trend model had a impact on the prediction of large-scale soil organic matter. For prediction accuracy of large scale soil organic matter, Y_1 was the best trend equation when residual value was estimated with “large-scale” samples and the root mean square error (RMSE) was equal to 6.06 g kg⁻¹; Y_1 was still the optimal one when residual value was estimated

with “large-scale + small-scale” samples along with $RMSE = 0.88 \text{ g kg}^{-1}$; furthermore, Y_2 was the not most ideal one, no matter residual value interpolated with “large-scale” samples or “large-scale + small-scale” samples. The trend model had a impact on the prediction of small-scale soil organic matter as well. In small scale, measured value and prediction value was not correlated significantly when residual value was predicted with “large-scale” samples alone, under conditions of Y_1 , Y_2 or Y_3 . The main reason was “large-scale” samples could not display spatial variability of small scale soil. Remarkably, the best trend equation was Y_1 , the second one was Y_3 and the worst one was Y_2 , regardless of estimating residual value based which dataset, when prediction was effectively. For example, in the landscape of Purple Soil Zone, $RMSE$ of Y_1 , Y_2 and Y_3 was 1.16 g kg^{-1} , 2.51 g kg^{-1} and 1.44 g kg^{-1} respectively when residual calculated with “large-scale + small-scale” samples.

Table 6. Prediction accuracy of Y_3 + residual simulation with different dataset

Scales	Landscape regions	residual			
		Large scale		Large scale+small scale	
		RMSE /g kg ⁻¹	R	RMSE /g kg ⁻¹	R
Large scale	Purple Soil Zone	6.25	0.521**	1.24	0.995**
	Quartz Sandstone Zone in middle mountain	5.98	0.254	1.44	0.963**
Small scale	Proluvial fan in mountain front	4.33	0.086	1.61	0.973**
	Diluvium Zone in low mountain	7.40	0.035	1.18	0.975**
		5.42	0.377	1.52	0.985**

Table 7. Comparison of the mapping precision in small scale zone between using small scale data only and using regression Kriging under the conditions of Y_1 or Y_3

Landscape regions	Using small scale data only		Y_1	Y_3
	RMSE /g kg ⁻¹	R	RMSE /g kg ⁻¹	RMSE /g kg ⁻¹
Purple Soil Zone	1.92	0.930**	1.16	1.44
Quartz Sandstone Zone in middle mountain	2.70	0.697**	1.27	1.61
Proluvial fan in mountain front	2.09	0.898**	0.92	1.18
Diluvium Zone in low mountain	4.46	0.655**	1.12	1.52

The above analysis made it clear that Y_1 described the trend value of soil organic matter depending on elevation factor effectively, however, Y_3 was weaker to express the trend value due to the addition of “small-scale” samples, and Y_2 was weakest to exhibit the trend owing to the trend equation fitted with “small-scale” samples only.

Adding small-scale residuals improved the prediction accuracy of large-scale data significantly, under the condition of Y_1 , for example, $RMSE$ was reduced to 0.88 g kg^{-1} and correlation coefficient (R) between measured value and predicted value was raised to 0.997, calculating the residual value with “large-scale + small-scale”

samples. Whereas, RMSE was 6.06 g kg^{-1} and R was 0.563 predicting residual value with “large-scale” samples. Y_2 and Y_3 also showed the same law.

In small scale, prediction accuracy had been improved if comprehended trend and residual, compared with using their own samples in various small-scale landscape regions, that was represented most obviously in both regions of diluvium zone in low mountain and proluvial fan in mountain front (Table 7).

Table 8. Comparison of the range estimating in small scale zone between using small scale data only and using regression Kriging under the condition of Y_1

Landscape regions	Using small scale data only		Y_1	
	Min /g kg ⁻¹	Max /g kg ⁻¹	Min /g kg ⁻¹	Max /g kg ⁻¹
Purple Soil Zone	12.16	25.42	11.65	27.91
Quartz Sandstone Zone in middle mountain	12.94	17.55	9.33	27.36
Proluvial fan in mountain front	19.93	30.5	19.21	32.70
Diluvium Zone in low mountain	14.18	24.25	10.11	32.07

Further analysis indicated that predictions with their own samples in various small-scale landscape regions narrowed the range of measured value markedly because of lacking of macro trend control, and the range of prediction value was always less than the one under the control of macro trend value (under the condition of Y_1 , for example, Table 8)

4 Conclusions

In a complex landscape region, soil attribute trend value related to elevation factor is a macro trend. It is suitable for fitting with large-scale data. Although adding small-scale samples in the trend fitting process increases the number of samples, it increases the interference of model fitting as well, so the prediction accuracy of soil attribute not be improved. The optimal way to predict soil attribute, in a complex region, is that: fitting the trend value with large-scale samples alone and expressing the nested spatial variability of residual value.

Although small-scale sample does not enhance the simulation of macro trends, it reveals the spatial variability details of the residuals in a smaller scale, and it enhances the spatial structure of soil property residuals in large-scale, which improve the prediction accuracy of large-scale soil property.

Predicting with their own samples in each small-scale landscape region, lacking the large-scale trends in control, will result in the smoothing compression of the measured data, and will reduce the forecast accuracy of small-scale soil property.

Acknowledgments. Funding provided by the Natural Science Foundation of China (No. 40801080). We gratefully thank the anonymous reviewers for their constructive and valuable advices for the paper.

References

1. Burrough, P.A.: Multiscale sources of spatial variation in soil: I. The application of fractal concepts to nested levels of soil variation. *Journal of Soil Science* 34, 577–597 (1983)
2. Bourennane, A., Salvador-Blanes, S., Cornu, S., et al.: Scale of spatial dependence between chemical properties of topsoil and subsoil over a geologically contrasted area (Massif central, France). *Geoderma*. 112, 235–251 (2003)
3. Xu, Y., Chen, Y.X., Shi, H.B., et al.: Scale effect of spatial variability of soil water-salt. *Transactions of the CSAE* 20(2), 1–5 (2004)
4. Stenger, R., Priesack, E., Beese, F.: Spatial variation of nitrate-N and related soil properties at the plot-scale. *Geoderma*. 105(3-4), 259–275 (2002)
5. Rastetter, E.B., King, A.W., Cosby, B.J.: Aggregating fine-scale ecological knowledge to model coarser-scale attributes of ecosystems. *Ecological Applications* 2, 55–70 (1992)
6. Oliver, M.A.: Some novel geostatistical application in soil science. In: *Geostatistical Methods: Recent Development and Application in SSH. HIP-IV, UNESCO, Paris* (1992)
7. Sylla, M., Stein, A., Van Breemen, N., et al.: Spatial variability of soil salinity at different scales in the mangrove rice agro-ecosystem in West Africa. *Agriculture, Ecosystems & Environment* 54(1-2), 1–15 (1995)
8. Cerri, C.E.P., Bernoux, M., Chaplot, V., et al.: Assessment of soil property spatial variation in an Amazon pasture: basis for selecting an agronomic experimental area. *Geoderma*. 123(1-2), 51–68 (2004)
9. Yang, Y.L., Sheng, J.D., Tian, C.Y., et al.: A study on relationship between the spatial variability of saline anthropogenic alluvial soil available nitrogen, phosphorus, potassium and cotton growth. *Scientia Agricultura Sinica* 36, 542–547 (2003) (in Chinese)
10. Welsh, J.P., Wood, G.A., Godwin, R.J., et al.: Developing strategies for spatially variable nitrogen application in cereals, Part II: wheat. *Biosystems Engineering* 84, 495–511 (2003)
11. Li, D.Y., Pan, G.X., Chen, L.S., et al.: Spatial distribution and variability of Topsoil Organic Carbon content at different scales in Liuan city, Anhui Province, China. *Journal of Ecology and Rural Environment* 24(4), 37–41 (2008)
12. Liu, Q., Sun, J.K., Chen, Y.P., et al.: Spatial variability of the soil heavy metal with different sampling scales. *Chinese Journal of Soil Science* 40(6), 1406–1410 (2009)
13. Zhao, H.X., Li, B., Liu, Y.H., et al.: The soil properties along landscape heterogeneity on different scales in Huangfuchuan watershed. *Acta Ecologica Sinica* 25(8), 2010–2018 (2005)
14. Zhao, J., Liu, H.J., Du, L.J., et al.: Analysis for spatial heterogeneity of organic matter content and available nutrients in blacksoil crop area with different scales. *Journal of Soil and Water Conservation* 20(1), 41–44 (2006)
15. Yu, J., Zhou, Y., Nie, Y., et al.: Spatial variability of soil nitrogen in different scales and nested simulation. *Scientia Agricultura Sinica* 40(6), 1297–1302 (2007)
16. Huo, X.N., Li, H., Zhang, W.W., et al.: Multi-scale spatial structure of heavy metals in Beijing cultivated soils. *Transactions of the CSAE* 25(3), 223–229 (2009)
17. Lu, R.K.: *Chemical analysis of soil in agriculture*. China Agricultural Science and Technology Press, Beijing (2000)

New Machine Learning Algorithm: Random Forest

Yanli Liu, Yourong Wang, and Jian Zhang

Basic Teaching Department, Tangshan College, Tangshan Hebei 063000, China
ly17937@126.com, yourong1214@163.com, zhjian8765@yahoo.com.cn

Abstract. This Paper gives an introduction of Random Forest. Random Forest is a new Machine Learning Algorithm and a new combination Algorithm. Random Forest is a combination of a series of tree structure classifiers. Random Forest has many good characters. Random Forest has been widely used in classification and prediction, and used in regression too. Compared with the traditional algorithms Random Forest has many good virtues. Therefore the scope of application of Random Forest is very extensive.

Keywords: random forest, accuracy, generalization error, classifier, regression.

1 Introduction

The traditional machine learning algorithms usually give low classifier accuracy, and easy got over-fitting. To improve the accuracy, many people research on the algorithm of combining classifiers. Many scholar start the research on improve the classification accuracy by means of combining classifiers. In 1996, Leo Breiman advanced Bagging algorithm which is one of the early stage algorithm [1]. Amit and Geman define a large number of geometric features and search over a random selection on these for the best split at each node[2]. In 1998, Dietterich put forward the random split selection theory[3]. At each node the split is randomly selected from the N best splits. Ho[4] has done much study on “the random subspace” method which grows each tree by a random selection of a subset of features. Breiman [5]generate new training sets by randomizing the outputs in the original training set. Among these, the idea, in Amit and Geman’s paper, influenced Breiman’s thinking about random forests.

Random forests are a combination machine learning algorithm. Which are combined with a series of tree classifiers, each tree cast a unit vote for the most popular class, then combining these results get the final sort result. RF posses high classification accuracy, tolerate outliers and noise well and never got overfitting. RF has been one of the most popular research methods in data mining area and information to the biological field. In China there are little study on RF, so it is necessary to systemic summarize the down to date theory and application about RF.

2 The Principle of Operation and Characters of Random Forest

2.1 Principle of Operation

2001, Leo Breiman definite random forests as:

Definition 2.1 A random forest is a classifier consisting of a collection of tree-structured classifiers $\{h(x, \Theta_k), k = 1, \dots\}$ where the $\{\Theta_k\}$ are independent

identically distributed random vectors and each tree casts a unit vote for the most popular class at input x .

This definition show RF is a combination of many tree-structure classifiers. In Breiman’s RF model, every tree is planted on the basis of a training sample set and a random variable, the random variable corresponding to the k th tree is denoted as Θ_k , between any two of these random variables are independent and identically distributed, resulting in a classifier $h(x, \Theta_k)$ where x is the input vector. After k times running ,we obtain classifiers sequence $\{h_1(x), h_2(x), \dots, h_k(x)\}$, and use these to constitute more than one classification model system ,the final result of this system is drawn by ordinary majority vote, the decision function is

$$H(x) = \arg \max_Y \sum_{i=1}^k I(h_i(x) = Y) \tag{1}$$

where $H(x)$ is combination of classification model, h_i is a single decision tree model, Y is the output variable, $I(\cdot)$ is the indicator function. For a given input variable, each tree has right to vote to select the best classification result. Specific process shown in Fig. 1.

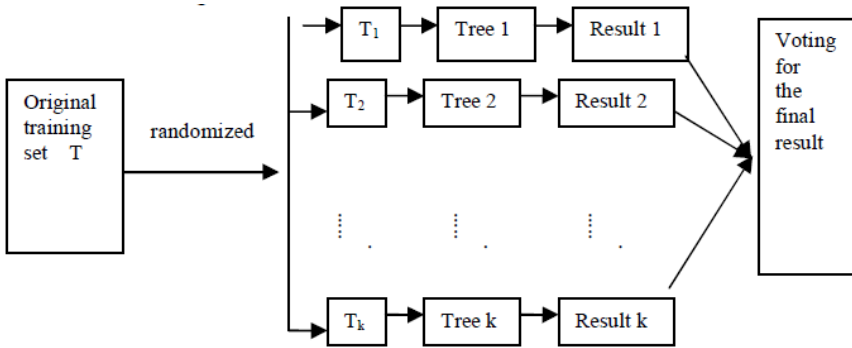


Fig. 1. Random forest schematic

2.2 Characters of Random Forest

In Random Forest, margin function is used to measure the extent to which the average number of votes at X, Y for the right class exceeds that for the wrong class, define the margin function as:

$$mg(X, Y) = av_k I(h_k(X) = Y) - \max_{j \neq Y} av_k I(h_k(X) = j) \tag{2}$$

The larger the margin value, the higher accuracy of the classification prediction, and the more confidence in classification.

Define the generalization error of this classifier as:

$$PE^* = P_{X,Y}(mg(X, Y) < 0) \tag{3}$$

when the number of decision tree is big enough, $h_k(X) = h(X, \Theta_k)$ obey the Strong Law of Large Number. Leo Breiman has proved two conclusions. One is RF do not over-fitting but really produce a limiting value of the generalization error. The reason is with the number of the decision trees increases, for almost surely all sequences Θ_1, \dots PE* converges to

$$P_{X,Y}(P_\theta(h_k(X, \theta) = Y) - \max_{j \neq Y} P_\theta(h(X, \theta) = j) < 0) \tag{4}$$

Another is the upper bound of the generalization error is exist, and

$$PE^* \leq \bar{\rho}(1 - s^2) / s^2 \tag{5}$$

s is the strength of the set of classifiers $\{h(\mathbf{x}, \theta)\}$, $\bar{\rho}$ is the mean value of the correlation [6]. It shows that the generalization error of RF depends on two aspects: one is the strength of the individual trees in the forest, another is the correlation between these trees. Obviously, the smaller this value, the better the results of random forest.

2.3 Out-of-Bag Estimation

In the process of constructing RF, the tree is planted on the new training set by using random features selection, the new training set is drawn from the original training set by bagging methods. There are two reasons for using bagging. The first is that the use of bagging seems to enhance accuracy when random features are used. The second is using bagging will bring out-of bag data, which can be used to give ongoing estimates of the PE*of RF, as well as estimates for the strength and correlation.

Given an original training set T with N samples, the kth training set is drawn from T with replacement by bagging, every T_k contains N samples. Then the probability of each sample can not be contain is $(1 - 1/N)^N$, when N large enough, $(1 - 1/N)^N$ is converges to e^{-1} . In other words, 36.8% samples of the T is not contained in T_k . This samples is called out-of -bag data. The algorithm of using these data to estimate the performances of classification is called OOB estimation. For each tree, there is an OOB estimate for its error. The estimation of generalization error of RF is the average of estimations of all tree error for every tree contained in the RF. Compared with cross-validation the OOB estimate is unbiased and runs faster. The accuracy of OOB estimate is favorable to cross-validation. Tibshirani, Wolpert and Macready Proposed using OOB estimate as an ingredient in estimates of generalization error [7-8]. Breiman has proofed the out-of-bag estimate is as accurate as using a test set of the same size as the training set .Therefore, using the OOB error estimate removes the need for a test set aside[9].Strength and correlation can also be estimated using out-of -bag methods. This gives an internal estimate what is helpful in understanding classification accuracy and how to improve it.

3 The Methods of Random Forests Construction

There are many methods to construct RF, for example bagging method, using random input selection, the effects of output noise, etc.

3.1 Using Input Variables to Construct R.F.

There are three methods to construct R.F. by using input variables, Forests-RI, Forests-RC and Categorical Variables. Because the mechanism of Categorical Variables is complex and the strength is not much better than other RF, so we won't explain this method here, and only introduce the other two methods.

Forest-RI is the simplest RF with random features. Forest-RI is formed by randomly selecting a small group of input variables at each node to split on. F the size of the group is fixed. Using CART methodology to plant tree, maximum size and do not prune. In Breiman's experiment two values of F were tried. One is F=1, another is the first integer less than \log_2^{M+1} , M is the number of inputs. The accuracy of Forest-RI is favorable with Adaboost. Forest-RI can be much faster than both Adaboost and Bagging. And the procedure is not overly sensitive to the value of F. It is surprising that when F=1, the procedure has good accuracy too.

When there are a few inputs, M is not big, taking F inputs from all as random selection might lead an increase in strength but higher correlation too. Defining new feature by random linear combination of specifying L input variables. Then there are much enough features. At each given node, L variables are randomly selected and

added together with coefficients $k_i, v = \sum_{i=1}^L k_i v_i, k_i \in [-1, 1]$. F linear combinations are

generated, and the best split can be found over these. We call this procedure Forest-RC. Breiman's study show Forest-RC has merits: 1) Forest-RC can deal with data set contain incommensurable input variables; 2) On the synthetic data sets Forest-RC does exceptionally well; 3) Compared with Forest-RI, the accuracy of Forest-RC is more favorably to Adaboost.

3.2 Using Output Construct Random Forest

There are two methods to construct RF using output. One is output smerring, putting Gauss noise in the procedure of output. Another is output flipping, changing one or several classifying labels of the output [10]. In this procedure, the variable remained relatively the same in the classification section is very important. The most obvious virtue of this idea is the RF process the ability of estimating the importance of each feature. The RF constructed by this method can be used to regression well as classification, and better than Bagging in strength. But output flipping depend on the selection of flip rate.

Using updated the weight to built RF [11-12], the merits of this idea is easy and run faster, and easy to realize by program. But this method obviously relay to the data itself and week learning, and can be easy influenced by noises. SRF is built by

randomly selected feature subspace [13-14]. In the given sample space, using this idea you can built as many tree as you want. The strength is much better than tree. With the complexity of the construct, the overall accuracy is almost monotonically increasing. SRF accuracy for multi-tree is optimal.

4 Random Forests for Regression

Random forest can be used to regression too. Specific construction method of RF regression model can be found in [15]. RF regression model can be briefly summarized as: given sample space x and classification labels Y , random forests for regression are formed by planting trees depending on the random variable Θ , relative to each category label, tree predictor $h(x, \Theta)$ can give a numerical result. The random forest predictor is formed by taking the average over k of the trees $h(x, \Theta_k)$. Similarly to the classification case, the following holds:

Lemma 4.1 As the number of trees in the forest goes to infinity, almost surely,

$$E_{x,y}(Y - av_k h(X, \theta_k))^2 \rightarrow E_{x,y}(Y - E_\theta(X, \theta_k))^2 \tag{6}$$

Random forests regression function is $Y = E_\theta(X, \theta_k)$. In practice, when k big enough $Y = av_k h(X, \theta_k)$ is usually used to instead of the regression function. Breiman has proofed the conclusion that assume for all $\Theta, E_y = E_x h(X, \theta)$ then

$$PE^*(forest) \leq \bar{\rho} PE^*(tree) \tag{7}$$

This pinpoints that low correlation between residuals and low error trees can give high accurate regression forest. To test effect of the RF regression, compare this regression with SVR [16] and linear regression. Do regression on the data set CPU.arff (Weka’s data set), the resulting parameters are shown in Table 1.

Table 1. The results of three regression models

Regression model name Parameter	RFR	SVR	Linear R
Correlation coefficient	0.9613	0.9398	0.9544
Mean absolute error	13.0878	19.7969	32.1855
Root mean squared error	50.3600	62.0144	46.0993
Relative absolute error	14.9775%	22.6554%	36.8327%
Root relative squared error	32.6195%	40.1683%	29.8597%

The results show that random forest regression better than the other two regression models. RF can deal with numerical data and data, but the other two only can deal with numerical variables and continuous variables but the other two only can deal with numerical data. So RFR can be more widely applied. The results show that random forest regression better than the other two regression models.

The study about RFR is still going on. 2006, Quantile Regression Forest defined by Nicolai, is derived from random forests[17]. Nicolai has proved mathematically that Quantile Regression Forest is consistent. Quantile regression forests can be seen as one of the applicability of the nearest neighbor classification and regression process[18]. In addition Brenc and Brown improve the robust RF regression algorithm based on the to booming algorithm[19].

5 The Application of Random Forest

RF can be used to deal with micro-information data, and the accuracy of RF is higher than those traditional predictions. So in recently 10 years, Random Forest has been got a rapid development, and widely used in many areas, such as bioinformatics, medicine, management science, economics. In bioinformatics, Smith et al. studied the tracking data on bacteria by RF, and compared with Discriminant Analysis method. Alonso et al. use biomarkers parasite to discriminate fish stocks [20-21]; In medicine, Using RF technology such as Lee to help lung CT images of lung nodules automatic detection, and also in the RF (CAC)[22]. In China, Jia FuCang, Li Hua, Using RF to the Dhoop magnetic resonance image segmentation, and that the RF has fast speed and high accuracy, is a promising multi-channel image segmentation method[23]. The main application in economic management field, is predicating the loss degrees of customers. Bart used RF in customer relationship management, found that the effect of RF is better than ordinary linear regression and Logistic model[24]. Coussement et al. compared the predictive ability of SVM, logistic model and the RF in loss of customers, found that RF is always better than the SVM[25]. Burez et al. applied weighted RF in loss of customers, comparing with the RF, and found the weighted RF has better prediction[26].

Today, the range of application of RF is very broad, in addition to the above mentioned application, the RF also used in ecology[27-28], remote sensing geography terms[29-30], customer's loyalty forecasting[31]; Lessmann etc. also use Random Forest predict horse racing winning, and that the predictions of the Random Forest is superior to traditional forecasting methods can bring in huge commercial profits[32].

6 Conclusions and Outlook

In summary, the RF as a combination of the tree classifier is an effective classification predicting tool. It has the following advantages: 1) the accuracy of random forests is not less than Adaboost, run faster, and does not produce over-fitting. 2) the OOB data can be used to estimate the the RF generalization error, correlation and strength, can also estimate the importance of individual variables. 3) the combination of bagging and the random selection of features to split allows the RF to better tolerate noise. 4) RF can handle continuous variables and categorical variables.

Recently, the RF theory is more mature and the application range is becoming wilder. But there many work to do to further improve RF, and use RF to much wider fields. Hopping the interested scholars can do further research.

References

1. Breiman, L.: Bagging Predictors. *Machine Learning* 24, 123–140 (1996)
2. Amit, Y., Geman, D.: Shape quantization and recognition with randomized trees. *Neural Computation* 9, 1545–1588 (1997)
3. Dietterich, T.: An Experimental Comparison of Three Methods for Constructing Ensembles of Decision Trees: Bagging, Boosting and Randomization. *Machine Learning*, 1–22 (1998)
4. Ho, T.K.: The random subspace method for constructing decision forests. *IEEE Trans. on Pattern Analysis and Machine Intelligence* 20(8), 832–844 (1998)
5. Breiman, L.: Using adaptive bagging to debias regressions, Technical Report 547, Statistics Dept. UCB (1999)
6. Breiman, L.: Random Forests. *Machine Learning* 45(1) (2001)
7. Tibshirani, R.: Bias, Variance, and Prediction Error for Classification Rules, Technical Report, Statistics Department, University of Toronto (1996)
8. Wolpert, D.H., Macready, W.G.: An Efficient Method to Estimate Bagging's Generalization Error. *Machine Learning* (1997) (in press)
9. Breiman, L.: Out-of-bag estimation [EB/OL] (June 30, 2010), <http://stat.berkeley.edu/pub/users/Breiman/OOBestimation.ps>
10. Breiman, L.: Prediction Games and Arcing Algorithms. *Neural Computation* 11, 1493–1517 (1999)
11. Bauer, E., et al.: An Empirical Comparison of Voting Classification Algorithms: Bagging, Boosting, and Variants. *Machine Learning* 36, 105–142 (1999)
12. Freund, Y., Shapire, R.: Experiments with a new boosting Algorithm. In: *Machine Learning: Proceedings of the Thirteenth International Conference*, pp. 148–156 (1996)
13. Ho, T.K.: Random: Decision Forests. In: *Proceeding of the 3rd International Conference on Document Analysis and Recognition*, Montreal, Canada, August 14–18, pp. 278–282 (1995)

Efficient Method of Formal Event Analysis

Ying Liu¹ and Zongtian Liu²

¹Dept of Knowledge Service Engineering, Korea Advanced Institute of Science and Technology

²School of Computer Engineering and Science Shanghai University, Shanghai, 200072, China
ztliu@shu.edu.cn

Abstract. Based on the philosophic idea about static concept, Professor R.Wille in Germany had provided the theory about formal concept analyses firstly in 1982, which opens out inwardness of static concept and relation between the concepts well. But the theory has evident limitation for representing dynamic event. The paper extends the theory to be a more generic model. On the basis of the definition of event proposed by us, more formalized description of event is given, a method to build formal event lattice from formal event context is suggested. Then the feasibility of the method is demonstrated from theory and instance. The model of Formal Event lattice repairs the limitation of formal concept lattice, so it can characterize the essence of the dynamic world very well. It will exert effect in the application in the fields of network resource management, Web servers, event ontology, robot control, marketing strategy, data mining and etc.

Keywords: Formal concept analysis, concept lattice, Formal event analysis, Event Lattice.

1 Introduction

1.1 Event Is the Basic Cell Making Up the Knowledge about the World

Most of philosophers believe that the world is of material and consists of objects and events [1]. Event is a concrete fact which does not go back, while object is a thing which possesses invariable particularity.

Many cognitive scientists research on event by simulating the cognitive process of human and exploring its semantic structure, such as Zacks of Stanford University regarded events as a stream of activities perceived by observers [2]

Some of linguists have given definition and structure of event. For example, some scientists suggested using the triple, SVO (Subject-Verb-Object), to represent an event, in which Subject corresponds to the initiator of the event, Verb to the activity and Object to the accepters [3].

In the field of Artificial Intelligence, event is defined for knowledge representation and information extraction. The main attention is about structure, reasoning, and recognition of events. For example, Nelson thought that an event is a larger whole incorporating objects and relations, and given an event representation model GER (Generalized event representations) through combining script-based knowledge expression method [4].

WordNet. is an excellent English new-style electronic dictionary designed by some psychologists, linguists and computer engineers in Princeton University, based on cognition philology, in which an event is defined as a something that happens at some place and some time. But “the tennis problem”, which Roger Chaffin point out, reflects its shortcoming in event knowledge representation. HowNet developed by Professor Zhendong Dong in Chinese Academy science is a common-sense knowledge base of Chinese and their English equivalents, in which event is identical with the concept “fact” and can be broadly classified as static or dynamic. But there is still “the tennis problem”.

In the fields of Information Retrieval, Information Extraction and Automatic Summarization, researchers also attach importance of the knowledge about event

Though the comprehending about event does not go all the way, the viewpoint to cognize event as a basic cell of knowledge is consistent.

1.2 The Limitation of Formal Concept Analysis for Event Analysis

Object is silent relatively. The set of objects that possess common attributes is called Concept. The inclusion relation between concepts determines the hierarchy of the concepts. Theoretically, this is a complete lattice.

By the philosophic thoughts of the concepts, Professor Wille in Germany proposed the theory of formal concept analysis firstly [5]. Following researches develop the theory, provide some constructing algorithms of concept lattices [6] [7], build fuzzy Concept Lattice Model [8] [9], and so on.

Concept lattice is a basic data structure of formal concept analysis. Each node of it is a formal concept, which consists of two parts: extension and intension. Extension is a set of objects in the concept, while intension is a set of common attributes of the objects.

Event is motional, perceptible, discrete and differ from static concept.

As dynamic characteristic of event, it is difficult for the method of Wille to be directly applied in formal event analysis, the essential reason is that attributes of object may be represented with a group of Boolean value, while those of event is more complex.

1.3 Contribution and Organization of This Paper

The paper proposes a formal event analysis method. According to the peculiarity of event, the formal concept analysis method is extended to be suitable for formal event.

As formal concept analysis, formal event analysis is an important theory, and may be applied in many fields. The set of events which possess common attributes make up an event class. The event classes and inclusion relation between them compose a complete lattice, called a formal event lattice, for short, event lattice.

2 The Definitions of Event and Event Class

The definitions of event and event class have been given by us in [10] as follows:

Definition 1. Event e can be defined as a 6-tuple formally:

$$e ::=_{def} \langle A, O, T, V, P, L \rangle$$

There A,O,T,V,P,V are six Factors of the event. Means actions happen in the event. It describes the process of event happens. These actions executed in sequence or concurrently while the event happens. O means objects taking part in the event, including all actors and entities involved in the event. We define two types of the objects, action initiators and action acceptor. T means the period that event lasting. The time period can be expressed as absolute time and relative time. V means environment of event, including nature environment and social environment, such as location and background of event. P means assertions on the procedure of actions execution in the event. Assertions include pre-condition, post-condition and intermediate assertions. Pre-condition represents the state that has to be satisfied for triggering the event. Post-condition represents the result states after event happens. Intermediate assertions represent some intermediate states during the event. L means language expressions, including Core Word Set, Core Words Expressions and Core Words Collocations. Core Words are high-frequency words in sentences of the event. Core Word Expressions describe the position relationships between no-core factors words and core words. Core Word Collocations denote the fixed collocations between core words and other words.

```

/*-----
Event class name: procreate
Objects: Role-1: Animal, initiator
         Role-2: Animal, acceptor
         Role-3: Human, acceptor
Action:  Role-1 is having a delivery, as well as Role-2
         is coming in the world, Role-3 is seeing after
         Role-1 and Role-2 .
         Degree: difficult.
Time:    T,T+Δ, Δ= a few hours
Environment: ground
Predicate: Rpre-condition: Role-1 is an adult female,
         Role-1 lives in body of Role-1, ole-1 and
         Role-2 are same class. Post-condition:
         Role-1 lives out body of Role-2, Role-2 is
         infantile has apperceived Role-2
Language: “生产”、“生育”、“下崽”、“procreate”、
         “inbreed” 。
-----*/

```

Fig. 1. An illustration of an event class

Definition 2. Event Class means a set of events with common characteristics, defined as

$$EC = (E, C_1, C_2, \dots, C_6)$$

$$C_i = \{c_{i1}, c_{i2}, \dots, c_{im}, \dots\} \quad (1 \leq i \leq 6, m \geq 0)$$

There E is an event set, called extension of the event class. Ci called intension of the event class. It denotes the common characteristics set of the event class. C_{im} denotes one of the common characteristics of event factor i.

An illustration of event class is shown in Fig. 1.

3 Extension of Formal Concept Analysis

The content of the definitions 3-6 are from reference [5], but the descriptions of them have been modified by us.

Definition 3. A triple (U,A,I) is called a formal context, if U is a set of objects, A is a set of attributes, and $I \subseteq U \times A$ is a binary relation between U and A, $(o,d) \in I$ means the object o possesses attribute d.

Definition 4. For a formal context (U,A,I), a pair of operators, f and g, for any $X \subseteq U$ and $B \subseteq A$ can be definite by

$$f(X) = \{ a \in A \mid (x,a) \in I, \forall x \in X \},$$

and

$$g(B) = \{ x \in U \mid (x,a) \in I, \forall a \in B \},$$

Where f(X) is the set of attributes shared by all the objects in X, and g(B) is the set of objects possessing all the attributes in B.

Definition 5. Let (U,A,I) be a formal context. A pair (X,B) is called a formal concept, for short, concept of (U,A,I), if and only if $X \forall U, B \subseteq A, f(X)=B$ and $X=g(B)$.

In which, X is called the extension and B the intension of the concept (X,B).

Definition 6. All the concepts from a formal context can form a complete lattice by extension or intension inclusion relation among them, that is called the concept lattice of (U,A,I) and is denoted by L(U,A,I).

But in actual problem, attributes are not only used to express that an object has or not the property, but, in more cases, are to express that an object may have one of the many values of an attribute. Such attributes are called many-valued attributes, For example, attribute color has attribute value “red”, “yellow”, “green” and so on.

Definition 7. A 4-tuple $K=(U,A,V, I)$ is called many-valued-formal-context, if U is a set of objects, A is a set of attributes, V is the value domain of A, and I is a ternary relation, $I \subseteq U \times A \times V$, there $(u,m,v) \in I$, or $I(u,a)=v$, means the value of object u on attribute m is v, $I(u)=v_1, v_2, \dots, v_n$, if $i \in A, v_i \in V_i$

In Wille’ theory, a many-valued attribute has to be changed into many binary-attributes. To be suitable of event analyses, we must break the limitation.

Definition 8. If the value domain of an attribute is a complete lattice, it is called complete lattice attribute.

Property 1.

(1) If the value domain of an attribute is binary numeral domain, then the attribute is a complete lattice attribute, because of binary numeral domain and relation $<$ on it is a complete lattice.

(2) If the value domain of every one of many attributes is binary numeral domain, then the product of the many attributes is a complete lattice attribute, because of product of binary numeral domains and relation $<$ on them is also a complete lattice. Actually, this is a binary denotation of a traditional finite set.

(3) If the value domain of an attribute is a power set of a limited set A , then the attribute is a complete lattice attribute, because of a power set and relation \subseteq on it is a complete lattice.

(4) If the value domain of an attribute is some limited interval numerals and union or intersect of any subset of them, then the attribute is a complete lattice attribute, because of the domain and relation \subseteq on it is a complete lattice.

(5) If the value domain of an attribute is first order predication formula domain, then the attribute is a complete lattice attribute, because of first order predication formula domain and relation \Rightarrow on it is a complete lattice.

(6) If the value domain of an attribute is an union set of two sets, the first set consists of some limited prime numbers including one and the second set consists of all the Least Common Multiples of any subset of the first set, then the attribute is a complete lattice attribute, because of the union set and multiple relation on it is a complete lattice.

(7) If the value domain of each one of many attributes is a complete lattice, then the product of the attributes is a complete lattice attribute, because of a product of many complete lattices is still a complete lattice.

Definition 9. A many-valued-formal-context $K=(U,A,V, I)$ is called complete lattice formal context, if and only if, for every $a \in A, \forall a \in V$ is a complete lattice value domain.

Definition 10. Let's define two mapping, f and g , on complete lattice formal context $K=(U,A,V,I)$:

$$\begin{aligned} \forall O \subseteq U: (f(O) = \bigwedge \{I(o) \mid o \in O\}) \\ \forall v \in V: (g(v) = \{o \mid o \in U \text{ and } v \leq I(o)\}) \end{aligned}$$

Definition 11. For a complete lattice formal context $K=(U,A,V,I)$, a pair (O,v) is called a class, if $O \subseteq U, v \in V, f(O)=v$ and $g(v)=O$, where O is referred as extension and v intension of the class.

Definition 12. Two classes, $C_1=(O_1,v_1)$ and $C_2=(O_2,v_2)$ is called $C_2 \leq C_1$, if $O_2 \subseteq O_1$, i.e. $v_1 \leq v_2$.

Definition 13. The all classes of K as well as the relations \leq among them form a complete lattice, called class lattice, denoted by $L(K)$.

For a many-valued-formal-context $K=(U,A,V,I)$, if the value domain of every one of many attributes is binary numeral domain, then the class lattice generated from the context is a traditional formal concept lattice[5].

For a many-valued-formal-context $K=(U,A,V, I)$, if the value domain of A is some limited interval numerals and union or intersect of any subset of them, then the class lattice generated from the context is an interval numeral formal concept lattice [11].

For a many-valued-formal-context $K=(U,A,V, I)$, if the value domain of the attribute is first order predication formula domain, then the class lattice generated from the context is a first order predication formula formal concept lattice[12][13].

4 Formal Event Analysis

We re-describe event of definition 1 by types of attribute value domains. Actions A in definition 1 can be denoted by an attribute Δ which value domain is a power set of a limited set, an element of the set represents the means that the event has the action attribute. Objects O by an attribute Θ which value domain is a product of many concept lattices. The period T by an attribute T which value domain is an Interval number domain. Environment V, assertions P and Language L by an attribute ξ which value domain is a product of two first order predication formula domains, representing the initiatory state and the end state of the event respectively So that, we definite a formal event as follows.

Definition 14. A formal event is denoted by a 4-tuple

$$fe=(\Delta, \Theta, T, \xi)$$

Definition 15. The information table composed of many formal events as many rows is called formal event context.

An example of formal event context is shown in Table 1.

Theorem 1. A formal event context is a complete lattice formal context.

Prove: According to property 1, the value domains of Δ, Θ, T, ξ are all complete lattice value domains, so the product of them is also complete lattice value domain, the formal event context is a complete lattice formal context.

Definition 16 the lattice generated from a formal event context by definition 10-13 is called formal event lattice, for shot, event lattice.

Table 1. A formal event lattice

Δ	Θ	T	ξ
e1 {slow,curve}	O=animal, $ O \geq 2$	[3,8]	$\forall o \in O: at(o,x) \wedge on(x, earth);$ $\exists o \in O: at(o,y) \wedge on(y, earth) \wedge x \neq y$
e2 {quick,curve }	O=buman, $ O \geq 2$	[4,5]	$\forall o \in O: at(o,x) \wedge on(x, earth);$ $\forall o \in O: at(o,y) \wedge on(y, earth) \wedge x \neq y$
e3 {slow,curve }	O=animal, $ O \geq 2$	[1,7]	$\forall o \in O: at(o,x) \wedge in(x, water);$ $\exists o \in O: at(o,y) \wedge on(y, water) \wedge x \neq y$
e4 {quick,beeline}	O=buman, $ O \geq 2$	[1,5]	$\forall o \in O: at(o,x) \wedge in(x, water);$ $\forall o \in O: at(o, y) \wedge on(y, water) \wedge x \neq y$

The event classes generated from it are:

(e1e2e3e4, {}, o ∈ animal, [1,8], $\forall o \in O: \text{at}(o,x) \wedge (\text{on}(x, \text{earth}) \vee \text{in}(x, \text{water})); \exists o \in O: \text{at}(o,y) \wedge (\text{on}(y, \text{earth}) \vee \text{in}(x, \text{water})) \wedge x \neq y$) // Animals move since 1 to 8

(e1e2e3; { curve }, o ∈ animal, [1,8], $\forall o \in O: \text{at}(o,x) \wedge (\text{on}(x, \text{earth}) \vee \text{in}(x, \text{water})); \exists o \in O: \text{at}(o,y) \wedge (\text{on}(y, \text{earth}) \vee \text{in}(x, \text{water})) \wedge x \neq y$) // Animals move along curve since 1 to 8

(e1e2; { curve }, o ∈ animal, [3,8], $\forall o \in O: \text{at}(o,x) \wedge \text{on}(x, \text{earth}); \exists o \in O: \text{at}(o,y) \wedge \text{on}(y, \text{earth}) \wedge x \neq y$) // Animals move on earth along curve since 3 to 8

(e1e3; { slow, curve }, o ∈ animal, [1,8], $\forall o \in O: \text{at}(o,x) \wedge (\text{on}(x, \text{earth}) \vee \text{in}(x, \text{water})); \exists o \in O: \text{at}(o,y) \wedge \text{on}(y, \text{earth}) \wedge x \neq y$) // Animals move on earth slowly along curve since 1 to 8

(e2e4, { quick }, o ∈ human, [1,5], $\forall o \in O: \text{at}(o,x) \wedge (\text{on}(x, \text{earth}) \vee \text{in}(x, \text{water})); \exists o \in O: \text{at}(o,y) \wedge (\text{on}(y, \text{earth}) \vee \text{in}(x, \text{water})) \wedge x \neq y$) // human move quickly since 1 to 5

(e3e4, {}, o ∈ animal, [1,7], $\forall o \in O: \text{at}(o,x) \wedge \text{in}(x, \text{water}); \forall o \in O: \text{at}(o,y) \wedge \vee \text{in}(x, \text{water}) \wedge x \neq y$) // Animals move in water since 1 to 7 not fall beband.

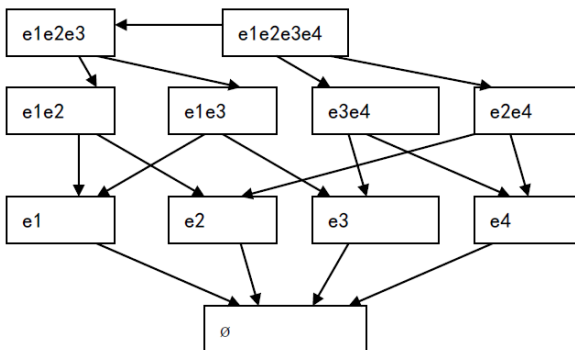
(e1; { slow, curve }, o ∈ animal, [3,8], $\forall o \in O: \text{at}(o,x) \wedge \text{on}(x, \text{earth}); \exists o \in O: \text{at}(o,y) \wedge \text{on}(y, \text{earth}) \wedge x \neq y$) // Animals move slowly on earth along curve since 3 to 8

(e2, { quick, curve }, o ∈ human, [4,5], $\forall o \in O: \text{at}(o,x) \wedge \text{on}(x, \text{earth}); \forall o \in O: \text{at}(o,y) \wedge \text{on}(y, \text{earth}) \wedge x \neq y$) // Animals move quickly on earth along curve since 4 to 5 not fall beband.

(e3, { slow, curve }, o ∈ animal, [1,7], $\forall o \in O: \text{at}(o,x) \wedge \text{in}(x, \text{water}); \exists o \in O: \text{at}(o,y) \wedge \vee \text{in}(x, \text{water}) \wedge x \neq y$) // Animals move in water along curve since 1 to 7

(e4, { quick, beeline }, o ∈ human, [1,5], $\forall o \in O: \text{at}(o,x) \wedge \text{in}(x, \text{water}); \forall o \in O: \text{at}(o,y) \wedge \vee \text{in}(x, \text{water}) \wedge x \neq y$) // Animals move in water along curve since 1 to 7, not fall be band.

The event lattice generated from K is follows.



5 Conclusion

The event-based idea has taken root into the many fields, such as Web servers, robot control, business management, communications, programming and etc. “Event-oriented”, “Event-driven” and “Event-based” have been as terminology

known by all. Formal event analysis theory will become the foundation of event research and application.

Similar to concepts, along with classifying relation, there are many kinds of non-classifying relations between events. How to research the non-classifying relations is one of our main tasks in future.

Acknowledgement. This work is supported by the projects of National Science Foundation of China (NSFC No.60575035, No.60975033), Shanghai Leading Academic Discipline Project (J50103) and Postgraduate Innovation Fund of Shanghai University (SHUCX091041).

References

1. Chen, X.: Why did John Herschel fail to understand polarization? The differences between object and event concepts. *Studies in History and Philosophy of Science* 34, 491–513 (2003)
2. Zacks, J.M., Tversky, B.: Event structure in perception and conception. *Psychological Bulletin* 127(1), 3–21 (2001)
3. Chang, J.: Event Structure and Argument Linking in Chinese. *Language and Linguistics* 4(2), 317–351 (2003)
4. Nelson, K., Gruendel, J.: *Event knowledge: structure and function in development*. Erlbaum, Hillsdale (1986)
5. Ganter, B., Wille, R.: *Formal Concept Analysis: Mathematical Foundations*. Springer, Berlin (1999)
6. Godin, R., Missaoui, R., April, A.: Experimental comparison of navigation in a Galois lattice with conventional information retrieval methods. *International Journal of Man-machine Studies* 38, 747–767 (1993)
7. Liu, Z.-T., Li, L.S., Zhang, Q.: Research on a Union Algorithm of Multiple Concept Lattices. In: Wang, G., Liu, Q., Yao, Y., Skowron, A. (eds.) *RSFDGrC 2003*. LNCS (LNAI), vol. 2639, pp. 533–540. Springer, Heidelberg (2003)
8. Alcalde, C., Burusco, A., Fuentes-Gonzalez, R., Zubia, I.: An application of the L-Fuzzy Concept Theory to the short-circuit detection. In: *IEEE International Fuzzy Systems Conference, FUZZ-IEEE 2007*, July 23–26, pp. 1–6 (2007)
9. Liu, Z.-T., Qiang, Y., Zhou, W., Li, X., Huang, M.-L.: A Fuzzy Concept Lattice Model and Its Incremental Construction Algorithm. *Chinese Journal of Computers* 30(2), 184–188 (2007) (in Chinese)
10. Liu, Z.-T., Huang, M.-L., et al.: Research on Event-oriented Ontology Model. *Computer Science* 36(11), 189–192, 199 (2009) (in Chinese)
11. Jaoua, A., Elloumi, S.: Galois connection, formal concepts and Galois lattice in real relations: application in a real classifier. *The Journal of Systems & Software* 60, 149–163 (2002)
12. Ferré, S., Ridoux, O.: A Logical Generalization of Formal Concept Analysis. In: Ganter, B., Mineau, G.W. (eds.) *ICCS 2000*. LNCS (LNAI), vol. 1867, pp. 371–385. Springer, Heidelberg (2000)
13. Bain, M.: Predicate Invention and the Revision of First-Order Concept Lattices. In: Eklund, P. (ed.) *ICFCA 2004*. LNCS (LNAI), vol. 2961, pp. 329–336. Springer, Heidelberg (2004)

Application of Monte Carlo Simulation in Reliability and Validity Evaluation of Two-Stage Cluster Sampling on Multinomial Sensitive Question

Qiao-qiao Du, Ge Gao^{*}, Zong-da Jin, Wei Li, and Xiang-yu Chen

School of Public Health, Medical College of Soochow University, 215123, Suzhou, China
gaoge@suda.edu.cn

Abstract. In this paper, randomized response technique (RRT) model was presented in application to investigating multinomial sensitive question with the sample selected by two-stage cluster sampling. Monte Carlo simulation was successfully performed in the assessment of reliability and validity of two-stage cluster sampling investigation on multinomial sensitive questions. The results show that the two-stage cluster sampling method and corresponding formulae were feasible.

Keywords: Monte Carlo simulation, Sensitive questions, RRT, Two-stage cluster sampling, Reliability assessment.

1 Introduction

The sample survey is commonly used in medical research and health survey research. In a sample investigation, if a question is sensitive or highly personal, it is likely to lead to either refusals to respond or untruthful answers by using the traditional method of direct interview because of the respondent's concern about revealing their privacy, such as cheating, sexual activity, drugs, AIDS, etc, namely, sensitive question [1]. In order to obtain reliable information about the sample, it is necessary to use a scientific available technology - Randomized Response Technique (RRT). The RRT was originally introduced by Warner in 1965, in Warner's randomized response model, population is divided into two mutually exclusive categories [2]. But in reality population always be divided into more than two categories, namely, multinomial sensitive question, such as investigation of male sexual behaviors[2], it can be divided into "anal sex", "oral sex" and "other sexual behavior".

This paper is a part of the study sampling design on sensitive questions investigation which is supported by National Natural Science Foundation of China [2]. In this paper, we not only focus on statistical method of one-sample RRT model of two-stage cluster sampling on multinomial sensitive question investigation, but also employ Monte Carlo simulation to examine the reliability and validity of that statistical method. Meanwhile, it provides a scientific and reliable method for the sensitive issues investigation in large and complex sampling [3].

^{*} Corresponding author.

2 Methods

2.1 RRT Model for Multinomial Sensitive Question

A RRT model for multinomial sensitive question was adopted in the investigation. Suppose this sensitive question was divided into K kinds of incompatible categories. In the model there is a randomization device which contains some balls written $0, 1, 2, \dots, K$ respectively. The proportion of each kind of ball is $P_0, P_1, P_2, \dots, P_k (P_0 + P_1 + P_2 + \dots + P_k = 1)$. Every respondent was instructed to pick out a ball from the device randomly, If the ball signed 0 the answer was the corresponding serial number of sensitive question they belong to. Otherwise, was the number picked out.

2.2 The Concept of Two-Stage Cluster Sampling

Suppose the population is composed of N_1 primary units, and the i th primary unit contains N_{i2} second-stage units ($i = 1, 2, 3, \dots, N_1$). On average, each primary unit contains \bar{N}_2 second-stage units. The j th second-stage unit of the i th primary unit contains N_{ij3} third-stage units ($i = 1, 2, 3, \dots, N_1, j = 1, 2, 3, \dots, N_{i2}$). On average, each second-stage unit contains \bar{N}_3 tertiary units. The population totally includes N tertiary units.

At the first stage, n_1 primary units were randomly selected from the population. At the second stage, n_{i2} second-stage units were randomly drawn from the i th selected primary unit ($i = 1, 2, 3, \dots, n_1$). On average, \bar{n}_2 secondary units were drawn from per chosen primary unit. The total tertiary units were investigated from the selected second-stage unit.

2.3 Statistical Formulae

Let p_k and $v(p_k)$ stand for the estimator of the population proportion and its variance in the k th categorize respectively. Furthermore, p_{i-k} and p_{ij-k} denote the sample proportion of i th chosen primary unit and the proportion of j th chosen second-stage unit of the i th chosen primary unit. According to the formulae given by WANG Jianfeng, GAO Ge [4], the estimator are shown to be

$$p_k = \frac{N_1}{N n_1} \sum_{i=1}^{n_1} \frac{N_{i2}}{n_{i2}} \sum_{j=1}^{n_{i2}} N_{ij3} P_{ij-k} = \frac{N_1}{N n_1} \sum_{i=1}^{n_1} \frac{N_{i2}}{n_{i2}} p_{i-k} \sum_{j=1}^{n_{i2}} N_{ij3} \quad k = 1, 2, \dots, K \tag{1}$$

$$V(p_k) = \frac{\sigma_1^2}{n_1} \left(1 - \frac{n_1}{N_1} \right) + \frac{\sigma_2^2}{n \bar{n}_2} \left(1 - \frac{\bar{n}_2}{\bar{N}_2} \right) \tag{2}$$

3 Monte Carlo Simulation

3.1 Simulated Population

In the study, sample size was calculated by the formula of the relevant literatures on two-stage cluster sampling [4]. This sampling was implemented to investigate behavioral features of MSMs in Beijing, in August 2010. In the first stage, 6 districts/counties were randomly selected from 16 districts/counties (primary units) of Beijing city ($N_1 = 16$, $n_1 = 6$). In the second stage, 28 chambers of MSMs were randomly sampled from the primary units selected in the first stage as cluster (secondary units), such as bathhouses nightclubs, bars etc. Then RRT model for multinomial sensitive question is presented for application to ask MSMs about the sexual behavior they have with male partners lately [5]. According to formula (1) (2), we could get the estimate of the population proportion of each sexual behavior (p_k) is "0.6095, 0.1685, 0.2190" and its variance ($v(p_k)$) is "0.0009, 0.0004, 0.0018" respectively.

To get both acceptable accuracy and reliability in parameter estimates, we could calculate sample size of simulated sampling from the relevant research [6]. At the first stage, 2 primary units were randomly selected from the simulated population. At the second stage, 6 secondary units were randomly drawn from the selected primary unit [7]. The total simulated respondents, which account for approximately 30% of the entire tertiary units.

Sampling procedure and statistical calculation were simulated by Monte Carlo method. The simulated population was composed of the field investigation [8]. First of all, all parameters were set, for instance, the total number of investigated population, the number of unit at all levels and so on. The simulated population consists in 6 primary units ($N_1 = 6$), Primary units contain 7, 5, 3, 5, 3 and 3 secondary units respectively. The simulated population totally contains 1523 tertiary units ($N = 1523$). Therefore, the proportion of the simulated population of each sexual behavior (p_k) is "0.6095, 0.1685, 0.2190" respectively.

3.2 Part of Program Code

```
Close all; Data=load('data.txt');
R=randint (1,2,[1,6]);
while(sign==0)
    if(R(1)==R(2))
        R=randint (1,2,[1,6]);
    Else
        sign=1;
```

```

end end
for i=1:1523
    if(Data(i,1)==R(1))
        if(Data(i,2)>num_yule_first)
            num_yule_first=Data(i,2);
        end end
for i=1:1523
    if(Data(i,1)==R(1))
        if (Data(i,2)==R_yule_first(1))
            num_human_11=num_human_11+1;
            if(Data(i,4)==1)
                num_ans11_1=num_ans11_1+1;
            end
            if(Data(i,4)==2)
                num_ans11_2=num_ans11_2+1;
            end end
p1=3/1523*(num_yule_first/2*p1_1*(num_human_11+num_huma
n_12)+num_yule_second/2*p2_1*(num_human_21+num_human_22
));
p2=3/1523*(num_yule_first/2*p1_2*(num_human_11+num_huma
n_12)+num_yule_second/2*p2_2*(num_human_21+num_human_22
));
p3=3/1523*(num_yule_first/2*p1_3*(num_human_11+num_huma
n_12)+num_yule_second/2*p2_3*(num_human_21+num_human_22
));
vaiance1=0.5*ss1_1*(1-1/3)+0.25*ss2_1*(1-3/7);
vaiance2=0.5*ss1_2*(1-1/3)+0.25*ss2_2*(1-3/7);
vaiance3=0.5*ss1_3*(1-1/3)+0.25*ss2_3*(1-3/7);

```

3.3 Simulated Investigation Results

Monte Carlo simulation sampling was repeated 30 times under matlab program. Each tertiary unit has produced a response value. Then we could get the estimators of the simulated population proportion and the simulated population proportion variance of that sensitive question according to formulae, all the results were given in Table 1.

Table 1. Monte Carlo simulation results

Simulation Times	\hat{p}_1	$\hat{v}(p_1)$	\hat{p}_2	$\hat{v}(p_2)$	\hat{p}_2	$\hat{v}(p_3)$
1	65.48	0.67	13.10	0.34	21.43	0.50
2	69.27	0.61	10.78	0.28	19.95	0.46
3	65.50	0.83	14.33	0.45	20.18	0.59
4	56.81	0.72	16.90	0.41	26.29	0.57
5	70.85	0.61	16.35	0.41	12.80	0.33
6	69.10	0.57	17.60	0.39	13.30	0.31
7	76.08	0.49	12.50	0.29	11.42	0.27
8	67.58	0.63	17.35	0.41	15.07	0.37
9	55.33	0.78	15.99	0.43	28.68	0.65
10	58.52	0.83	21.43	0.58	20.05	0.55
11	64.01	0.79	11.81	0.36	24.18	0.63
12	62.80	0.59	9.96	0.23	27.24	0.50
13	63.17	0.60	24.07	0.47	12.76	0.29
14	58.54	0.74	19.51	0.48	21.95	0.52
15	57.76	0.66	25.43	0.51	16.81	0.38
16	69.83	0.54	11.98	0.27	18.18	0.38
17	65.64	0.62	18.28	0.41	16.08	0.37
18	63.21	0.69	12.50	0.32	24.29	0.54
19	56.13	0.73	19.58	0.46	24.29	0.54
20	70.83	0.72	11.11	0.34	18.06	0.51
21	68.48	0.59	14.13	0.33	17.39	0.39
22	61.61	0.66	16.96	0.39	21.43	0.47
23	57.73	0.84	13.54	0.40	28.73	0.71
24	66.20	0.66	15.73	0.39	18.08	0.43
25	57.99	0.99	15.09	0.47	26.92	0.73
26	62.94	0.64	10.31	0.25	26.75	0.54
27	64.64	0.79	9.39	0.29	25.97	0.66
28	58.05	0.64	23.09	0.47	18.86	0.41
29	62.55	0.63	23.59	0.49	13.85	0.32
30	71.30	0.76	10.65	0.35	18.05	0.55

3.4 Reliability and Validity Assessment

The principles of reliability and validity are fundamental cornerstones of the scientific method [7, 8, 9]. Reliability is defined as the extent to which a measurement is repeated under identical conditions. Validity refers to the degree to which a test measures what it purports to measure [10]. In this paper, Monte Carlo simulation was used to evaluate the reliability and validity. Establishing good quality studies need both high reliability and high validity [11].

Assessment of reliability: The one-sample Z-test is used to test whether the proportion of every simulation-based investigation on the sensitive question, which can be considered as the simulated sample proportion, is significantly different from the simulated population proportion [12]. Table 2 gives a summary of statistical analyses on data from simulation survey. The p-value statistic for a two-tailed test ranges from 0.07 to 0.84. It is clear that there are no statistically significant

differences between each simulated sample and the simulated population. Therefore all results of simulated investigation are very close to the simulated population proportion, showing that the method and formula for one-sample RRT model for application to multi-classified sensitive question under two-stage cluster sampling are highly reliable [13].

Assessment of validity: The total results of simulated investigation are very close to the simulated population proportion [12]. In addition, the simulated population proportion was obtained by researchers through conducting an actual investigation. In conclusion, our method and corresponding formulae presented in this study are of quite high validity [11-13].

Table 2. Assessment of reliability and validity of sensitive question

Simulation Times	$k=1$		$k=2$		$k=3$	
	z -value	p -value	z -value	p -value	z -value	p -value
1	0.56	0.57	-0.61	0.54	0.06	0.95
2	1.23	0.22	-0.08	0.94	-1.13	0.26
3	1.04	0.30	-1.07	0.28	-0.12	0.91
4	0.52	0.61	-0.36	0.72	-0.08	0.93
5	-0.42	0.68	0.01	0.99	0.62	0.53
6	0.99	0.32	0.11	0.91	-1.09	0.28
7	1.95	0.05	-0.76	0.45	-1.41	0.16
8	0.83	0.41	0.07	0.94	0.79	0.43
9	-0.56	0.57	0.13	0.89	0.85	0.39
10	-0.21	0.83	0.58	0.56	-0.10	0.92
11	0.37	0.71	-0.80	0.42	0.36	0.72
12	0.27	0.79	-1.33	0.18	0.77	0.44
13	0.31	0.75	1.01	0.31	-1.18	0.24
14	-0.22	0.82	0.37	0.71	0.13	0.90
15	-0.32	0.75	1.16	0.25	-0.55	0.58
16	1.17	0.24	-0.88	0.38	-0.36	0.72
17	0.60	0.55	0.21	0.83	-0.65	0.52
18	0.30	0.76	-0.73	0.47	0.40	0.69
19	-0.49	0.62	0.39	0.70	0.40	0.69
20	1.14	0.25	-0.93	0.35	-0.34	0.73
21	0.96	0.34	-0.45	0.65	-0.46	0.64
22	0.12	0.90	0.02	0.98	0.02	0.98
23	-0.29	0.77	-0.50	0.62	0.83	0.41
24	0.65	0.51	-0.17	0.86	-0.36	0.72
25	-0.26	0.80	-0.25	0.80	0.63	0.53
26	0.28	0.78	-1.22	0.22	0.69	0.49
27	0.44	0.66	-1.30	0.19	0.55	0.58
28	-0.29	0.77	0.87	0.38	-0.27	0.79
29	0.23	0.81	0.93	0.35	-1.00	0.32
30	1.16	0.24	-1.00	0.32	-0.33	0.74

4 Discussion

Monte Carlo Simulation method (Monte Carlo Simulation MCS) is a kind of numerical calculation method which is based on statistical sampling theory and studying random variable through computer [14]. As the increasing complexity and dimension of the questions, the algorithm complexity of traditional deterministic numerical methods increases exponentially, it seems more and more difficult. But Monte Carlo method is only related with repeated times; it can describe the characteristics of the random nature things intuitively [12-13].

Recently, with the rapid development of computer simulation technology, Monte Carlo method has developed into an important research tool [15]. The systematic and flexibility of the method offers the possibility for the promotion and application in other areas. In many of the useful applications, the mathematical problem itself arises in a problem of probability in physics, biology, pharmacy, operational research, general statistics, economics, or econometrics [16-17]. In the our research, Monte Carlo sampling simulation, employed under two-staged complex survey designs for sensitive question, had certain innovation and high-applied value.

Most of the literatures on theory of RRT are restricted to simple random sampling, especially for the research on sensitive questions [12-16]. Moreover, objects of investigation on sensitive questions confined to a small range are selected by a simple random design [15]. However, evaluations of reliability and validity on sensitive questions in survey using the RRT have been seldom reported [17]. In our present research, these weaknesses have been successfully overcome [11].

The method and formulae for one-sample RRT model of two-stage cluster sampling on sensitive question investigation show higher reliability and validity [13]. Thus, two-stage cluster sampling appeared in this paper is presented as an effective method for obtaining real data of sensitive questions in a wide range of area [15]. This would be expected to not only allow local policy makers to better formulate public health policy and guide efficient allocation of resources, but also provide the scientific basis for effective prevention and control of HIV/AIDS among high risk group [18]. At the same time, our research results will fill the research blank for the statistical survey method and calculation formulae [15-18].

References

1. Wang, M., Gao, G.: Quantitative sensitive question survey in cluster sampling and its application. In: *Recent Advance in Statistics Application and Related Areas*, pp. 648–652 (2008)
2. Li, W., Gao, G., He, Z.: *Statistical Methods of Two-Stage Sampling on Simmons Model for Sensitive Question Survey with and Its Application*. *Studies in Mathematical Sciences*, 46–51 (2011)
3. He, Z.L., Gao, G., Wang, L.: Multiplication models of quantitative sensitive questions survey in two-stage sampling and its application. In: *Data Processing and Quantitative Economy Modeling*, pp. 6–10 (2010)

4. Wang, J., Ge, G., Fan, Y., et al.: The estimation of sampling size in multi-stage sampling and its application in medical survey. *Applied Mathematics and Computation*, 239–249 (2006)
5. Shaul, K.B., Elizabeta, B., Benzion, B.: A two-stage sequential sampling scheme for Warner's randomized response model. In: *Communications in Statistics Theory and Methods*, pp. 2373–2387 (2003)
6. Wen, L., Ge, G., Lei, W.: Stratified random sampling on the Simmons model for sensitive question survey. In: *Data Processing and Quantitative Economy Modeling*, pp. 21–24 (2010)
7. Liu, P., Gao, G., He, Z., Ruan, Y.H., Li, X.D., Yu, M.R.: Two-stage sampling on additive model for quantitative sensitive question survey and its application. *Progress in Applied Mathematics*, 67–72 (2011)
8. Su, L.: *Advanced Mathematical Statistic*. Beijing University Press (2007)
9. Gao, G., Fan, Y.: Stratified cluster sampling and its application on the Simmons RRT model for sensitive question survey. *Chinese Journal of Health Statistics*, 562–565 (2008)
10. He, Z., Gao, G.: Multiple choices sensitive questions survey in two-stage sampling and its application. In: *Recent Advance in Statistics Application*, pp. 1160–1164 (2009)
11. Kim, J.M., Warde, W.D.: A stratified Warner's randomized response model. *Journal of Statistical Planning and Inference*, 155–168 (2004)
12. Zannette, A.U., Chantay, M.D.: Sensitive topics: Are there model differences? *Computers in Human Behavior*, 76–87 (2009)
13. van den Hout, A., van der Heijden, P.G.M., Gilchrist, R.: The logistic regression model with response variables subject to randomized response. *Computational Statistics & Data Analysis*, 6060–6069 (2007)
14. Li, X.D., Gao, G., Yu, M.R.: Stratified random sampling on the randomized response technique for multi-class sensitive question survey. In: *Recent Advance in Statistics Application and Related Areas*, pp. 800–803 (2009)
15. Lau, J.T.F., Lin, C., Hao, C., et al.: Public health challenges of the emerging HIV epidemic among men who have sex with men in China. *Public Health*, 260–265 (2011)
16. Zhang, D., Bi, P., Lv, F., et al.: Changes in HIV prevalence and sexual behavior among men who have sex with men in a northern Chinese city: 2002–2006. *Journal of Infection*, 456–463 (2007)
17. Fido, A., Al Kazemi, R.: Survey of HIV/AIDS knowledge and attitudes of Kuwaiti family physicians. *Family Practice*, 682–684 (December 2002)
18. Gage, A.J., Ali, D.: Factors associated with self-reported HIV testing among men in Uganda. *AIDS Care*, 153–165 (2004)

New Lax-Friedrichs Scheme for Convective-Diffusion Equation*

Haixin Jiang** and Wei Tong

College of Science, Jiujiang University
Jiujiang, China, 332005
jianghaixin@163.com

Abstract. Two different types of generalized Lax-Friedrichs scheme for the convective-diffusion equation $u_t + au_x = \varepsilon u_{xx}$ ($a \in R, \varepsilon > 0$) are given and analyzed. For the convection term, both of two schemes use generalized Lax-Friedrichs scheme. For the diffusion term, explicit central difference scheme is used. The propagation of chequerboard mode is considered for two schemes and several numerical examples are presented only for the first scheme, which display how different parameters control oscillations. For low and high frequency modes, applying discrete Fourier analysis, some results have been obtained about stability condition of schemes, and clarify the reasons of oscillations and the interrelation among the amplitude error.

Keywords: generalized Lax-Friedrichs (LxF) scheme, discrete Fourier analysis, oscillations, frequency modes.

1 Introduction

In [1], to compute the numerical solution of the hyperbolic conservation laws

$$u_t + f(u)_x = 0, x \in R, t > 0, \quad (1)$$

where $u = (u_1, \dots, u_m)^T$, and $f(u) = (f_1, \dots, f_m)^T$, we consider the generalized Lax-Friedrichs(LxF) scheme of the viscosity form

$$u_j^{n+1} = u_j^n - \frac{\nu}{2}[f(u_{j+1}^n) - f(u_{j-1}^n)] + \frac{q}{2}(u_{j+1}^n - 2u_j^n + u_{j-1}^n), \quad (2)$$

where the mesh ratio $\nu = \tau/h$ is assumed to be a constant, τ and h are step sizes in time and space, respectively, u_j^n denotes an approximation of $u(jh, n\tau)$, the term $q \in (0, 1]$ is the coefficient of numerical viscosity. When $q = 1$, it is the classical Lax-Friedrichs(LxF) scheme.

With the flux function $f = au$, (1) is the linear advection equation as follow

$$u_t + au_x = 0, \quad x \in R, t > 0, \quad (3)$$

and the scheme (2) turns into the generalized LxF scheme of equation (3)

* Supported by Natural Science Foundation of Jiujiang University.

** Corresponding author.

$$u_j^{n+1} = u_j^n - \frac{\nu a}{2}(u_{j+1}^n - u_{j-1}^n) + \frac{q}{2}(u_{j+1}^n - 2u_j^n + u_{j-1}^n). \tag{4}$$

By adding a diffusion term εu_{xx} (ε is a positive constant) to the right of (3), we obtain a convective-diffusion equation.

$$u_t + au_x = \varepsilon u_{xx}, \quad a \in R, \varepsilon > 0. \tag{5}$$

There are two different finite difference schemes of the convective-diffusion equation (5). For the convective term, we still use the generalized LxF scheme. Then, we have two following ways to approximate the diffusion term: one uses explicit central difference scheme, i.e.

$$u_j^{n+1} = u_j^n - \frac{\nu a}{2}(u_{j+1}^n - u_{j-1}^n) + \frac{q}{2}(u_{j+1}^n - 2u_j^n + u_{j-1}^n) + \frac{\varepsilon \tau}{h^2}(u_{j+1}^n - 2u_j^n + u_{j-1}^n), \tag{6}$$

and the other one uses implicit central difference scheme, i.e.

$$u_j^{n+1} = u_j^n - \frac{\nu a}{2}(u_{j+1}^n - u_{j-1}^n) + \frac{q}{2}(u_{j+1}^n - 2u_j^n + u_{j-1}^n) + \frac{\varepsilon \tau}{h^2}(u_{j+1}^{n+1} - 2u_j^{n+1} + u_{j-1}^{n+1}). \tag{7}$$

Scheme (6) also can be written in this form

$$u_j^{n+1} = u_j^n - \frac{\nu a}{2}(u_{j+1}^n - u_{j-1}^n) + \left(\frac{q}{2} + \mu\right)(u_{j+1}^n - 2u_j^n + u_{j-1}^n), \tag{8}$$

where $\mu = \frac{\varepsilon \tau}{h^2}$. Scheme (8) is the generalized LxF schemes of the convective-diffusion equation (5). As observed in [1], we discussed the discretization of initial data

$$u(x, 0) = u_0(x), \quad x \in [0, 1], \tag{9}$$

with M grid points and $h = 1/M$, while M is even, and $u_0(0) = u_0(1)$. The numerical solution value at the grid point x_j is denoted by u_j^0 . We express this grid point value u_j^0 by using the usual discrete Fourier sums, as in [5, Page 120], and obtain

$$u_j^0 = \sum_{k=-M/2+1}^{M/2} c_k^0 e^{i\xi k j}, \quad i^2 = -1, \quad j = 0, 1, \dots, M-1, \tag{10}$$

where $\xi = 2\pi kh$. And the coefficients c_k^0 are expressed as

$$c_k^0 = \frac{1}{M} \sum_{j=0}^{M-1} u_j^0 e^{-i\xi k j}, \quad k = -M/2 + 1, \dots, M/2, \tag{11}$$

First, for the special case that

$$c_k^0 = \begin{cases} 1, & \text{if } k = M/2, \\ 0, & \text{otherwise.} \end{cases} \tag{12}$$

i.e. the initial datas are taken to be just the highest Fourier mode which is a single chequerboard mode oscillation

$$u_j^0 = e^{i2\pi\frac{M}{2}jh} = e^{i\pi j} = (-1)^j. \tag{13}$$

Second, For single square signal, as follow

$$u(x, 0) = \begin{cases} 1, & 0 < x^{(1)} < x < x^{(2)} < 1, \\ 0, & otherwise. \end{cases}$$

we have two different types of discretizations, if they are discretized with an odd number of grid points, the chequerboard mode is present. In contrast, if they are discretized with an even number of grid points, the chequerboard mode is suppressed.

According to the discussions in [1], the solution to (??) also can be expressed in the standard form of discrete Fourier series,

$$u_j^n = \sum_{k=-M/2+1}^{M/2} c_k^n e^{i\xi j}, \tag{14}$$

where $\xi = 2\pi kh$. Combining (14) and (??), the coefficients c_k^n are obtained analogously and expressed as,

$$c_k^n = [1 + q'(\cos \xi - 1) - i\nu a \sin \xi]^n c_k^0. \tag{15}$$

First, if we take initial data as a single chequerboard mode like (13), i.e. $u_j^0 = e^{i2\pi\frac{M}{2}jh} = e^{i\pi j} = (-1)^j$, we have

$$u_j^n = (1 - 2q')^n (-1)^j. \tag{16}$$

2 Discrete Fourier analysis

In this section, we use the method of discrete Fourier analysis to discuss the generalized LxF scheme (8) for the convective-diffusion equation. While using $q' = q + 2\mu$, the variable q' in scheme (??) acts as q in scheme (4). Thus, most important conclusions are obtained easily through the related results about scheme (4) of the linear advection equation in [1].

We denote a Fourier mode using the scaled wave number $\xi = 2\pi kh$ by $e^{i\xi}$. Then using it as initial data for a linear finite difference scheme results after n times steps in the solution

$$u_k^n = \lambda_k^n e^{i\xi} = (\lambda(k))^n e^{i\xi}, \quad i^2 = -1, \tag{17}$$

where λ_k^n is the amplitude. The modulus of the ratio

$$\lambda(k) = \lambda_k^{n+1} / \lambda_k^n$$

is the amplitude of the mode for one time step. For the scheme (8) we have with $\nu = \tau/h$ in particular, where $\mu = \frac{\varepsilon\tau}{h^2}$.

$$\lambda(k) = 1 + (q + 2\mu)(\cos \xi - 1) - i\nu a \sin \xi \tag{18}$$

and

$$|\lambda(k)|^2 = [1 + (q + 2\mu)(\cos \xi - 1)]^2 + (\nu a)^2 \sin^2 \xi. \tag{19}$$

In [1], considering $\xi \approx 0$, $\xi \approx \pi$ and $|\lambda| \leq 1$, we obtained the conditions $0 < \nu^2 a^2 \leq q \leq 1$ are necessary and sufficient for stability of scheme (1.4). Therefore, the conditions $0 < \nu^2 a^2 \leq q' \leq 1$, i.e.

$$0 < \nu^2 a^2 \leq q + 2\mu \leq 1, \tag{20}$$

are necessary and sufficient for stability of scheme (8), too.

The exact solution of the Fourier mode $e^{i\xi}$ for $x = h$ after one time step τ is $e^{i(\xi - 2\pi a k \tau)} = e^{-i2\pi a k \tau} e^{i\xi} = \lambda_{exact}(k) e^{i\xi}$. The exact amplitude $\lambda_{exact}(k)$ has modulus 1. If the modulus of $\lambda(k)$ is less than one, the effect of the multiplication of a solution component with $\lambda(k)$ is called *numerical dissipation* and then the amplification error is called *dissipation error*.

Further, comparing the exponents of $\lambda(k)$ and $\lambda_{exact}(k)$ there is a *phase error*

$$\arg \lambda(k) - (-\nu a \xi),$$

where $\nu = \tau/h$, $\xi = 2\pi k h$. The *relative phase error* is then defined as

$$E_p(k) := \frac{\arg \lambda(k)}{-\nu a \xi} - 1.$$

A mode is a low frequency mode if $\xi \approx 0$ and a high frequency mode if $\xi \approx \pi$.

2.1 Low Frequency Modes

We first look at the low frequency modes

$$(U^s)_j^n := \lambda_k^n e^{i\xi j}, \quad \xi \approx 0$$

Then substituting $(U^s)_j^n$ into (??),

$$\lambda(k) = 1 + q'(\cos \xi - 1) - i\nu a \sin \xi \tag{21}$$

and

$$\begin{aligned} |\lambda(k)|^2 &= [1 + q'(\cos \xi - 1)]^2 + (\nu a)^2 \sin^2 \xi \\ &= 1 - (1 - \cos \xi)[2q' - q'^2(1 - \cos \xi) - \nu^2 a^2(1 + \cos \xi)], \end{aligned} \tag{22}$$

where $q' = q + 2\mu$. We express (22) by using $\cos \xi = 1 - \frac{1}{2}\xi^2 + O(\xi^4) + \dots$

$$|\lambda|^2 = 1 - [q' - \nu^2 a^2]\xi^2 + O(\xi^4) \tag{23}$$

So, we can see that the dissipation error is of order $O(\xi)$, if $q' > \nu^2 a^2$.

As $q' = q + 2\mu \leq 1$, $\xi \in [0, \pi/2]$, from (2.11) we obtain

$$\frac{d(|\lambda(k)|^2)}{dq'} = 2[1 - q'(1 - \cos \xi)] \cdot (\cos \xi - 1) < 0. \tag{24}$$

This implies that when $0 < \nu^2 a^2 \leq q' \leq 1$, the dissipation becomes stronger as q' larger. If μ is fixed, we also can say the dissipation becomes stronger as q larger. Thus scheme (8) with $q + 2\mu = 1$ has the largest numerical dissipation for low frequency modes.

The phase of the low frequency modes in approximated by Taylor expansion at $\xi = 0$

$$\begin{aligned} \arg \lambda &= \arctan\left(\frac{-\nu a \sin \xi}{1 - q'(1 - \cos \xi)}\right) = \arctan\left(\frac{-\nu a(\xi - \frac{1}{6}\xi^3 + \frac{1}{120}\xi^5 + \dots)}{1 - q'(1 - \frac{1}{2}\xi^2 + \dots)}\right) \\ &= \arctan\left(-\nu a \xi + \frac{1 - 3q'}{6} \nu a \xi^3 + \dots\right). \end{aligned} \tag{25}$$

According to lemma 4.1 [2, page 97], we obtain

$$\arg \lambda = -\nu a \xi \left(1 + \frac{3q' - 1 - 2\nu^2 a^2}{6} \xi^2 + \dots\right). \tag{26}$$

Then,

$$E_p(k) = \frac{3q' - 1 - 2a^2\nu^2}{6} \xi^2 + \dots,$$

so the relative phase error $E_p(k)$ is of order $O(\xi^2)$, at least (in some cases, $3q' - 1 - 2a^2\nu^2 = 0$). Therefore, oscillations caused by this relative phase error can be suppressed by the stronger dissipation of order $O(\xi)$.

2.2 High Frequency Modes

For high frequency modes, i.e. $\xi \approx \pi$, the situation is very different. We use $\xi = \pi + \xi'$, i.e. $\xi' = 2\pi k'h$ with $kh = 1/2 + k'h$, and thus $\xi' \approx 0$. We write the modes in the form

$$(U^h)_j^n = \lambda_k^n e^{i\xi j} = \lambda_k^n e^{i(\pi+\xi')j} = (-1)^{j+n} \lambda_{k'}^n e^{i\xi' j}, \tag{27}$$

with $\lambda_{k'}^n = (-1)^{j+n} e^{i\pi j} \lambda_k^n$ and set

$$(U^o)_j^n := \lambda_{k'}^n e^{i\xi' j}.$$

The factor $(U^o)_j^n$ can be regarded as a perturbation amplitude of the checker-board modes ($e^{i\pi} \lambda_k^n = (-1)^{j+n}$). The dissipation (amplitude error) depends only on $\lambda_{k'}^n$. Then substituting $(U^h)_j^n$ into (??) yields

$$\lambda' := \lambda_{k'}^{n+1} / \lambda_{k'}^n = -1 + q'(1 + \cos \xi') - i\nu a \sin \xi'.$$

Therefore, we have

$$\begin{aligned}
 |\lambda'|^2 &= [-1 + q'(1 + \cos \xi')]^2 + (\nu a)^2 \sin^2 \xi' \\
 &= 1 + 4(\nu^2 a^2 - q') \cos^2(\xi'/2) + 4(q'^2 - \nu^2 a^2) \cos^4(\xi'/2). \tag{28}
 \end{aligned}$$

Similar to low frequency modes, we express (28) by using $\cos \xi = 1 - \frac{1}{2}\xi^2 + O(\xi^4) + \dots$

$$\begin{aligned}
 |\lambda'|^2 &= [-1 + q'(1 + \cos \xi')]^2 + (\nu a)^2 \sin^2 \xi' \\
 &= (1 - 2q')^2 + (q' - 2q'^2 + \nu^2 a^2)\xi'^2 + O(\xi'^4) \\
 &= 1 - [1 - (1 - 2q')^2] + (q' - 2q'^2 + \nu^2 a^2)\xi'^2 + O(\xi'^4). \tag{29}
 \end{aligned}$$

If $0 < \nu^2 a^2 \leq q' < 1$, the dissipation error is $O(1)$ (particularly if $q' = 1$, the dissipation error is of order $O(\xi')$).

Obviously, for $\xi' \approx 0$, i.e. the high frequency modes, we consider

$$\frac{d(|\lambda'|^2)}{dq'} = -2(1 + \cos \xi')[1 - q'(1 + \cos \xi')]. \tag{30}$$

If $q' > 1/2$, we obtain $\frac{d(|\lambda'|^2)}{dq'} > 0$, which means $|\lambda'|^2$ is an increasing function of q' ; on the other hand, if $q' < 1/2$, we have $\frac{d(|\lambda'|^2)}{dq'} < 0$ and $|\lambda'|^2$ is a decreasing function of q' .

That means, under $0 < \nu^2 a^2 \leq q' < 1$, if q' is closer to $1/2$, high frequency modes decay stronger, see Figure 2 and Figure 3. By the way, the results are in sharp contrast with the situation for low frequency modes.

Furthermore, let us look at the relative phase error. We compute

$$\begin{aligned}
 \arg \lambda' &= \arctan\left(\frac{-\nu a \sin \xi'}{-1 + q'(1 + \cos \xi')}\right) = \arctan\left(\frac{-1}{2q' - 1}\nu a \xi' - \frac{q' + 1}{6(2q' - 1)^2}\nu a \xi'^3 - \dots\right) \\
 &= -\frac{\nu a \xi'}{2q' - 1} - \frac{2q'^2 + q' - 2\nu^2 a^2 - 1}{6(2q' - 1)^3}\nu a \xi'^3 + O(\xi'^5). \tag{31}
 \end{aligned}$$

Then for the high frequency modes $(U^h)_j^n$, we have by recalling that $\xi = \pi + \xi'$

$$\begin{aligned}
 (U^h)_j^n &= (-1)^{j+n} \lambda_k^n e^{i\xi' j} = |\lambda'|^n e^{in(-\pi + \arg \lambda')} \cdot e^{ij(\pi + \xi')} \\
 &= |\lambda'|^n e^{i(j\xi - 2\pi k n \tau)} \cdot e^{in(-\pi + \arg \lambda' + \nu a \xi)}. \tag{32}
 \end{aligned}$$

Therefore, the relative phase error of high frequency modes at each time step is using (31)

$$\begin{aligned}
 E_p(k) &= -\frac{-\pi + \arg \lambda' + \nu a \xi}{\nu a \xi} \\
 &= \frac{\pi(1 - \nu a)}{\nu a \xi} + \frac{(2 - 2q')\xi'}{(2q' - 1)\xi} + \frac{(2q'^2 + q' - 2\nu^2 a^2 - 1)\xi'^3}{6(2q' - 1)^3 \xi} + O(\xi'^5) \tag{33}
 \end{aligned}$$

We note that $\xi \approx \pi$. Therefore the relative phase error has $O(1)$, while $1 - \nu a > 0$. This error is huge, and strong numerical dissipation is needed to suppress it.

When $0 < \nu^2 a^2 \leq q' < 1$, the dissipation error is of order $O(1)$. As the parameter q' is closer to $1/2$, the dissipation error becomes larger. The dissipation error can suppress the numerical oscillations caused by the relative phase error when $q' = 1/2$.

3 The Relation of Numerical Oscillations and Several Parameters

In this part, we only discuss high frequency modes. (1)The relation between numerical oscillations and the coefficient of numerical viscosity q : From (30), we have

$$\frac{d(|\lambda'|^2)}{dq} = \frac{d(|\lambda'|^2)}{dq'} \cdot \frac{dq'}{dq} = -2(1 + \cos \xi')[1 - (q + 2\mu)(1 + \cos \xi')]. \quad (34)$$

If $q' = q + 2\mu > 1/2$, $\xi' \approx 0$, thus

$$\frac{d(|\lambda'|^2)}{dq} > 0. \quad (35)$$

It implies that the dissipation becomes stronger for high frequency modes as the parameter q decreases, and then oscillations become weaker. That is, the numerical dissipation is more effective in controlling numerical oscillations as q decreases. In contrast, if $q' = q + 2\mu < 1/2$, $\xi' \approx 0$, we obtain $\frac{d(|\lambda'|^2)}{dq} < 0$. It means the dissipation becomes stronger for high frequency modes as the parameter q increases. (2)The relation between numerical oscillations and the coefficient of physical viscosity ε : Similar to (1) above, since $q' = q + 2\mu = q + \frac{2\varepsilon\tau}{h^2}$, we obtain

$$\frac{d(|\lambda'|^2)}{d\varepsilon} = \frac{d(|\lambda'|^2)}{dq'} \cdot \frac{dq'}{d\varepsilon} = -\frac{4\tau}{h^2}(1 + \cos \xi')[1 - (q + \frac{2\varepsilon\tau}{h^2})(1 + \cos \xi')].$$

If $q' = q + 2\mu > 1/2$, i.e. $q + \frac{2\varepsilon\tau}{h^2} > 1/2$, $\xi' \approx 0$, we obtain

$$\frac{d(|\lambda'|^2)}{d\varepsilon} > 0.$$

The same to the parameter q , the dissipation becomes much more stronger for high frequency modes as ε decreases, and then oscillations become weaker. As ε decreases, the numerical dissipation is effective in controlling numerical oscillations. If $q' = q + 2\mu < 1/2$, the conclusion is opposite. (3) The relation between numerical oscillations and width of mesh h . Now the mesh ratio $\nu = \tau/h$ is assumed to be a constant c and h is step size in space. Similar to (1) above, since $q' = q + 2\mu = q + \frac{2\varepsilon c}{h}$, we have

$$\begin{aligned} \frac{d(|\lambda'|^2)}{dh} &= \frac{d(|\lambda'|^2)}{dq'} \cdot \frac{dq'}{dh} = -2(1 + \cos \xi')[1 - q'(1 + \cos \xi')] \cdot \frac{-2\varepsilon c}{h^2} \\ &= \frac{4\varepsilon c}{h^2}(1 + \cos \xi')[1 - (q + \frac{2\varepsilon c}{h})(1 + \cos \xi')]. \end{aligned}$$

If $q' = q + 2\mu > 1/2$, i.e. $q + \frac{2\epsilon c}{h} > 1/2$, $\xi' \approx 0$, we obtain

$$\frac{d(|\lambda'|^2)}{dh} < 0.$$

This implies that the dissipation becomes much more stronger for high frequency modes as the parameter h increases, and then oscillations become weaker. That is, the numerical dissipation is effective in controlling numerical oscillations as h increases. However, if $q' = q + 2\mu < 1/2$, the conclusion is opposite, too.

References

1. Li, J., Yang, Z., Zheng, Y.: Characteristic decompositions and interactions of rarefaction waves of 2-D Euler Equations. *Journal of Differential Equations* 250, 782–798 (2011)
2. Li, J., Li, Q., Xu, K.: Comparison of the Generalized Riemann Solver and the Gas-Kinetic scheme for Compressible Inviscid Flow Simulations. *Journal of Computational Physics* 230, 5080–5099 (2011)
3. Dou, L., Dou, J.: The Grid: Time-domain analysis of lossy multiconductor transmission lines based on the LaxCWendroff technique. *Analog Integrated Circuits and Signal Processing* 68, 85–92 (2011)
4. Zhu, P., Zhou, S.: Relaxation Lax-Friedrichs sweeping scheme for static Hamilton-Jacobi equations. *Numer. Algor.* 54, 325–342 (2010)
5. Li, J., Liu, T., Sun, Z.: Implementation of the GRP scheme for computing spherical compressible fluid flows. *Journal of Computational Physics* 228, 5867–5887 (2009)
6. Li, J.-Q., Tang, H.-Z., Warnecke, G., Zhang, L.-M.: Local Oscillations in Finite Difference Solutions of Hyperbolic Conservation Laws. *Mathematics of Computation*, S0025-5718(09)02219-4
7. Morton, K.W., Mayers, D.F.: *Numerical Solution of Partial Differential Equations*, 2nd edn. Cambridge University Press (2005)
8. Breuss, M.: The correct use of the Lax-Friedrichs method. *M2AN Math. Model. Numer. Anal.* 38, 519–540 (2004)
9. Breuss, M.: An analysis of the influence of data extrema on some first and second order central approximations of hyperbolic conservation laws. *M2AN Math. Model. Numer. Anal.* 39, 965–993 (2005)
10. Thomas, J.W.: *Numerical Partial Differential Equations: Finite Difference Methods*. Springer (1995)
11. Warming, R.F., Hyett, B.J.: The modified equation approach to the stability and accuracy of finite difference methods. *J. Comput. Phys.* 14, 159–179 (1974)
12. Tadmor, E.: Numerical viscosity and the entropy condition for conservative difference schemes. *Math. Comp.* 43, 369–381 (1984)

Span of T-Colorings Multigraphs

Juan Du

Department of Public Course, Environment Management College of China, 06004
Qinhuangdao, China
dujuan1978@eyou.com

Abstract. Suppose G is a graph and T is a set of nonnegative integers. A T -coloring of G is an assignment of a positive integer $f(x)$ to each vertex x of G so that if x and y are joined by an edge of G , then $|f(x) - f(y)|$ is not in T . Here, the vertices of G are transmitters, an edge represents interference, $f(x)$ is a television or radio channel assigned to x , and T is a set of disallowed separations for channels assigned to interfering transmitters. The span of a T -coloring of G equals $\max |f(x) - f(y)|$, where the maximum is taken over all edges $\{x, y\} \in E(G)$. The minimum order, and minimum span, where the minimum is taken over all T -colorings of G , are denoted by $\chi_T(G)$, and $sp_T(G)$, respectively. We will show several previous results of multigraphs, and we also will present a new algorithm to compute $sp_T(G)$ of multigraphs.

Keywords: T -coloring, multigraphs, span, interference, algorithm.

1 Introduction

T -colorings of graphs arose in connection with frequency assignment problem in communications. In this problem, there are n transmitters x_1, x_2, \dots, x_n situated in a region. To each transmitter x_i , a channel $f(x_i)$ (a fixed positive integer) is to be assigned. Some of the transmitters interfere because of proximity, meteorological, or other reasons. Two interfering transmitters must be given frequencies such that the absolute value of the difference of their frequencies does not belong to a forbidden set T of nonnegative integers. Our objective is to make a frequency assignment that is efficient according to certain criteria, while satisfying the above constraint.

One level of interference frequency assignment problem is defined as follows: $V = \{x_1, x_2, \dots, x_n\}$, and $\{x_i, x_j\}$ is in E if and only if x_i and x_j interference. It is assumed that 0 belongs to T . The requirement can be summarized by the following equation:

$$\{x, y\} \in E(G) \Rightarrow |f(x) - f(y)| \notin T \quad (1)$$

We are in the case of ordinary vertex coloring when $T = \{0\}$.

For any T-set, it is easy to find a T-coloring. If α is the largest entry of T, then one can use a different channel for each transmitter, choosing channels from $\{1, \alpha + 2, 2\alpha + 3\}$. However this will not be very efficient in terms of the difference between the smallest and largest channels used. Our objective is to find efficient assignments for various graphs and T-sets. So let us at this point define some criteria for efficient channel assignment.

Cozzens and Roberts [1] introduced the following definition and notation. Given a graph G and T-set T, the order of a T-coloring f of G is the number of distinct values of $f(x), x \in V(G)$. The span of a T-coloring f of G equals $\max |f(x) - f(y)|$, where the maximum is taken over all edges $\{x, y\} \in E(G)$. The minimum order, and minimum span, where the minimum is taken over all T-colorings of G, are denoted by $\chi_T(G)$, and $sp_T(G)$, respectively.

Although interference can be due to geographic proximity, meteorological factor, etc., Hale in [2] pointed out that we define interference constrains as a function of frequency alone. Constrains of type are important when physical distance between transmitters are either unknown or uncontrollable, and only frequency may be adjusted.

In this landmark paper, Hale [2] describes a unifying theory and terminology for all frequency assignment problems. Hale defines the kth general constraint as

$R(k) = [T(k), d(k)]$, where $T(k) \subset Z_+, d(k) \in Z_+, k = 0, 1, \dots, K-1$, and $0 = T(0) \subset T(1) \subset \dots \subset T(K-1)$ and $d(0) > d(1) > \dots > d(K-1)$. Then if $|f(x_i) - f(x_j)| \notin T(k), \forall x_i, x_j \in V \rightarrow f$ feasible for $V, T(k)$.

Thus, the kth adjacent constraint is a special instance of the general kth constraint where $T(k) = \{0, 1, \dots, k\}$ and $d(k) \in Z_+$. Notice that in general, $T(k)$ need not be a contiguous set of integers; in fact it may be any set of integers that satisfies the nesting property defined above.

In making frequency assignments, we sometimes consider several different levels of interference. For instance, transmitters at most 10 miles apart might interfere at one level, while transmitters at most 50 miles apart might interfere at a second level. To take into account these different levels of interference, we consider K different graphs, $G_0, G_1, G_2, \dots, G_{k-1}$, each on the same vertex set V, the set of transmitters, with an edge between transmitters x_u and x_v appearing in graph G_i if and only if x_u and x_v interfere at level i . In UHF television, K is 5. (See Hale [2], MiddleKamp [3], and Push et al. [4]). For each level i , we have a disallowed set of separations $T(i)$ for transmitters interfering at level i . Typically

$$G_0 \supseteq G_1 \supseteq G_2 \cdots \supseteq G_{k-1} \tag{2}$$

and

$$T(0) \subseteq T(1) \subseteq T(2) \cdots \subseteq T(K-1) \tag{3}$$

The results on T-colorings of one level interference will now be extended to the general (and more practical) case where interference may occur on different level. Here, the transmitters are the vertices of a multigraph and each interference level is represented by a separate set of edges on the common set of vertices.

A multigraph is a graph without the restriction that the edges be distinct. That is, a pair of vertices may have more than one edge connecting them. If the edge set E of a multigraph is partitioned into K distinct sets, then we have a multigraph that is a family of K graphs $G_0, G_1, G_2, \dots, G_{K-1}$ on the same vertex set which we represent by $G(V, G_0, G_1, G_2, \dots, G_{K-1})$.

We seek a function which assigns to each transmitters a frequency, a positive integer, so that f is simultaneously a $T(i)$ -coloring of G_i for all i , i.e., so that for $i = 0, 1, \dots, K - 1$,

$$\{x, y\} \in E(G_i) \Rightarrow |f(x) - f(y)| \notin T(i) \tag{4}$$

For instance, if $K = 2$, and $T(0) = \{0\}$, $T(1) = \{0, 1\}$, then if x and y interfere at level 0, they must get different frequencies but if they interfere at level 1, they must get not only different but also non-adjacent frequencies. If $G = G(V, G_0, G_1, G_2, \dots, G_{K-1})$ satisfies (2), we call it a nested graph, and if it satisfies (4), we call it a T-coloring. The order and span of a T-coloring are defined as before, as are $\chi_T(G)$ and $sp_T(G)$.

2 Previous Results

In this section, we quote some known results about T-colorings of multigraphs. See [5], [6] for results.

Theorem 1 (Cozzens and Wang [7]). Let $G = G(V, G_0, G_1, G_2, \dots, G_{K-1})$ be a nested graph and let $T = (T(0), T(1), \dots, T(K - 1))$ satisfy (3), then

- (i): $\chi_T(G) = \chi(G_0)$
- (ii): $sp_T(G) \geq esp_T(G) \geq \chi(G_0) - 1$
- (iii): $\max sp_{T(m)}(G_m) \leq sp_T(G) \leq sp_{T(K-1)}(G_0)$
- (iv): If each $T(i)$ is r_i -initial, $i = 0, 1, \dots, K - 1$; $0 \leq m \leq K - 1$, then $\max [(r_m + 1)(\chi(G_m) - 1)] \leq sp_T(G) \leq (r_{K-1} + 1)(\chi(G_0) - 1)$
- (v): If $\chi(G_i) = \chi(G_j)$, for all i, j , and if $T(i)$ is r_i -initial for all i , then $sp_T(G) = sp_{T(K-1)}(G_0) = (r_{K-1} + 1)(\chi(G_0) - 1)$.
- (vi): If $\chi(G_i) = \chi(G_j)$, for all i, j , and if $T(i)$ is r_i -initial for all i , and if G_0

is chordal, then the greedy algorithm finds T-colorings of order χ_T and span sp_T in $o(n^2)$ time, where $n = |V|$.

Raychaudhuri [8] developed the following intriguing result for 2-level multigraphs with the restriction that graph G_0 be complete.

Theorem 2 (Raychaudhuri [8]). Consider the 2-level nested multigraph $G(V, G_0, G_1)$, where G_0 is complete and $T(0) = \{0\}, T(1) = \{0, 1\}$. If the weighted graph G' is defined as: $V(G') = V(G), E(G') = E(G_0)$, and every edges $\{x, y\} \in G'$ has a weighted $W(x, y)$ assigned as

$$W(x, y) = \begin{cases} 2 & \text{if } \{x, y\} \in E(G_1) \\ 1 & \text{else} \end{cases}$$

Then $sp_T(G) = \text{length of the shortest Hamiltonian path in } G'$, where the length of a Hamiltonian path is the sum of the weights of the edges contained in the path.

Raychaudhuri [8-9] also observe that it follows from Theorem 2 that, using results of Goodman and Hedetniemi [10] and of Boesch et al. [11], one can compute sp_T in $o(n^2)$ time in the situation of Theorem 2 if in addition G_1 is a tree or a forest.

3 A New Algorithm to Compute sp_T of Multigraphs

To compute sp_T of multigraphs, we can use Theorem 2, but it runs only for 2-level multigraphs, so we will give a more widely useful algorithm to compute sp_T of multigraph $G(V, G_0, G_1, G_2, \dots, G_{k-1})$, where $T(k)$ is an any set of nonnegative integers. Suppose H is any Hamiltonian path of a multigraph $G(V, G_0, G_1, G_2, \dots, G_{k-1})$ and W_H is the weight of path H.

Color H with: $f(v_1) = c_1 = 0, f(v_i) = c_i (i = 2, 3, \dots, n)$, there

$$c_{i+1} = \min \left\{ z \mid z - c_j > 0, \text{ and } z - c_j \notin T(k_{i+1,j}) \right\} \tag{5}$$

$$(i = 1, 2, \dots, n-1, j = 1, 2, \dots, i, z \in Z_+)$$

and $k_{j,i+1}$ is the order of the edges v_j, v_i . Let $W_H = c_n, W(G) = \min W_H$.

Theorem 3. Suppose $v_i \rightarrow f_i, i = 1, 2, \dots, n$, and $T(k)$ is an any set of nonnegative integers. $f_1 = 0, f_i = c_i (c_i \in Z_+)$, and $c_{i+1} = \min \left\{ z \mid z - c_j > 0, \text{ and } z - c_j \notin T(k_{i+1,j}) \right\} (i = 1, 2, \dots, n-1, j = 1, 2, \dots, i, z \in Z_+)$, H is any Hamiltonian path of G, W_H is the weight of this Hamiltonian path. Let $W_H = f_n$, and $W(G) = \min W_H(G)$. So

- (i) f is a T-coloring of multigraph G .
- (ii) $sp_T(G) = W(G)$.

Proof. $T(0) \subseteq T(1) \subseteq T(2) \dots T(K-1)$

(i) we will prove $\{c_1, c_2, \dots, c_n\}$ is a T -coloring of G .

$$\because c_i - c_j > 0$$

$$\therefore |c_i - c_j| \notin T(k, j)$$

And, according this way, we color v_i with c_i , and color v_j with c_j , for $\forall 1 \leq i, j \leq n$, we can get

$$\{v_i, v_j\} \in E(G_k) \Rightarrow |c_i - c_j| \notin T(k) \quad (0 \leq k \leq K-1)$$

So, f is a T -coloring of G .

(ii) first we will prove $sp_T(G) \leq W(G)$.

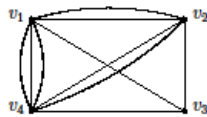
Suppose $f = \{f_1^*, f_2^*, \dots, f_n^*\}$ is a T-coloring of G , and satisfying (5) and $f_n^* = W(G)$. So $sp_T(G) \leq sp_T(f) = f_n^* - f_1^* = f_n^*$, that is $sp_T(G) \leq \min W_H(G)$.

then we will prove $sp_T(G) \geq W(G)$

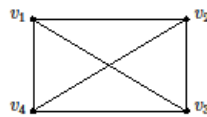
Suppose f' is any T-coloring of G that satisfying (5),

let $f' = \{c'_1, c'_2, \dots, c'_n\}$, and without loss generality, we let $0 = c'_1 < c'_2 < \dots < c'_n$ and the vertex which is colored with c'_i is v_i , then since $L' = \{v_1, v_2, \dots, v_n\}$ is a Hamiltonian path of G , we can get $sp_T(f') = c'_n \geq W_L(G) \geq \min W_H(G)$.

example

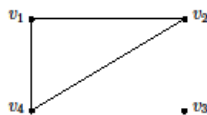


graph G



graph G_0

$$T(0) = \{0\}$$



graph G_1

$$T(1) = \{0, 1\}$$



graph G_2

$$T(2) = \{0, 1, 2\}$$

Consider the following cases:

- (1). Hamiltonian path $H_1(v_1, v_2, v_3, v_4)$: give v_1 a color 0, since the order of edges $\{v_1, v_2\}$ is 2, color v_2 with 2, so color v_3 with 3, v_4 with 4. $W_{H_1} = 4$.

- (2). Hamiltonian path $H_2(v_1, v_2, v_4, v_3)$: give v_1 a color 0, so color v_2 with 2, v_4 with 4, v_3 with 5. $W_{H_2} = 5$.
- (3). Hamiltonian path $H_3(v_1, v_3, v_2, v_4)$: $W_{H_3} = 4$.
- (4). Hamiltonian path $H_4(v_1, v_3, v_4, v_2)$: $W_{H_4} = 4$.
- (5). Hamiltonian path $H_5(v_1, v_4, v_2, v_3)$: $W_{H_5} = 6$.
- (6). Hamiltonian path $H_6(v_1, v_4, v_3, v_2)$: $W_{H_6} = 5$.
- (7). Hamiltonian path $H_7(v_2, v_1, v_3, v_4)$: $W_{H_7} = 4$.
- (8). Hamiltonian path $H_8(v_2, v_1, v_4, v_3)$: $W_{H_8} = 6$.
- (9). Hamiltonian path $H_9(v_2, v_3, v_1, v_4)$: $W_{H_9} = 5$.
- (10). Hamiltonian path $H_{10}(v_2, v_4, v_1, v_3)$: $W_{H_{10}} = 6$.
- (11). Hamiltonian path $H_{11}(v_3, v_1, v_2, v_4)$: $W_{H_{11}} = 5$.
- (12). Hamiltonian path $H_{12}(v_3, v_2, v_1, v_4)$: $W_{H_{12}} = 6$.

So, $sp_T(G) = \min W_H = 4$.

References

1. Cozzens, M.B., Roberts, F.S.: T-colorings of graphs and the channel assignment problem. Congr. Numer. 35, 191–208 (1982)
2. Hale, W.K.: Frequency assignment: theory and applications. Proc. of the IEEE 68(12), 1497–1514 (1980)
3. Middlekamp, L.C.: UHF taboos-history and development. IEEE Trans. Consumer Electron. CE-24, 514–519 (1978)
4. Pugh, G.E., Lucas, G.L., Krupp, J.C.: Optimal allocation of TV channels-a feasibility study, Tech. Rep. DSA No.261., Decision-Science Applications, Inc., Arlington, VA (August 1981)
5. Bonias, I.: T-Colorings of complete graphs, Ph.D. Thesis, Department of Mathematics, Northeastern University, Boston, MA (1991)
6. Tesman, B.: T-colorings, list T-coloring, and set T-colorings of graphs, Ph.D. Thesis, Department of Mathematics, Rutgers University, New Brunswick, NJ (October 1989)
7. Cozzens, M.B., Wang, D.-I.: The general channel assignment problem. Congr. Number. 41, 115–129 (1984)
8. Raychaudhuri, A.: Intersection assignments, T-colorings, and powers of graphs, Ph.D. Thesis, Department of Mathematics, Rutgers University, New Brunswick, NJ (1985)
9. Raychaudhuri, A.: Further results on T-colorings and frequency assignment problems. SIAM. J. Discrete Math. (to appear)
10. Goodman, S., Hedetniemi, S.: On the Hamiltonian completion problem. In: Bar, R., Harary, F. (eds.) Graphs and Combinatorics. Lecture Notes in Math., vol. 406, pp. 262–274. Springer, Berlin (1974)

11. Lucarelli, G., Milis, I., Paschos, V.T.: On the Maximum Edge Coloring Problem. In: Bampis, E., Skutella, M. (eds.) WAOA 2008. LNCS, vol. 5426, pp. 279–292. Springer, Heidelberg (2009)
12. Boesch, F.T., Chen, S., McHugh, J.A.M.: On covering the points of a graph with point disjoint paths. In: Bari, R., Harary, F. (eds.) Graphs and Combinatorics. Lecture Notes in Math., vol. 406, pp. 201–212. Springer, Berlin (1974)
13. Tesman, M.A.: Set T-colorings of forbidden difference graphs to T-colorings. *Congressus Numerantium* 74, 15–24 (1980)
14. Murphey, R.A., Panos, Resende, M.G.C.: Frequency assignment problems. In: Handbook of Combinatorial Optimization. Kluwer Academic Publishers (1999)
15. Roberts, F.S.: T-colorings of graphs: recent results and problems. *Discrete Mathematics* 93, 229–245 (1991)
16. Janczewski, R.: Greedy T-colorings of graphs. *Discrete Mathematics* 309(6), 1685–1690 (2009)
17. Aicha, M., Malika, B., Habiba, D.: Two hybrid ant algorithms for the general T-colouring problem. *International Journal of Bio-Inspired Computation* 2(5), 353–362 (2010)
18. Villegas, E.G., Ferré, R.V., Paradells, J.: Frequency assignments in IEEE 802.11 WLANs with efficient spectrum sharing. *Wireless Communications and Mobile Computing* 9(8), 1125–1140 (2009)
19. Malaguti, E.: The Vertex Coloring Problem and its generalizations. *4OR: A Quarterly Journal of Operations Research* 7(1), 101–104 (2009), doi:10.1007/s10288-008-0071-y
20. Leila, N., Malika, B.: Resolution of the general T-coloring problem using an MBO based algorithm. In: 2011 10th IEEE International Conference on Cognitive Informatics & Cognitive Computing (ICCI*CC), pp. 438–443 (2011)

Dual-Scaled Method for the Rheology of Non-newtonian Boundary Layer and Its High Performance FEM

Lei Hou¹, Hanling Li¹, Ming Zhang¹, Weijia Wang¹, Dezhi Lin¹, and Lin Qiu²

¹Dept. Math, Shanghai University, Shanghai, China

²Dept. Math, Shanghai Jiao Tong University, Shanghai, China
houlei@shu.edu.cn, qiulin916@yahoo.com.cn

Abstract. In this paper, a macro-micro dual-scaled method is used to model the rheology of non-Newtonian fluids. In the micro scale, stochastic analysis of boundary layer data resulting from engineering test is presented, and bending deformation of cellular porous materials is introduced. In the macro scale, coupled PDEs: Cauchy fluid equation and P-T/T stress equation are used for modeling free surface and over-stretched element, which are non-Newtonian fluid domain. Semi-discrete finite element method is used to solve the macroscopic equations, and three solving schemes are compared. The call of high performance function library, the NAG, is introduced.

Keywords: dual-scaled, non-Newtonian boundary layer, cellular porous materials, semi-discrete FEM, high performance function library (NAG).

1 Introduction

Rheology of non-Newtonian fluid plays a great role in the study of material property, for its science and technology concept has a direct influence upon its commercial value. For example, the cellular porous material, which is widely used in the automobile industry, shows both solid property and fluid property on boundary layer during its high velocity impact process. Another example is the mixing of rubber compound [1].

In this paper, coupled PDEs: Cauchy fluid equation and P-T/T [2] stress equation, is used to model the rheology of complex contact boundary layer of non-Newtonian fluid. The P-T/T equation is a Maxwell equation with an exponential impact which is the principle character of non-Newtonian fluid. Rheology of cellular porous materials is introduced in section 2. Rheology of rubber can be found in [3,4].

The semi-discrete finite element method is used to solve the coupled PDEs. To improve the solving precision, stability and speed, several solving schemes are compared. Lagrangian interpolation function with 9-point biquadrate element in the space domain is adopted in [5], while Hermite's interpolation function with 4-point bicubic element is discussed in [6]. Comparison of Euler time difference scheme with Runge-Kutta time difference scheme can be found in [5]. In this paper, the numerical results using Euler scheme, modified Euler scheme and Crank-Nicolson scheme are presented. High-performance computing platforms and computing software, such as

the NAG function library and Ls-Dyna solver, are applied. High-performance finite element solution of the macroscopic equations is discussed in section 3.

Since 1753 when Euler first proposed the continuum model, the method of using both macro and micro scale to study the mathematical laws behind the physical phenomena is widely applied: the fluid, which is composed of numerous molecular in the microscopic point of view, can be regarded as a composition of a large number of fluid particles, according to the macroscopic point of view; the fluid particle contained numerous molecular, therefore make it possible to describe its property using statistical average of molecular. The macro scale in this paper is the coupled PDEs, while in the micro scale, the statistical method is used to extract the required initial and boundary conditions for solving the PDEs. The micro-scale simulation is given in section 2. [7] used molecular dynamics (micro) and mass-energy conservation equation (macro) to study complex fluids, which is an important reference for this article.

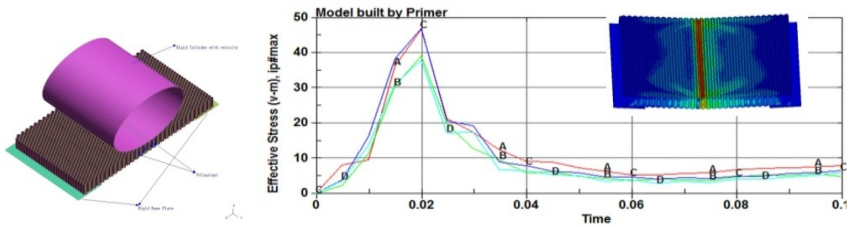


Fig. 1. Simulation of aluminum honeycomb

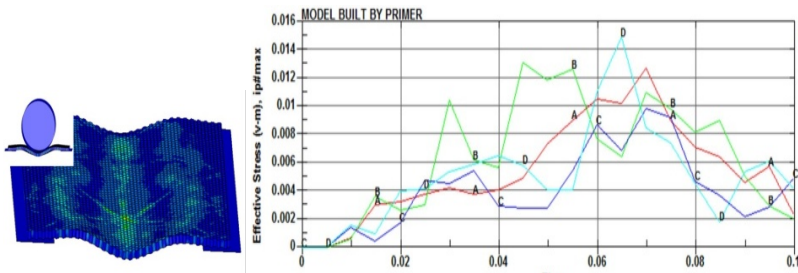


Fig. 2. Simulation of rubber honeycomb

2 Microscopic Structure

2.1 Cellular Porous Structure of Boundary Layer

Crash safety is a transient problem: in the short-term of the collision process, materials withstand large deformation. At present, substitute polymer materials (such as cellular materials) for metal materials in some part of the automobile is the mainstream method to slow down the impact[8], and to control plastic deformation. The study of cellular porous materials can be attributed to the numerical calculation of complex non-smooth contact surface.

2.2 Finite Element Simulation

The finite element simulation of honeycomb porous materials shows that material properties have a significant effect on distribution of deformation and stress: the stress of aluminum materials is relatively severe, and grid deformation is mainly concentrated in the compressed position, where folding phenomenon is obvious (Fig. 1); rubber material can get relatively stable deformation and uniform stress distribution (Fig. 2).

3 Macroscopic Equations

3.1 Coupled Solid-Liquid Equations

In this paper, the domain between the contact surface and non-contact surface is regarded as fluid domain[9,10], which is formed of the free surface element and overstretched elements. It is usually singular. The extension and simple shear rate resistance are analyzed by studying the non-Newtonian PT/T equation (1) [6]. The large deformation caused by macroscopic and microscopic distribution of stress field τ is calculated by studying the Cauchy conservation eq. (2) [5,11].

$$\rho \underline{\dot{u}} = [\nabla \cdot \underline{\tau} - \rho \underline{u} \cdot \nabla \underline{u}] \tag{1}$$

$$\lambda \underline{\dot{\underline{\tau}}} = [2\eta \underline{D} - \exp\left\{\frac{\epsilon \lambda}{\eta_0}(\tau_{xx} + \tau_{yy})\right\} \tau] - \lambda [\underline{u} \cdot \nabla \underline{\tau} - \nabla \underline{u} \cdot \underline{\tau} - (\nabla \underline{u} \cdot \underline{\tau})^T + \xi \left[\underline{D} \cdot \underline{\tau} + (\underline{D} \cdot \underline{\tau})^T \right]] \tag{2}$$

The semi-discrete finite element method is used to solve the above PDEs. In order to improve the precision and stability, time difference schemes and finite element basis function should be carefully selected. Discussion of basis function on space domain can be found in [5,6]. The discrete form of coupled PDEs can be found in [12].

3.2 Time Difference Scheme

In this paper, the numerical results of Euler difference scheme, modified Euler difference scheme and Crank-Nicolson (C-N) difference scheme are presented. The precision and stability of the three schemes are shown in table 1. $\lambda^B_j (j = 1, 2, \dots, N)$ are all the eigen values of stiff matrix B[11].

Table 1. Precision and stability of the three schemes

	Euler	modified Euler	Crank-Nicolson
precision	one order	two order	two order
stability	poor	depends on $\Delta t * \max\{ \lambda^B_j \}$	absolute stable

The amplification factors of modified Euler scheme (μ_{j_E}) and C-N scheme ($\mu_{j_{CN}}$) are shown in eq. (3). Obviously, $|\mu_{j_{CN}}| \leq 1$, therefore assured the absolute convergence of C-N scheme. The modified Euler scheme is stable only when $0 \leq \Delta t \lambda_j^B \leq 2$.

$$\mu_{j_E} = \frac{1 - \frac{\Delta t \lambda_j^B}{2}}{1 + \frac{\Delta t \lambda_j^B}{2}} \mu_{j_{CN}} = \frac{1}{2} \left[\left(1 - \Delta t (\lambda_j^B)^2 \right) + 1 \right] \tag{3}$$

C-N scheme and modified Euler scheme

Table 1 shows that Δt should be selected as small as possible when using modified Euler scheme. However, a too small Δt will greatly increase computing time, and will cause error in calculation. Such problem also exists in C-N scheme. In this paper, $\Delta t = 0.4$ is selected to compare the two schemes. Table 2 shows the numerical results of $\Delta t * \max_j \{|\lambda_j^B|\}$. u_0 is the initial velocity of eq. (3), and T is the time span. The stability of modified Euler scheme is poorer than C-N scheme.

Table 2. Numerical results of $\Delta t * \max_j \{|\lambda_j^B|\}$

	$u_0 = 0.2$	$u_0 = 0.4$	$u_0 = 0.6$	$u_0 = 0.8$	$u_0 = 0.9$	$u_0 = 1$
T=0.4	0.000564	0.000564	0.000764	0.000800	0.000632	0.000823
T=8	0.57281	0.60669	0.64220	0.67935	0.75854	0.71813
T=12	49.6	137.8	3513.5	135.9	1224.4	74.3

C-N scheme and Euler scheme.

In general, C-N scheme works better than modified Euler scheme does. However, the FEM for PDEs usually evolves large scale of calculation, which makes it important to consider the computation time. Fig. 3 to Fig. 5 are numerical results of eq. (3) using different schemes. One grid is used as a numerical experiment to illustrate the numerical results of coupled PDEs. Details about coupling process can be found in [11].

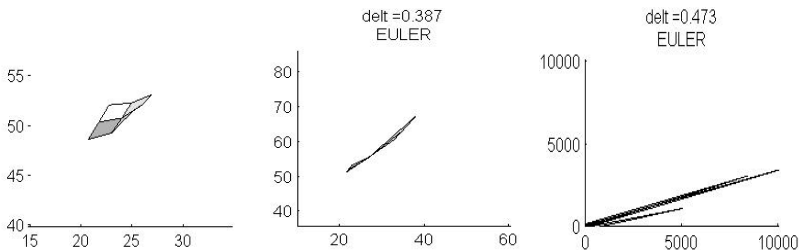


Fig. 3. Positions of grid at T=10, using Euler scheme; $u_0 = 5$, and Δt is shown in each figure

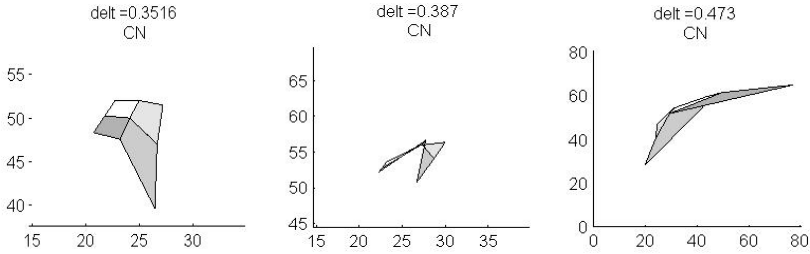


Fig. 4. Positions of grid at $T=10$, using C-N scheme; $u_0 = 5$, and Δt is shown in each figure

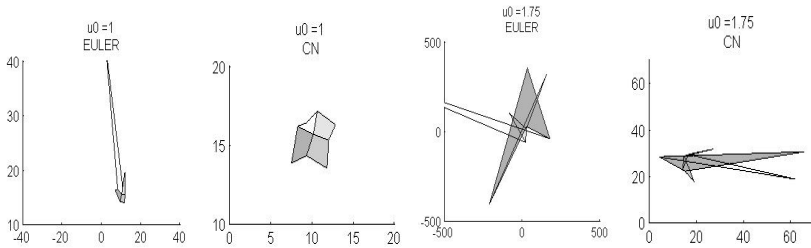


Fig. 5. Positions of grid at $T=13$; $\Delta t = 1$; time scheme and u_0 is shown in each figure

Fig. 3 and Fig. 4 show the influence of Δt . C-N scheme maintains a good mesh shape, while Euler scheme may sometimes results in grid rupture. Fig. 5 shows the numerical results under various initial velocities: C-N scheme exerts wider application than Euler scheme does. Note that the range of coordinates is inconsistent.

3.3 High Performance Function Library

The function library of the NAG, i.e., Numerical Algorithms Group, is called during the solving of PDEs using semi-discrete finite element method. It contains functions for generating Gauss weight and abscises (d01bb), and for calculating the integral (d01fb). The latter take too much time to pass the parameters when the whole matrix is to be integrated, therefore it is necessary to write a custom summing function.

Besides, the NAG provides grid generating function on 2-D domain, which is under the class of D06, and functions that apply specific properties (such as symmetric, positive definiteness, sparsity, etc.) of coefficient matrix, which is under the class of F.

4 Stochastic Analysis

Apart from solving mathematic model using the FEM, to analyze the data resulting from numerical results (especially the data of boundary layer) is also a great way to further study the non-Newtonian materials. The stochastic analysis focuses on the explicit influence factors to the system: the impact angle (IA) [12], the plastic strain (PS), and the Von-Mises stress (VMS) [13]. As frequency may characterize the

component (e.g., IA) contributing to the deformation result (e.g., PS, VMS) in the system, histogram graph of the frequency is presented in this paper. Fig. 6 shows the trend of IA. The trend of PS and VMS can be found in [12]. Fig. 7 is a 3-d histogram of IA and PS. A trend towards normal distribution can be observed.

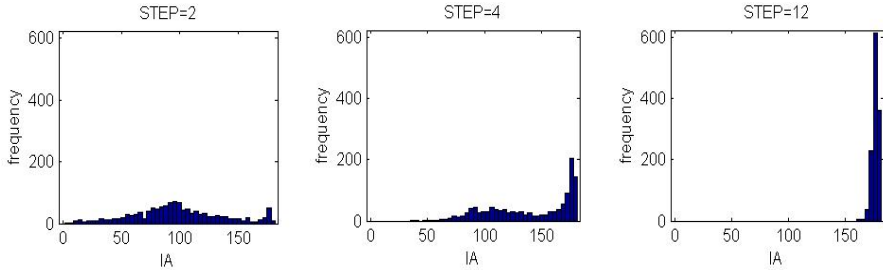


Fig. 6. The frequency of the impact angle

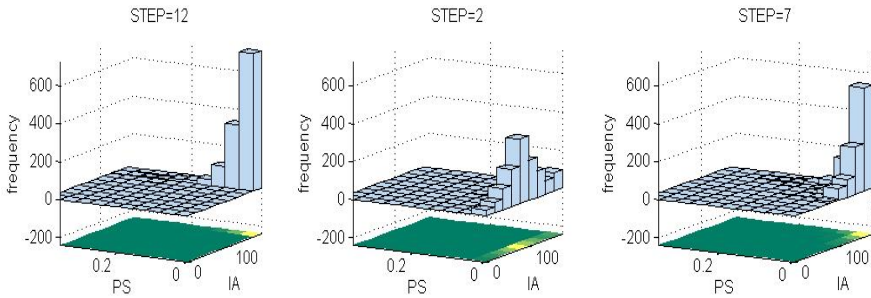


Fig. 7. 3-d histograms of plastic strain

5 Conclusions

In this paper, a macro-micro dual-scaled method is used to model the rheology of non-Newtonian fluids. In the micro scale, the FEM simulation of boundary layer of non-Newtonian material is presented. Compared with aluminum honey comb, rubber honey comb can get relatively more stable deformation and uniform stress field. Normality is observed in the boundary layer.

In the macro scale, coupled PDEs (i.e., Cauchy fluid equation and P-T/T stress equation) are used to describe the deformation and stress of boundary layer of non-Newtonian fluid. Semi-discrete finite element method is used to solve the coupled PDEs. The space domain is discrete by Lagrangian interpolation function with 9-point biquadrate element. The time domain is discrete by three schemes: Euler scheme modified Euler scheme and Crank-Nicolson scheme. Numerical results show that the stability of Crank-Nicolson scheme is better than Euler scheme and modified Euler scheme, and the precision of Crank-Nicolson is better than Euler scheme, which coexists with theoretical analysis.

References

1. Hou, L., Harwood, R.: Nonlinear Properties in the Newtonian and Non-Newtonian Equations. *Non-linear Analysis*. Elsevier Sciences 30(4), 2497–2505 (1997)
2. Thien, N.P., Tanner, R.I.: A New Constitutive Equation Derived From Network Theory. *J. Non-Newton Fluid* 2(4), 353–365 (1977)
3. Hou, L., Nassehi, V.: Evaluation of Stress Effects Flow in Rubber Mixing. *Nonlinear Analysis*. Elsevier Sciences 47(3), 1809–1820 (2001)
4. Hou, L., B PR, D. WA: Physics of Plasmas. *American Institute of Physics* 2, 473–481 (1996)
5. Hou, L., Cai, L.: Nonlinear Property of the Visco-Elastic-Plastic Material in the Impact Problem. *Journal of Shanghai University (English Edition)* 13(1), 23–28 (2009)
6. Hou, L., Ding, H., Li, H., Qiu, L.: The High Performance Computing on the Crash-Safety Analysis. In: Zhang, W., Chen, Z., Douglas, C.C., Tong, W. (eds.) *HPCA 2009*. LNCS, vol. 5938, pp. 169–176. Springer, Heidelberg (2010)
7. Ren, W., Weinan, E.: Heterogeneous Multiscale Method for the Modeling of Complex Fluids and Micro-Fluidics 204(1), 1–26 (2005)
8. Hou, L., Wang, H., Cai, L.: Engineering Mathematical Study for the Visco-Elastic Impact Problem. In: *Proceedings of the World Congress on Engineering*, pp. 2701–2706 (2009)
9. Hou, L., Li, H.L., Lin, D.Z.: The Stochastic Boundary-Layer in Thenon-Newtonian Problem. In: *Proceedings of the World Congress on Engineering*, pp. 1872–1876 (2010)
10. Hou, L., Qiu, L.: Computation and Asymptotic Analysis in the Non-Newtonian Impact Problem. *ACTA Math. Applic. Sinica (English)* 25(1), 117–127 (2009)
11. Hou, L., Lin, D.Z., Li, H.L.: Computational Modelling on the Complex Boundary Conditions in the Impact Problem. In: *International Conference on Computer and Network Technology*, pp. 231–235 (2011)
12. Hou, L., Li, H.L., Wang, H.: Stochastic Analysis in the Visco-Elastic Impact Condition. *International Review of Chemical Engineering* 2(1), 178–183 (2010)
13. Frank, P.: The Work of Richard Von Mises: 1883-1953. *Science* 119, 823–824 (1954)

Factor Model Averaging Quantile Regression and Simulation Study

Zhimeng Sun*

School of Statistics,
Central University of Finance and Economics,
South College Road.39, 100081, Beijing, P.R. China
zmsun82@163.com

Abstract. In this paper, a model averaging approach is developed for the linear factor regression model in light of smoothed focused information criterion. With respect to factors, a frequentist model averaging estimation of the regression parameter is proposed based on quantile regression techniques, and the model averaging estimator thus is non-sensitive to outliers and robust. We show that the asymptotic properties of the proposed estimator is asymptotically normal and root-n consistent. A simulation study is conducted to investigate the finite properties of the proposed estimator.

Keywords: Model averaging, Factor, Quantile regression.

1 Introduction

Consider the following linear regression model

$$Y = X^\top \beta + \varepsilon, \quad (1)$$

where Y is a scalar response, $X = (X_1, \dots, X_p)^\top$ is a p -dimensional vector of explanatory variables, $\beta = (\beta_1, \dots, \beta_p)^\top$ is the vector of unknown parameters, ε is the random error of the model, \top represents transpose.

When building a linear model for real data, model selection or variable selection at the initial stage of modeling is of great importance to the whole process of data analysis. It has been always one of the hot topics of statistics since the seventies of last century and a great many well-known criteria of model selection have been developed such as AIC (Akaike (1973)), BIC (Schwarz (1978)), Lasso (Tibshirani 1996), SCAD (Fan and Li 2001), FIC (Claeskens and Hjort (2003)), among others. To a large extent, these model selection methods provide some effective solution to the problem of choosing which variables to be included in the model in order to select out a relative better model under some

* This work is supported by the National Natural Science Foundation of China (No.11171011); Foundation of Academic Discipline Program at Central University of Finance and Economics; Fund of Third Stage Foundation of 211 Project at Central University of Finance and Economics.

condition. However, an increasing amount of research shows that these kinds of model selection approaches also have some inherent disadvantages, for example, it may ignore some uncertainty in the stage of model selection and thus underestimate the variance of the resulted estimator (Leeb and Pötscher (2003, 2008)); It also can not avoid the risk of selecting a very poor model which may lead to perishing prediction results (Leung and Barron (2006)). One available way to overcome this difficulty is weighting estimators across several models rather than entirely relying upon a single selected model. This method is known as model averaging which can provide insurance against selecting a very poor model and avoid model selection instability (Yang (2001), Leung and Barron (2006)). As a consequence, the model averaging method can improve coverage probabilities. In recent years, model averaging has been widely considered in a growing many of literatures, such as Buckland, Burnham and Augustin (1997), Hansen (2007), Hjort and Claeskens (2003, 2006), Liang, Zou, Wan and Zhang (2011), Zhang and Liang (2011), among others. Most model averaging methods are designed for crucial explanatory variables individually. However, in many fields of applied science, it is often desirable to identify significant explanatory variables in a grouped manner. The most common example is the multifactor analysis-of-variance (ANOVA) problem, in which each factor may have several levels and can be expressed through a group of dummy variables. Another example is when one wants to produce more flexible functions and thus employs an additive nonparametric model with linear approximation method like polynomial or spline approximation other than linear models. Some regularization methods have been proposed for automatic factor selection. This kind of methods include group lasso (Yuan and Lin (2006)), adaptive group lasso (Wang and Leng (2008)), and penalized method using composite absolute penalties (Zhao, Rocha and Yu (2009)).

On the other hand, most model averaging methods, which are based on either least squares or likelihood function, may break down when there are wicked outliers since outliers sometimes can overly determine or even damage the fit of the model. Thus, exploring robust model averaging method which is non-sensitive to outliers is a meaningful question and as far as we know no literature has concerned about this problem. We try to consider this problem through a way of quantile techniques. The quantile regression method gradually emerging as a unified statistical methodology for estimating models of conditional quantile functions. By complementing the exclusive focus of classical least-squares regression on the conditional mean, quantile regression offers a systematic strategy for examining how covariates influence the location, scale, and shape of the entire response distribution. Discussions on quantile regression are referred to Koenker (2005), Kai, Li and Zou (2010), Huang (2010), Pang, Lu and Wang (2012), among others. To the best of our knowledge, there is not yet any research on factor model averaging method through quantile regression techniques, and it is a challenging task. In this paper, we propose a new estimation procedure for factor model factor model averaging based on quantile regression. It is also mentioned that the proposed procedure works finely when the random

error has heavy tails, even if when the random error has a infinite variance, while the least squares method disastrously breaks down.

The rest of this paper is organized as follows. In section 2, we set the model framework and provide quantile based factor model averaging estimation procedure. In section 3, we show the asymptotic distributions of the estimators and construct a confidence interval. In Section 4, we conduct a simulation to detect the finite-sample properties of the proposed estimators. Technical proofs of the theorems are given section 5.

2 Model Framework and Estimation Procedure

Without loss of generality, assume that the explanatory variables $X = (X_1, \dots, X_p)^\top$ of model (1) can be grouped into K factors as (Z_1, \dots, Z_K) , where $Z_j = (X_{j1}, \dots, X_{jp_j})$ is a group of p_j explanatory variables for $j = 1, \dots, K$ and $\sum_{j=1}^K p_j = P$. Then the regression model (1) can be represented by

$$Y = \sum_{j=1}^K Z_j \beta_{(j)} + \varepsilon,$$

where $\beta_{(j)} = (\beta_{j1}, \dots, \beta_{jp_j})^\top$ is the unknown regression coefficient vector associated with the j th factor. To conduct quantile regression, assume that the random error of the model has zero τ th quantile with a fixed constant $\tau \in (0, 1)$. When $\tau = \frac{1}{2}$, the problem then becomes the well known median regression.

Suppose that the regression coefficient β contains two parts $\beta = (\check{\beta}^\top, \bar{\beta}^\top)^\top$, the first part $\check{\beta}^\top = (\beta_{(1)}^\top, \dots, \beta_{(\check{K})}^\top)^\top$ contains coefficients associated with factors we surely wish to be in the model. For example, it may contains coefficients $\beta_{(j)}$ which are rejected to be zero under the null hypothesis $H_0 : \beta_{j1} = \beta_{j2} = \dots = \beta_{jp_j} = 0$ via a F test with some significant level. The second part $\bar{\beta}^\top = (\beta_{(\check{K}+1)}^\top, \dots, \beta_{(K)}^\top)^\top$ contains coefficient associated with factors which we may be potentially included in the model. Similarly, it may contains coefficients $\beta_{(j)}$ which are accepted to be zero under H_0 via a F test. It can be seen that the sizes of $\check{\beta}^\top$ and $\bar{\beta}^\top$ are supposed to be \check{K} and $\bar{K} = K - \check{K}$ respectively. We consider a local misspecification framework where the true value of the parameter vector β is $\beta_0 = (\check{\beta}_0^\top, \delta^\top / \sqrt{n})^\top$, $\delta = (\delta_{(1)}^\top, \dots, \delta_{(\bar{K})}^\top)^\top$ and $\delta_{(j)} = (\delta_{j1}, \dots, \delta_{jp_j})^\top$. The $O(1/\sqrt{n})$ framework is chosen here to assess large-sample approximations to distributions of the factor model average estimators. It is canonical in the sense that it leads to the most fruitful large-sample approximations, with squared model biases and estimator variances as exchangeable currencies, both of size $O(1/n)$.

To derive a model averaging estimator of β_0 , we smooth estimators of β_0 achieved at overall $2^{\check{K}}$ sub-models. Thus, we need to derive an estimator of β_0 for every sub-model firstly. Keeping this in mind, denote $\beta_S = (\check{\beta}^\top, \bar{\beta}_S^\top)^\top$ the

parameter vector in the S th sub-model, where $\ddot{\beta}^\top = (\beta_{(1)}^\top, \dots, \beta_{(\bar{K})}^\top)$, and $\bar{\beta}_S$ contains S factors of $\bar{\beta}$ those are in the S th sub-model. Denote $dim(\ddot{\beta})$ and $dim(\bar{\beta}_S)$ the numbers of parameters contained in $\ddot{\beta}$ and $\bar{\beta}_S$ respectively. Let π_S be the projection matrix of size $dim(\bar{\beta}_S) \times (p_{\bar{K}+1} + \dots + p_K)$, mapping $\bar{\beta}$ to $\bar{\beta}_S$. Thus, the explanatory variables in the S th sub-model are $\Pi_S X$, where $\Pi_S = \text{diag}(I_{dim(\ddot{\beta})}, \pi_S)$, $I_{dim(\ddot{\beta})}$ is a identity matrix of size $dim(\ddot{\beta}) \times dim(\ddot{\beta})$.

Let $\{Y_i, X_i\}, i = 1, \dots, n$ be independent identically distributed observations from the focused model. Then the quantile regression estimator of β_S is defined as follows

$$\hat{\beta}_S(\tau) = \arg \min_{\beta_S} \sum_{i=1}^n \rho_\tau(Y_i - (\Pi_S X)_i^\top \beta_S),$$

where $\rho_\tau(u) = u(\tau - I_{\{u < 0\}})$, $(\Pi_S X)_i$ is the i th observation of $\Pi_S X$.

Then, a class of the factor model average estimators of β_0 is defined as follows:

$$\hat{\beta} = \sum_S w(S) \Pi_S^\top \hat{\beta}_S$$

where $w(S)$ some weights function.

3 Asymptotic Results

The following theorems summarize the asymptotic behavior of the estimators. Before stating these theorems, assume that the following conditions hold.

A1: $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i X_i^\top = \Sigma$, where Σ is strictly positive definite matrixes.

A2: The model errors ε has zero τ th quantile and a continuous, positive density $f(\cdot)$ in a neighborhood of zero.

The following Theorem 1 gives asymptotic normality of $\hat{\beta}_S(\tau)$ for each sub-model.

Theorem 1. Under conditions A1 and A2, as $n \rightarrow \infty$, we have

$$\sqrt{n} \{ \hat{\beta}_S(\tau) - (\ddot{\beta}_0^\top, 0^\top)^\top \} \xrightarrow{L} -(\Pi_S \Sigma \Pi_S^\top)^{-1} \Pi_S G + (\Pi_S \Sigma \Pi_S^\top)^{-1} \Pi_S \Sigma (0^\top, \delta^\top)^\top$$

where $G \sim N(0, \tau(1 - \tau)\Sigma/f^2(0))$, " \xrightarrow{L} " denotes convergence in distribution.

Based on Theorem 1, we can derive the FIC value of the S th sub-model with respect to the j th element of β_0 say $\beta_{0jt}, j = 1, \dots, K, t = 1, \dots, p_j$. Let $R_S = \Pi_S^\top (\Pi_S \Sigma \Pi_S^\top)^{-1} \Pi_S$, $B = \tau(1 - \tau)\Sigma/f^2(0)$ be the variance of G , I_{jt} be a p -dimensional vector with the j th element to be 1 and other elements to be 0, $\hat{\beta}_{full}(\tau)$ be the estimator of β under the full model, $\hat{\delta}$ be the estimator of δ under the full model. Then $\hat{\delta} = \sqrt{n}[0, I] \hat{\beta}_{full}(\tau)$. It follows from Theorem 1 that

$$\hat{\delta} \xrightarrow{d} \Delta \equiv -[0, I] \Sigma^{-1} G + \delta$$

and

$$\Delta \sim N(\delta, [0, I] \Sigma^{-1} B \Sigma^{-1} [0, I]^\top)$$

Then the FIC value of the S th sub-model with respect to β_{0jt} is defined as

$$FIC_S = I_{jt}^\top \left\{ R_S B R_S + (R_S \Sigma - I)(0^\top, \hat{\delta}^\top)^\top (0^\top, \hat{\delta}^\top)(R_S \Sigma - I)^\top \right\} I_{jt} - I_{jt}^\top (R_S \Sigma - I) \begin{pmatrix} 0 & 0 \\ 0 & I_{\bar{d}} \end{pmatrix} \Sigma^{-1} B \Sigma^{-1} \begin{pmatrix} 0 & 0 \\ 0 & I_{\bar{d}} \end{pmatrix} (R_S \Sigma - I)^\top I_{jt}$$

Based on the FIC value of the S th sub-model, we can define the smoothed FIC factor model averaging estimator of β_{0jt} as

$$\hat{\beta}_{0jt} = \sum_S w(S|\hat{\delta}) \hat{\beta}_{0jtS}$$

with weight function $w(S|\hat{\delta}) = \exp(-\frac{FIC_S}{I_{jt}^\top \hat{\Sigma}^{-1} B \hat{\Sigma}^{-1} I_{jt}}) / \sum_S \exp(-\frac{FIC_S}{I_{jt}^\top \hat{\Sigma}^{-1} B \hat{\Sigma}^{-1} I_{jt}})$,

where $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X_i X_i^\top$.

The following theorem involves asymptotic results of the general model average estimator $\hat{\mu}$.

Theorem 2. Under conditions C1 and C2, if the weight function has at most a countable number of discontinuities, we have

$$\sqrt{n}(\hat{\beta}_{0jt} - \beta_{0jt}) \xrightarrow{L} -I_{jt}^\top \Sigma^{-1} G + I_{jt}^\top \left[Q(\Delta)(0^\top, \Delta^\top)^\top - (0^\top, \Delta^\top)^\top \right], n \rightarrow \infty$$

where $Q(\cdot) = \sum_S w(S|\cdot) R_S \Sigma$.

Based on Theorem 2, it is easily seen that if $\hat{\Omega}^2$ is a consistent estimator of $I_{jt}^\top \Sigma^{-1} B \Sigma^{-1} I_{jt}$, then

$$\left[\sqrt{n}(\hat{\beta}_{0jt} - \beta_{0jt}) - I_{jt}^\top \{ Q(\Delta)(0^\top, \Delta^\top)^\top - (0^\top, \Delta^\top)^\top \} \right] / \hat{\Omega} \xrightarrow{d} N(0, 1)$$

As a consequence, a confidence interval of asymptotic $1 - \alpha$ level of μ_0 can be constructed as

$$\left(\hat{\beta}_{0jt} - d_n - \frac{u_{\frac{\alpha}{2}} \hat{\Omega}}{\sqrt{n}}, \hat{\beta}_{0jt} - d_n + \frac{u_{\frac{\alpha}{2}} \hat{\Omega}}{\sqrt{n}} \right)$$

where $u_{\frac{\alpha}{2}}$ is the upper $\frac{\alpha}{2}$ th quantile of $N(0, 1)$ and $d_n = I_{jt}^\top \{ Q(\hat{\delta})(0^\top, \hat{\delta}^\top)^\top - (0^\top, \hat{\delta}^\top)^\top \} / \sqrt{n}$.

4 Simulation Study

To evaluate the performance of the proposed estimator, we conducted a simulation experiment. In the simulation, we considered the following model

$$Y_i = Z_{1i}\beta_{(1)} + Z_{2i}\beta_{(2)} + Z_{3i}\beta_{(3)} + Z_{4i}\beta_{(4)} + \varepsilon_i, i = 1, \dots, n$$

where $\beta_{(1)} = (\beta_1, \beta_2)^\top$, $\beta_{(2)} = (\beta_3, \beta_4)^\top$, $\beta_{(3)} = (\beta_5, \beta_6)^\top$, $\beta_{(4)} = (\beta_7, \beta_8)^\top$, and the four factors specify $Z_{(1)} = (X_1, X_2)$, $Z_{(2)} = (X_3, X_4)$, $Z_{(3)} = (X_5, X_6)$, $Z_{(4)} = (X_7, X_8)$, where $X_{1i} \sim N(0, 1)$, $X_{2i} \sim N(0, 0.5)$, $X_{3i} = I[W_i > Q(W_i, 0.66)]$, $X_{4i} = I[W_i < Q(W_i, 0.33)]$, $X_{5i} = I[U_i > Q(U_i, 0.75)]$, $X_{6i} = I[Q(U_i, 0.25) < U_i < Q(U_i, 0.75)]$, $X_{7i} \sim N(0, 2)$, $X_{8i} \sim U[-2, 2]$, $\varepsilon_i \sim N(0, 1)$, $W_i \sim N(0, 0.8)$, $U_i \sim U[0, 2]$, $Q(W_i, w)$ represented the w th percentile of W_i . We set the first factor $Z_{(1)}$ to be in all candidate models. The other three factor may be or may not be present, so we had $2^3 = 8$ sub-models to be selected or averaged across. We generated $M = 1000$ random samples of size $n = 50, 100$ and 200 respectively. We calculated the coverage probability (CP) of claimed confidence interval resulted from the proposed smoothed FIC (SFIC) method. As a comparison, we calculated the coverage probability of confidence interval resulted from AIC, BIC and FIC. We also calculated the mean square error of the estimators

$$MSE(\hat{\beta}_n) = \frac{1}{M} \sum_{m=1}^M (\hat{\beta}_{j,m} - \beta_j)^2, j = 1, 2, 3, 4.$$

of the estimators. The simulation results for $\beta_1, \beta_2, \beta_3, \beta_4$ are reported in Table 1, the results of $\beta_5, \beta_6, \beta_7, \beta_8$ are similar and are not reported. From Table 1, we can clearly see the following facts:

- 1) For fixed n , the CP of SFIC is uniformly larger and the MSE of SFIC is uniformly smaller than AIC, BIC and FIC, which reflects that the proposed SFIC procedure performs better than AIC, BIC and FIC in terms of CP or MSE.
- 2) The proposed SFIC estimators become more accurate as the sample size increases.

Table 1. CP and MSE of the estimators

	β_1		β_2		β_3		β_4	
	CP	MSE	CP	MSE	CP	MSE	CP	MSE
n=50								
S-FIC	0.968	0.0323	0.961	0.0434	0.967	0.0860	0.955	0.1031
AIC	0.942	0.0381	0.874	0.0589	0.440	0.1388	0.440	0.1487
BIC	0.943	0.0377	0.792	0.0603	0.189	0.1194	0.184	0.1323
FIC	0.940	0.0420	0.946	0.0591	0.954	0.1367	0.945	0.1548
n=100								
S-FIC	0.957	0.0153	0.955	0.0196	0.966	0.0469	0.965	0.0494
AIC	0.953	0.0172	0.872	0.0257	0.424	0.0703	0.423	0.0751
BIC	0.951	0.0178	0.765	0.0290	0.139	0.0551	0.132	0.0611
FIC	0.953	0.0185	0.945	0.0267	0.956	0.0696	0.954	0.0700
n=200								
S-FIC	0.970	0.0069	0.965	0.0094	0.958	0.0247	0.960	0.0239
AIC	0.950	0.0081	0.893	0.0117	0.397	0.0367	0.409	0.0357
BIC	0.953	0.0080	0.776	0.0141	0.077	0.0273	0.081	0.0268
FIC	0.965	0.0079	0.948	0.0127	0.951	0.0357	0.963	0.0341

5 Proofs of Theorems

The following lemmas are needed to prove the theorems.

Lemma 1. Let

$$G_n(\beta_S) = \sum_{i=1}^n [\rho_\tau(\varepsilon_i - \frac{1}{\sqrt{n}}X_i^T \Pi_S^T \beta_S + \frac{1}{\sqrt{n}}X_i^T \beta_0) - \rho_\tau(\varepsilon_i)].$$

Then, under the conditions C1 and C2, it holds that

$$G_n(\beta_S) \xrightarrow{d} \frac{f(0)}{2} \beta_S^T \Pi_S \Sigma \Pi_S^T \beta_S + W^T \Pi_S^T \beta_S - W^T \beta_0 + \frac{f(0)}{2} \beta_0^T \Sigma \beta_0 - f(0) \beta_S^T \Pi_S \Sigma \beta_0.$$

for fixed β_S and β_0 , where $W \sim N(0, \tau(1 - \tau)\Sigma)$. And

$$\arg \min_{\beta_S} G_n(\beta_S) \xrightarrow{d} (f(0) \Pi_S \Sigma \Pi_S^T)^{-1} (\Pi_S W - f(0) \Sigma \beta_0)$$

Proof of Theorem 1. Note that $\hat{\beta}_S(\tau) = \arg \min_{\beta_S} \sum_{i=1}^n \rho_\tau(Y_i - (\Pi_S X)_i^T \beta_S)$, so $\hat{\beta}_S(\tau)$ minimizes

$$\sum_{i=1}^n [\rho_\tau(\varepsilon_i + X_i^T \beta_0 - (\Pi_S X)_i^T \beta_S) - \rho_\tau(Y_i - X_i^T \beta_0)]$$

Lemma 1 implies that

$$\sqrt{n} \left\{ \hat{\beta}_S(\tau) - \begin{pmatrix} \hat{\beta}_0 \\ 0 \end{pmatrix} \right\} \xrightarrow{d} -(\Pi_S \Sigma \Pi_S^T)^{-1} \Pi_S W / f(0) + (\Pi_S \Sigma \Pi_S^T)^{-1} \Pi_S \Sigma \begin{pmatrix} 0 \\ \delta \end{pmatrix}.$$

The proof of Theorem 1 is completed by replacing $W/f(0)$ with G .

Proof of Theorem 2. It can be verified that

$$\begin{aligned} \sqrt{n}(\hat{\beta}_{0jt} - \beta_{0jt}) &= \sum_S w(S|\hat{\delta}) \left\{ -I_{jt}^T R_S W_n / f(0) + I_{jt}^T R_S \Sigma \begin{pmatrix} 0 \\ \delta \end{pmatrix} - I_{jt}^T \delta + o_p(1) \right\} \\ &= I_{jt}^T \sum_S w(S|\hat{\delta}) R_S \Sigma \begin{pmatrix} -[I, 0] \Sigma^{-1} W_n / f(0) \\ -[0, I] \Sigma^{-1} W_n / f(0) + \delta \end{pmatrix} - I_{jt}^T \delta + o_p(1) \\ &= -I_{jt}^T \Sigma^{-1} W_n / f(0) + I_{jt}^T \left[Q(\hat{\delta}) \begin{pmatrix} 0 \\ \hat{\delta} \end{pmatrix} - \begin{pmatrix} 0 \\ \hat{\delta} \end{pmatrix} \right] + o_p(1) \\ &\xrightarrow{d} -I_{jt}^T \Sigma^{-1} G + I_{jt}^T \left[Q(\Delta) \begin{pmatrix} 0 \\ \Delta \end{pmatrix} - \begin{pmatrix} 0 \\ \Delta \end{pmatrix} \right] \end{aligned}$$

This complete the proof.

References

1. Akaike, H.: Maximum Likelihood Identification of Gaussian Autoregressive Moving Average Models. *Biometrika* 22, 203–217 (1973)
2. Buckland, S.T., Burnham, K.P., Augustin, N.H.: Model Selection: an Integral Part of Inference. *Biometrics* 53, 603–618 (1997)
3. Fan, J.Q., Li, R.: Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties. *Journal of the American Statistical Association* 96, 1348–1360 (2001)
4. Hjort, N.L., Claeskens, G.: Frequentist Model Average Estimators. *Journal of the American Statistical Association* 4, 879–899 (2003)
5. Hjort, N.L., Claeskens, G.: Focussed Information Criteria and Model Averaging for Cox's Hazard Regression Model. *Journal of the American Statistical Association* 101, 1449–1464 (2006)
6. Huang, Y.J.: Quantile Calculus and Censored Regression. *The Annals of Statistics* 38(3), 1607–1637 (2010)
7. Kai, B., Li, R., Zou, H.: Local Composite Quantile Regression Smoothing: an Efficient and Safe Alternative to Local Polynomial Regression. *J. Roy. Statist. Soc. Ser. B* 72, 49–69 (2010)
8. Koenker, R.: *Quantile Regression*. Cambridge University, London (2005)
9. Liang, H., Zou, G.H., Wan, A.T.K., Zhang, X.Y.: On Optimal Weight Choice in a Frequentist Model Average Estimator. *Journal of the American Statistical Association* 106(495), 1053–1066 (2011)
10. Leeb, H., Pötscher, B.M.: The Finite Sample Distribution of Post-Model-Selection Estimators and Uniform Versus Non-uniform Approximations. *Econometric Theory* 19, 100–142 (2003)
11. Leung, G., Barron, A.: Information Theory and Mixing Least-Squares Regressions. *IEEE Transactions on Information Theory* 8, 3396–3410 (2006)
12. Leeb, H., Pötscher, B.M.: Can One Estimate the Unconditional Distribution of Post-Model-Selection Estimators. *Econometric Theory* 24, 338–376 (2008)
13. Pang, L., Lu, W.B., Wang, H.J.: Variance Estimation in Censored Quantile Regression via Induced Smoothing. *Computational Statistics and Data Analysis* 56, 785–796 (2012)
14. Schwartz, G.: Estimating the Dimension of a Model. *The Annals of Statistics* 6, 461–464 (1978)
15. Tibshirani, R.: Regression Shrinkage and Selection via the Lasso. *J. Roy. Statist. Soc. Ser. B* 58, 267–288 (1996)
16. Wang, H., Leng, C.: A note on Adaptive Group Lasso. *Computational Statistics and Data Analysis* 52, 5277–5286 (2008)
17. Yang, Y.: Adaptive regression by mixing. *Journal of the American Statistical Association* 96, 574–586 (2001)
18. Yuan, M., Lin, Y.: Model Selection and Estimation in Regression with Grouped Variables. *J. Roy. Statist. Soc. Ser. B* 68, 49–67 (2006)
19. Zhang, X.Y., Liang, H.: Focused Information Criterion and Model Averaging for Generalized Additive Partial Linear Models. *The Annals of Statistics* 39(1), 174–200 (2011)
20. Zhao, P., Rocha, G., Yu, B.: The Composite Absolute Penalties Family for Grouped and Hierarchical Variable Selection. *The Annals of Statistics* 37, 3468–3497 (2009)

Analytical Studies and Experimental Examines for Flooding-Based Search Algorithms

Hassan Barjini, Mohamed Othman*, Hamidah Ibrahim, and Nur Izura Udzir

Department of Communication Technology and Networks, Universiti Putra Malaysia,
43400 UPM Serdang, Selangor D.E., Malaysia

`hassan.barjini@gmail.com`, `mothman@fsktm.upm.edu.my`

Abstract. Flooding search has no knowledge about network topology and files distribution, thus it offers an attractive method for file discovery in dynamic and evolving networks. Although flooding can achieve high coverage but it produces exponentially redundant messages in each hop. To improve this searching scheme and reduce redundant messages, there are many flooding-based search algorithms in literature. This paper provided analytical study for flooding and flooding-based algorithms. The analytical results provided the best hop threshold point for the optimum growth rate coverage and redundant messages from flooding and flooding-based algorithms. The simulation experiments validated analytical results.

Keywords: Peer-to-peer, Flooding-based searching, Redundant messages.

1 Introduction

Flooding can achieve high coverage but it produces exponentially redundant messages in each hop. Consequently, the growth of redundant messages limits system scalability and causes unnecessary traffic in networks. In order to tackle the negative aspect of flooding, there are many solutions in the literature. The major solutions are those who used flooding algorithm as a basic part of their schemes. These schemes named as flooding-based algorithms. Flooding-based algorithm can be classified as:

1. Time-To-Live (TTL) Limit-Based Flooding (TLBF). Those who limit flooding by number of TTL; (Expanding Ring, Blocking Expanding Ring, Iterative Deepening, ...)
2. Probabilistic Limit-Based Flooding (PLBF). Those who limit flooding by choosing sample neighbors; (Random Walk, Random Breadth-first-search, Teeming, ...)

* The author is also an associate researcher at the Lab of Computational Science and Informatics, Institute of Mathematical Research (INSPEM), Universiti Putra Malaysia.

3. Hybrid Limit-Based Flooding (HLBF). Those who limit flooding by either utilizing hybrid overlay or super-peer techniques (Dynamic Query, AntSearch, ...)

This paper covers analytical studies for flooding and flooding-based algorithms. The analytical results provided the best threshold point of hop for optimum coverage growth rate and redundant messages in these algorithms. The simulation experiments validated analytical results.

The rest of this paper is organized into six sections: Section 2 reviews related work. Section 3 investigate flooding across hops. Section 4 proposed performance metric. Section 5 performance evaluations, and section 6 presents the conclusion.

2 Related Work

Researchers have attempted to alleviate the overshooting messages of flooding search and improve acceptable versions of flooding. Many alternative schemes have been proposed to address this problem. Expanding Ring (ER) is the pioneer choice of these endeavors. This technique confines searching scope by limiting TTL value. Although this scheme mitigate loads and traffic. However it still produce many duplicate messages and give no guarantee of successful queries [1][2].

Random Walk (RW) [2] reduces overshooting messages. In this scheme, no nodes are visited more than once, so it gains minimum search cost, loads, and traffic. Despite these merits, it is almost non-deterministic, non-reliable, and has high variable performance. To overcome RW's faults, there is an extended version of RW, which is called Random Breadth-First-Search [3] or Modified- Breadth-First-Search [4] or Teeming [5]. In this scheme, at each step, the node propagates the query messages only to a random subset of its neighbors. Although when compared to flooding, the overshooting messages dramatically decrease, the algorithm is probabilistic and a query might not reach some large network segments [6].

Some of the other attempts implement a hierarchical structure and use super peers [7]. Gnutella2 [8] and KaZaa/FastTrack [9] are based on the hierarchical structure. These techniques divide peers into two groups: super peers and leaf peers. Super peer acts as a server who receives queries from its clients or other super peers. Clearly, this technique reduces traffic and improves the search efficiency. However, the main drawback of this approach is its vulnerability to single point failure and the limitation of the number of clients supported by each super peer.

3 Flooding across Hops

Flooding conducted in a hop-by-hop fashion. By increasing of hops, it gains new peers, and generates more messages. Part of these messages is redundant messages. This section investigated the trend of coverage growth rate and redundant messages in flooding.

3.1 Trend of Coverage Growth Rate in Flooding

Assume an overlay network as a Random Graph $G_{n,p}$. Each node is represented as a peer, and they are connected to each other by edges. The degree of each peer represents the number of its immediate neighbors. Assume that the graph has n nodes with the average degree d , suppose d is greater than 3 the number of messages broadcasting from each peer in hop i as given in [10], is:

$$M_{F,i} = d(d - 1)^{i-1} \tag{1}$$

Thus, the total messages broadcasting up to hop t is equal to:

$$TM_{F,t} = \sum_{i=1}^t d(d - 1)^{i-1} \tag{2}$$

Loop nodes or cyclic paths are grouping of nodes linked together as a ring fashion. In Gnutella and other internet topology, there are many cyclic paths. If there are no cyclic paths [11] or loop nodes in the topology, then the total number of new peers visited so far is equal to:

$$TP_{F,t} = \sum_{i=1}^t d(d - 1)^{i-1} \tag{3}$$

Assume $A = d - 1$ Thus, to observe the coverage growth rate of messages [12] in hop t :

$$CGR_{F,t} = \frac{TP_{F,t}}{TP_{F,t-1}} = 1 + \frac{(d - 1)^{t-1}}{(\sum_{i=0}^{t-2} (d - 1)^i)} = \frac{A^t - 1}{A^{t-1} - 1} \tag{4}$$

The discrete derivative of a function $f(n)$, with respect to n , define as:

$$\Delta_n f(n) = f(n) - f(n - 1) \tag{5}$$

Thus, derivative of CGR_t with respect to t , lead to:

$$\Delta_t (CGR_{F,t}) = CGR_{F,t} - CGR_{F,t-1} = -\frac{A^{t-2}(A - 1)^2}{(A^{t-1} - 1) \times (A^{t-2} - 1)} \tag{6}$$

The value of $\Delta_t (CGR_{F,t})$ is always negative. Thus, the (4) is always in descending order. By increasing the value of t (hops), the value of $CGR_{F,t}$ decreases, hence the maximum value of $CGR_{F,t}$ is visited in second hop. Therefore, we can show that:

$$(CGR_{F,2}) > (CGR_{F,3}) > (CGR_{F,4}) > \dots \tag{7}$$

3.2 Trend of Redundant Messages in Flooding

Redundant messages in each topology are generated by loop nodes. Assume that there is a loop in each hop of the defined topology, and that loop is started from second hop.

$$P_{F,2} = [d(d - 1)^2] - 1 \tag{8}$$

By considering a loop, the number of new peers in the third and fourth hop respectively are:

$$P_{F,3} = [(d(d-1)^2 - 1)(d-1)] - 1 = d(d-1)^3 - (d-1) - 1 \quad (9)$$

$$P_{F,4} = d(d-1)^4 - (d-1)^2 - (d-1) - 1 \quad (10)$$

Thus, the number of new peers in hop t becomes:

$$P_{F,t} = d(d-1)^t - \sum_{i=0}^{t-2} (d-1)^i \quad (11)$$

The number of redundant messages can be defined as the difference between the number of messages and the number of new peers visited in hop t .

$$R_{F,t} = M_{F,t} - P_{F,t} = \frac{A^{t-1} - 1}{(A-1)} \quad (12)$$

The derivative of (12) refer to (5) can be shown as:

$$\Delta_t R_{F,t} = R_{F,t} - R_{F,t-1} = A^{(t-2)} \quad (13)$$

The value of (12) is always positive. Thus, it is always in ascending order. By increasing the value of t (hops), the value of R_t increases. Hence, the minimum value of R_t is visited in second hop. We can show that:

$$R_{F,2} < R_{F,3} < R_{F,4} < \dots \quad (14)$$

The analytical study identified two important points:

1. The coverage growth rate of messages for flooding in defined topologies has a maximum value in the second hop, by increasing the hops the value of coverage growth rate decreases.
2. The redundant messages for flooding in the same topologies at low-hops are very low but by increasing hops their values increase exponentially.

4 Proposed Performance Metric

Inspired by the reverse trend in coverage growth rate and redundant messages in flooding schemes (7) and (14), we propose a new metric, which called the *critical metric*.

4.1 Critical Metric

Here is the description of *critical metric*. The definition is:

$$CM_{x,t} = \frac{R_{x,t}}{CGR_{x,t}} \quad (15)$$

This metric is valid for low-hops, because in high-hops the value of coverage growth rate becomes close to 0. This metric can evaluate the performance of flooding and flooding-based algorithms in each hop.

4.2 Critical Metrics in Flooding Algorithms

The coverage growth rate of flooding is equal to:

$$CGR_{F,t} = \frac{A^t - 1}{A^{(t-1)} - 1} = \frac{A^t}{A^{(t-1)} - 1} - \frac{1}{A^{(t-1)} - 1} > \frac{A^t}{A^{(t-1)} - 1} > A \quad (16)$$

The value of the redundant messages in flooding is equal to:

$$R_{F,t} = \frac{A^{(t-1)} - 1}{A - 1} < A^{(t-1)} \quad (17)$$

Thus the *critical metric* for the flooding algorithm is equal to:

$$CM_{F,t} = \frac{R_{F,t}}{CGR_{F,t}} = \frac{A^{(t-1)}}{A} = A^{(t-2)} \quad (18)$$

Assume that the TTL value used in flooding is k , thus the total *critical metric* for all hops up to k is equal to:

$$TCM_{F,k} = \sum_{i=1}^k CM_{F,i} = \sum_{i=1}^k A^{(i-2)} = \frac{A^k - 1}{A^2 - 1} \quad (19)$$

4.3 Critical Metrics in ER Algorithm

The ER is successive flooding; assuming that the *TTL* value in ER started from 1 and is incremented by 1 up to l , where l is less than k and greater than 2. Thus the *critical metric* for the ER from 1 to l is equal to:

$$TCM_{ER,l} = \sum_{t=1}^l \sum_{i=1}^t CM_{F,i} = \frac{A^{(l+1)} - A - l(A - 1)}{(A^2 - A)(A - 1)} \quad (20)$$

4.4 Critical Metrics in BER Algorithm

The BER is an extended version of the ER, which is not rebroadcast from source node, but rather in each new round it is rebroadcast from the nodes of the last attempts. Hence, its *critical metric* for ($l < k$) referring to (19) is equal to:

$$TCM_{BER,l} = \sum_{i=1}^l CM_{F,i} = \frac{A^l - 1}{A^2 - A} \quad (21)$$

4.5 Critical Metrics in RW Algorithm

RW is modified version of flooding. It is a Probabilistic Limit-Based Flooding (PLBF) type. In the RW, the average degree of d is equal to 2. The value of the *critical metric* is equal to:

$$TCM_{RW,l} = \sum_{i=2}^l A^{(i-2)} = \sum_{i=2}^l 1^{(i-2)} = l - 2 \quad (22)$$

4.6 Critical Metrics in Teeming Algorithm

Modified Breadth-First-Search or teeming is also a Probabilistic Limit-Based Flooding (PLBF) type. There is a fixed probability denoted by θ for selecting a particular neighbors.

$$TCM_{T,l} = \sum_{i=2}^l (A \times \theta)^{(i-2)} \quad (23)$$

The result shows that the value of *critical metric* in the teeming algorithm is almost $\theta^{(t-2)}$ of the *critical metric* in the flooding.

4.7 Analytical Results

Our analytical study shows critical metrics in TLBF methods decreases linearly refer to (19), (20) and (21). However, refer to (22) and (23) in PLBF methods for random walker and teeming it decreases exponentially. It proved that PLBF has a better performance compared to TLBF.

$$TCM_{F,l} > TCM_{ER,l} > TCM_{BER,l} > TCM_{RW,l} > TCM_{T,l} \quad (24)$$

5 Performance Evaluations

We used two metrics for evaluation our experiments:

1. Queries success rate
2. Number of redundant messages

The queries success rate is defined as the probability that a query is successful. This metric evaluates the efficiency and quality of the search algorithm.

When a multiple message with the same message *id* is sending to peer by its multiple neighbors, all except for the first message, are considering as a redundant message. The redundant messages are absolute overhead and it is the rate of system scalability.

5.1 Experimental Results

The evaluation compared the performance of flooding, ER, BER and teeming algorithms with number of redundant messages and query success rate. The value of probability θ in the teeming algorithm is set to 0.3 or 30% (Teeming_30). The research mainly uses two topology traces which collected by clip2 [13].

Figure 1 present the trend of the redundant messages for all algorithms. They show that the teeming algorithm reduces redundant messages by almost 90% compared to flooding. It shows that the BER and ER reduce the redundant messages by over 70% compared to flooding.

The trend in the success rate for each algorithm in different topologies is presented in Figure 2. As expected, the ER and BER improved the success rate

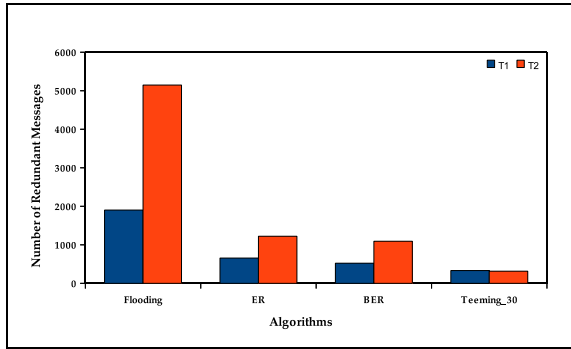


Fig. 1. Number of redundant messages for each algorithm at different topologies

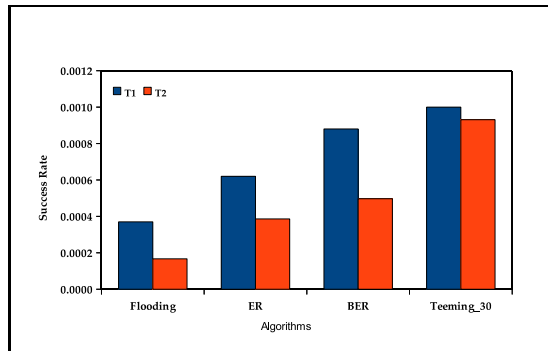


Fig. 2. Success Rate for each algorithm at different topologies

by 2.0 and 2.7 times compared with flooding, the reason refers to their limitation of the *TTL* values. The result illustrates the exact difference between the ER and BER algorithms. The ER collected many duplicate messages because it must start from a source peer in each round of processing, whereas the BER does not follow this procedure.

6 Conclusion

In this paper we provided an analytical study for flooding and flooding-based search algorithms. The study consider the main problem of flooding search. Flooding search has two different characteristics. First in low-hops it produces low redundant messages and high coverage growth rate in network. Second in high-hops it generates high redundant messages and low coverage growth rate. There is no exact threshold point of hop between these low-hops and high-hops in flooding and flooding-based algorithms.

The paper proposed a new metric called *critical metric*. *Critical metric* represents the state of the redundant messages and coverage growth rate at each hop of a flooding and flooding-based algorithms. Thus it can evaluate the performance of a search algorithm in each hop. Consequently this metric is suitable tool for comparing the performance of flooding and flooding-based search algorithms along the hops. This metric is also proper tool for comparison the performance of search in each flooding and flooding-based algorithms.

Acknowledgements. This work was partially supported by the Malaysian Ministry of High Education under the Fundamental Research Grant Scheme, FRGS No: FRGS/1/11/SG/UPM/01/1.

References

1. Barjini, H., Othman, M., Ibrahim, H., Udzir, N.: Shortcoming, problems and analytical comparison for flooding-based search techniques in unstructured p2p networks. *Peer-to-Peer Networking and Applications* 5(1), 1–13 (2012)
2. Lv, Q., Cao, P., Cohen, E., Li, K., Shenker, S.: Search and replication in unstructured peer-to-peer networks. In: *Proceedings of the 16th International Conference on Supercomputing*, pp. 84–95. ACM (2002)
3. Zeinalipour-Yazti, D., Kalogeraki, V., Gunopulos, D.: pfusion: A peer-to-peer architecture for internet-scale content-based search and retrieval. *IEEE Transactions on Parallel and Distributed Systems* 18(6), 804–817 (2007)
4. Kalogeraki, V., Gunopulos, D., Zeinalipour-Yazti, D.: A local search mechanism for peer-to-peer networks. In: Nicholas, C., Grossman, D., Kalpakis, K., Qureshi, S., van Dissel, H., Seligman, L. (eds.) *CIKM 2002*, pp. 300–307. Association for Computing Machinery (2002)
5. Dimakopoulos, V., Pitoura, E.: On the performance of flooding-based resource discovery. *IEEE Transactions on Parallel and Distributed Systems* 17(11), 1242–1252 (2006)
6. Zeinalipour-Yazti, D., Kalogeraki, V., Gunopulos, D.: Information retrieval techniques for peer-to-peer networks. *Computing in Science and Engineering*, 20–26 (2004)
7. Yang, B., Garcia-Molina, H.: Designing a super-peer network. In: *Proceedings of the 19th International Conference on Data Engineering*, pp. 49–60 (2003)
8. Gnutella2 (2003), <http://www.gnutella2.com> (January 2006)
9. KaZaa, <http://www.kazaa.com/>
10. Aberer, K., Hauswirth, M.: An overview on peer-to-peer information systems. In: *Workshop on Distributed Data and Structures, WDAS 2002*, pp. 1–14 (2002)
11. Zhu, Z., Kalnis, P., Bakiras, S.: Demp: A distributed cycle minimization protocol for peer-to-peer networks. *IEEE Transactions on Parallel and Distributed Systems* 19(3), 363–377 (2007)
12. Jiang, S., Guo, L., Zhang, X., Wang, H.: Lightflood: Minimizing redundant messages and maximizing scope of peer-to-peer search. *IEEE Transactions on Parallel and Distributed Systems* 19(5), 601–614 (2007)
13. Clip2: Distributed search solutions, <http://public.yahoo.com/~lguo/download/gnutella-trace/xmlfiles.tar.gz>, <http://public.yahoo.com/~lguo/download/gnutella-trace/sample-query.dat> (last visited January 19, 2012)

Innovative Study on the Multivariate Statistical Analysis Method of Chromatography Economy Analysis

Shibing You^{*}, Yuan Hui, Xue Yu, and Lili Bao

Economics and Management School, Wuhan University, China, 430072
{sbyou, lilibao}@whu.edu.cn,
{huiyuan027, yxfish123}@126.com

Abstract. The chemistry chromatography analysis method was introduced into the multivariate statistics subject in the economy research in order to clarify and simplify the complex economic phenomena. The innovative multivariate statistic method was proposed based on the analysis and the application of the traditional multivariate statistics theory. The present work introduced the basic concept, principle and basic theory of chemistry chromatography analysis method. With the replacement of the economy implication, the chromatography economy analysis method was expected to apply in the future statistic science.

Keywords: Chromatography Economy Analysis, Multivariate Statistical Analysis, Economic Meaning Replacement.

1 Introduction

Due to the complexity of the social economy phenomenon, the studied objects always need to be sorted and handled by some standard in order to clarify the properties and characteristics of the most social economic phenomenon. Therefore, the classification process is a crucial link in applied economic research.

The multivariate statistical analysis is a method which studies the mutual dependence of many random variables and the internal statistical regularity to classify and simplify the research objects. The method has been widely applied to many fields, such as the geology, hydrology, meteorology, medicine, industry, agriculture, and economics and so on. Of course, there is much further development space for the current multivariate statistical method. Especially, the new multivariate statistical method shall be explored to satisfy the appeal for the more complicated social economy classification.

The chromatography analysis is a method to separate the compositions in chemistry field. You et al. [1-4] initially introduced the method into social economics field, elaborated the primary research thought and innovatively proposed the application concept of the chromatography economy analysis method. In this paper, we introduced the chromatography analysis method is how to apply in the multivariate statistical analysis.

^{*} Corresponding author.

2 Comparison between the Traditional Statistical Analysis Method and the Chromatography Economy One

2.1 The Characteristic Analysis of Traditional Multivariate Statistical

The traditional multivariate statistical analysis includes many methods, such as clustering analysis, distinguishable analysis, principal components analysis, factor analysis, correspondence analysis and canonical correlation analysis. They have 3 characteristics.

The first characteristic is the relatively strong subjectivity. For example, on the basis of the distance of the sample or index, we classify the clustering and distinguishable analysis. However, the definitions of the distance have the disadvantage of treating the indexes equally and the inconsistent results for different distance algorithms. Moreover, the standard has not been formed currently for the evaluation to the distance algorithms [5].

The second characteristic is low resolution. As for the traditional multivariate statistical method, there are the subjective factors and the cross of the nature and characteristics of the separated objects during the classification process so that the resolution of the separation and classification results may be relatively low [6]. For example, the principal components analysis and the factor analysis may integrate the multiple indexes into the comprehensive index serving as the representative of the features of any economic thing through applying the mathematics and measurement analysis methods to the indexes [7].

The third characteristic is the lack of effectiveness. Because the traditional multivariate statistical method focuses on some certain characteristic factor of the separated object so as to carry out the induction and discrimination and achieve the separation purpose by utilizing the mathematics and measurement methods. on the other hand, the interactions between the social and economic things as well as the social and economic things and their living environments have been neglected so as to become the simple classification being away from the social and economic system and reflect more representations of the things rather than involve less investigation to the deep attributes of the separating objects and be lack of the judgment to the essence.

2.2 The Advantage for the Chromatography Economy Analysis Method

First, the qualitative and quantitative analysis may be carried out using chromatography economy analysis method [8]. Among the many separation technologies in the chemical field, the chromatography analysis method is the unique separation technology which may carry out the qualitative and quantitative analysis simultaneously. The current multivariate statistical methods lay particular stress on the quantitative analysis so that the separation may be carried out by studying the mathematics and measurement features of the separating object and the qualitative and quantitative separation of the social and

economic things may possibly be carried out synchronously by introducing the chromatography analysis method into the multivariate statistical science. As shall be one major breakthrough to the traditional multivariate statistics.

Second, the chromatography economy analysis method is combined with other methods. Each chemical separation technology has its own defects; the chromatography analysis method is also not the exception, whose judging ability to the separated objects is relatively poor. However, the advantages of the chromatography analysis lie in its being combined with other technologies.

Third, the chromatography economy analysis method has wide application range. The current multivariate statistical methods are generally lack of the broad applicability. On the other hand, the chromatography economy analysis method may possibly break the boundary and limitation of the multivariate statistical methods and widely applicable to the separation and classification to various kinds of social and economic things.

Last, the chromatography economy analysis method has high separation efficiency. The chromatography analysis method may separate the mixture including more than ten substances and may also separate successfully the mixed components whose properties are extremely close. In addition, the very trace samples may also measure. Thus, it has the high separation efficiency in contrast to other methods. The fundamental analysis established with the chromatography economy analysis method

2.3 Basic Concept

Stationary Phase and Mobile Phase

In the chromatography economy analysis method, the stationary phase concept may be standardized as the phase in the economy chromatographic column blocks the advancing of the economic things for separating. The basic requirement for the stationary phase is the good selectivity to the components. Namely, there are very strong attraction to some certain components and relatively weak or even no attraction to other components so that the various components in the sample may be separated each other.

The mobile phase concept may be standardized as, the phase in the economic chromatographic column drives the advancing of the economic things for separating, which is the driving force setting playing the separation role and the internal decomposition and refining to the core characteristics of the separating objects in the separation sub-system. The basic requirements for the mobile phase are to push of the components away from the stationary phase and advancing. The design of the mobile phase may directly influence the separation speed, the separation efficiency or even the measurement success.

Distribution Ration

In the chromatography economy analysis method, the distribution ratio concept may be standardized as, in some certain economic environment, the differential contrast of the preferences for various components in the stationary phase and the mobile phase

or the ratio of the action intensities for the separating components in the selected stationary phase and the mobile phase.

Chromatographic Peak

The chromatographic peak is a differential outflow curve occurring while the components being in the outflow mode. As shown in the Fig.1, while the outflow component concentration reaches the maximum value, the chromatographic peak location may occur.

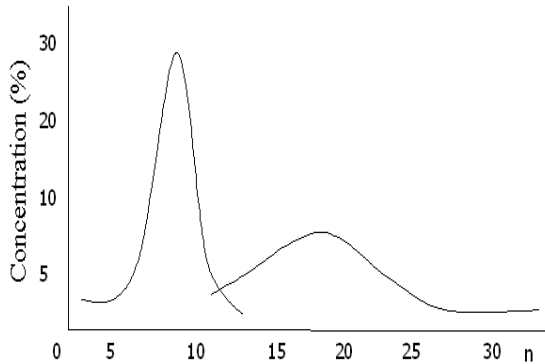


Fig. 1. Dependence of concentration on retention time

In the chromatography economy analysis method, the chromatographic peak concept may be standardized as the curve diagram including some certain parameters relationship may be formed while the separated economic things flow out of the column outlet through the separation in the chromatographic column.

Retention Time

In the chromatography economy analysis method, the retention time concept may be standardized as, the maximum time interval from the time of coming out the first sample for the separated economic things to the time when the sampling concentration being up to the maximum value. As for the different components, the larger the gap between their retention times is, the better the separation effect is.

2.4 The Principle of the Chromatography Economy Analysis Method

2.4.1 The Principle of the Chromatography Analysis in the Chemistry

The chemical chromatography separation is a very complicated process, which is the comprehensive representation of the thermodynamics processes and the dynamics processes in the chromatography system. The thermodynamic properties of the components, the mobile phase and the stationary phase may cause the various components have different distribution coefficients in the mobile phase and the stationary

phase. The distribution coefficient sizes may reflect the components' dissolution – volatilization or adsorption – desorption ability in the stationary phase. Any component corresponded the large distribution coefficient has the strong dissolution or adsorption ability in the stationary phase so that its moving speed inside the column may be slow. Conversely, any component corresponded the small distribution coefficient has the weak dissolution or adsorption ability in the stationary phase so that its moving speed inside the column may be quick. After some certain period, the differential speed transition may occur inside the column for the components due to the difference between the distributions coefficients so as to achieve the separation goal.

After the separation of the mixtures has been realized through the chromatographic column, the qualitative and quantitative analysis may be carried out to the separated components through the detecting system so as to determine the components and contents in the mixture.

2.4.2 The Basic Principle of the Chromatography Economy Analysis Method

The basic action principle of the chromatography economy analysis method may be abstracted as the different acting forces in the social economy for different things due to their own unique properties and characteristics as well as their location environments so that the things may be separated.

The “chromatographic column” shall be designed for the separation of the social and economic things so that the separation may be carried out to the social and economic things under the pushing force of the mobile phase and through the chromatographic column. Then the measurement and analysis may be carried out to the separated social and economic things to determine the properties of the separated components.

The thermodynamics and dynamics processes are two absolutely necessary processes in the chromatography economy analysis. While the separation is carried out to the economic things, the different interactions may occur for the things with different properties or characteristics as well as the stationary phase and the mobile phase in the chromatographic column. In other words, the components have different preferences as for the two phases to form the thermodynamics process. On the other hand, the mobile phase pushes the components to advance inside the chromatographic column and the different selected components have different advancing speeds to form the dynamics process.

In the early stage research on the chromatography economics, the regularly advancing cases for the components are mainly taken into account in the dynamics process for simplifying the analysis. However, the diffusion process (namely the components may flows in the opposite direction to the mobile phase moving or irregular motion cases due to the random disturbance), which may be the main constituent part in the future research. Such diffusion cases may fully take the influence of the random disturbance factors into account so that the research on the chromatography economics may furthermore truly reflect the social and economic phenomenon.

3 The Content in the Chromatography Economy Analysis System

One complete set of chromatography economy analysis method system includes several sub-systems, including the most important two sub-systems, the separation sub-system (the chromatographic column) and the detection sub-system. These two sub-systems may perform some functions in the entire system, respectively. The goal of the separation sub-system is to realize the separation of the economic things with different properties and characteristics, in which there are the mobile phase and the stationary phase and the chemical chromatographic column serves as the setting basis and which is the main research object in the previous stage research. On the other hand, the detection sub-system may carry out the quantitative and qualitative detection and analysis to the separated samples or indexes so as to determine the specific characteristics of the various groups of samples or indexes after the separation has been completed. The vaporization system is the pressure system to continuously fill gas into the chromatographic column. And the recording system may record the various data for the components, such as the outflow curves for the samples coming out of the column outlet, the retention time and so on.

The chemical chromatography analysis principle has been fully applied in the internal framework of the separation system. And the social and economic things for separating are taken as the “components”, which may be selected continuously in the “chromatographic column” according to the distribution ratio. For example, the proper mobile phase is set for the classification of the population in the economy so that it may serve as the driving force to prompt the economy population for separating advancing in the imaginary chromatographic column, and the proper stationary phase may be set so as to select some individuals in contrast with the stationary phase to block the advancing. Because the various economic populations have different characteristics, the natural separation may be formed.

Of course, the appeals for the classification in the social and economic system are far more than such simplicity for the classification to the population. On the other hand, there are a lot of uncertainties in such complicated system so as to provide one challenge for the standardization to the chromatography economy analysis method and form the breakthrough point for the chromatography economy analysis method to transcend the current multivariate statistical analysis method. The chromatography analysis method has been applied relatively maturely in the chemical field and has been widely utilized in other natural science fields. Such analysis method includes many branches and a lot of concepts, which may be introduced into the social and economic fields and the different chromatography separation sub-systems may be designed in view of the different economic sub-systems. Thus, it may be predicted that such method will have broad applicability and the chromatography economics will play roles in many fields such as the market research, the statistical grouping, the industrial economics, the regional economics, the financial securities analysis, the complex social phenomenon analysis, the macro and micro economic situations analysis, the sampling theory and application, the international relationship analysis, the consumer behavior analysis, the enterprise culture, the enterprise financial

analysis, the human resources analysis, the risk management, the disaster warning and evaluation and so on.

4 The Primary Problems to be Solved with the Chromatography Economy Analysis Method

The chromatography analysis method introduced into the multivariate statistical analysis field is the major research topic integrating the greatly different subjects. The challenges and difficulties in the subject amalgamation are faced. Currently, there are 3 aspects of specific difficulties.

The first difficulty is how to extract the core information. As for the complicated economic things for classification, we shall first handle them and extract the basic characteristics and attributes to get rid of the redundant information which may not influence the classification essence and retain the core information [9]. The principal components analysis thought may be used for reference and uniformization may be carried out to the samples and variables to select the representative samples or variables and simplify the separation process. In addition, this process is also similar to the “extraction” process in the chemistry. The preliminary treatment may be carried out to the mixture and the complicated mixture may be purified. Any impurity not belonging to the separation goal may be removed. The further research and exploration shall be carried out to realize this process.

The second difficulty is how to establish the model of mathematics. The qualitative separation classification may be carried out directly to some things according to their properties. However, the separation classification to more things shall be established based on some certain mathematic analysis [10]. For achieving the relatively good separation effect, the accurate statistics may be required to serve as the separation standard so as to eliminate any interference from the subjective factors [11]. Mathematics is one important science tools and the important contents being involved during the chromatography economy analysis method establishment process. The chromatography economy analysis process includes the analysis to various “reaction” data. We necessarily study its probability theory basis so that the chromatography economy analysis shall have the scientific objective basis.

The third difficulty is how to design the separation system accurately. While the separation sub-system is constructed, the design and selection of the mobile phase and the stationary phase may be the key factor influencing the classification results. The reacting object and the reaction speed may be different due to the different characteristics and attributes of the various things so as to possibly result in the hysteretic nature. Thus, we may possibly obtain the ideal testing time and the correct classification results through the design and testing again and again while the chromatography separation method is utilized.

Although the current multivariate statistical analysis method has been widely applied, the appeals for the classification in the economic field may not fully met and realized. And it is necessary to perfect the multivariate statistical science by introducing the new statistical method. Moreover, the chromatography analysis

method has been one kind of separation analysis method which has been widely applied in the natural science field and has the very strong practice advantages. And its development has been mature and the application prospects are good. Besides it has abundant practice experiences. it not only is feasible for being applied in the social and economic fields but also becomes the necessity for the subject amalgamation trend. The multivariate statistical analysis development may be promoted by introducing the chromatography economy analysis method analysis to provide one kind of scientific and effective new method for meeting the appeals for the classification in the practical economic life.

References

1. You, S.B., Wu, B., Shen, P., Mei, M., Su, Z.H.: Theoretical prospect of the complex economic phenomenon classification method innovation – reference and thinking based on the chemical chromatography analysis method. *Tong Ji & Jue Ce* 7, 4–7 (2011)
2. You, S.B., Wu, B., Mei, M.: Foundation research on applying the chromatography analysis principle in the economic field – taking the financial securities investment and fast moving consumer goods industry as examples. *Tong Ji & Jue Ce* 11, 4–7 (2011)
3. Shen, P., Zhang, P., Mao, K.Y., Li, G.Q., You, S.B.: Chromatography economy analysis method replacement series research, distribution ratio. *Tong Ji & Jue Ce* 17, 4–7 (2011)
4. You, S.B., Bao, L.L., Zhong, S.Y., Guan, X., Wang, L.X.: Chromatography economy analysis method replacement series research, trays theory. *Tong Ji & Jue Ce* 1, 4–6 (2012)
5. Florez, L.R.: Effects of missing data in credit risk scoring. A comparative analysis of methods to achieve robustness in the absence of sufficient data. Houndmills, Basingstoke, Hants, RG21 6XS, United Kingdom, pp. 486–501 (2010)
6. Noorossana, R., Eyvazian, M., Amiri, A., et al.: Statistical monitoring of multivariate multiple linear regression profiles in phase i with calibration application. *Quality and Reliability Engineering International* 26, 291–303 (2010)
7. Ren, Q., Li, S., Qiao, D., et al.: Application of key factor analysis method for elastic coefficient of highway freight transportation, pp. 281–297. Chengdu, China (2010)
8. Su, L.Q., Zheng, Y.J.: Chromatography analysis method, pp. 23–26. Tsinghua University Press, Beijing (2009)
9. Yu, F.J., Tsou, C.S., Huang, K.I., et al.: An economic-statistical design of x control charts with multiple assignable causes. *Journal of Quality* 17, 327–338 (2010)
10. Liu, Z.L., Li, C.J.: Quantitative Analysis on the Creation of Dualistic Interdisciplinary Science. *Studies in Dialectics of Nature* 20, 42–45 (2004)
11. Zhou, Y., Wan, A.T., Xie, S., et al.: Wavelet analysis of change-points in a non-parametric regression with heteroscedastic variance. *J. of Econometrics* 159, 183–201 (2010)

Optimization of Lifting Points of Large-Span Steel Structure Based on Evolutionary Programming

Xin Wang, Xu Lei, Xuyang Cao, Yang Zhou, and Shunde Gao

School of Mechanical Engineering, Dalian University of Technology, 116024 Dalian, China
{wangxbd21, leixu_just_do_it}@163.com, Saner@126.com,
{huijunishixian, gaoshunde}@163.com

Abstract. To design the lifting points of large-span steel structure when the various compatibility equations are undefined in the lifting process, the programs based on improved evolutionary programming are developed by MATLAB. Lifting points design is to determine the comprehensive optimal strategy on number and distribution of lifting points, among which the minimum strain energy theory is mentioned and the secondary development technology of ANSYS-APDL is used. The performance and efficiency of the algorithms in different optimization ideas (hierarchy optimization and synchronic optimization) and methods (the particle swarm optimization and evolutionary programming) are compared, the results indicate that the improved evolutionary programming method based on synchronic optimization idea is satisfactory and provides a new but more effective method.

Keywords: Large-span steel structure, Lifting points design, Synchronic optimization, Evolutionary programming, Single-point mutation.

1 Introduction

As an important measurement of the national civil construction level, the steel structure technology has been highly valued by all countries, and great progress has been made in recently twenty years. With the development of the technology, the large-span and high complexity steel structure is being built increasingly. As a result, early high-altitude bulk-way assembly technology and block installation method can not meet the requirements, while the integral hoisting construction method is becoming increasingly popular for its safety, short time limit and less expense. Scholars both at home and abroad have advanced some problems existing in the integral hoisting construction process. Nine key technical points, including optimizing numbers and distribution of lifting points, are indicated in [1].

There are two main methods to design the lifting points. One is to do theoretical calculation on simplified model. The structure is regarded as rigid body in [2], and the best location of the lifting points is determined by solving static force balance equations. The other is to do some analysis by finite element software, ANSYS etc. The finite element model is built and analyzed, consequently satisfactory number and distribution of lifting points can be obtained. A method introduced in [3] indicates

that the strain energy in all different combinations of lifting points are calculated, and then the optimal combination which makes structure achieve the minimum strain energy can be found. In these two methods, the former is convenient to implement, but the structure is simplified so much that it's far from the actual one and the results have almost no significance to construction; while in the latter, the characteristics of the structure can be fully considered and the results are more reliable.

As to the principle, according to which we determine the comprehensive optimal strategy on number and distribution of lifting points, it's not hard to find that non-structural factors such as the performance of construction equipment and the site cooperation are mostly taken into account while the force characteristics of the structure itself are not, through consulting the related literatures [1].

For a structure, there are lots of strategies for locating lifting points if the number is uniquely given, but the optimal strategy can be obtained after optimal calculation. If the position vector is determined, the number of lifting points is obviously made sure. The method called hiberarchy optimization program determines the number first and then achieves the optimal strategy. While another method called synchronic optimization program takes the position vector as optimization object.

This paper firstly analyzes the problems in designing the lifting points of large-span steel structure based on the force characteristics of the structure in Section 1, then presents a mathematical model in Section 2. An improved evolutionary programming is introduced in Section 3. In Section 4, an example is given; programs based on improved evolutionary programming are developed by MATLAB, which are guided by hiberarchy optimization idea and synchronic optimization idea respectively; the secondary development technology of ANSYS-APDL is used. The results show that the improved evolutionary programming is efficient and practical in designing the lifting points of large-span steel structure.

2 Mathematical Model

There are many appearing in large-span steel structure construction process [4]; the stability of the structure and/or the components is the most basic problem. To meet the demand of stability, several lifting points should be set. Three different aspects about the lifting points design are presented in [5]: determining the number, determining the distribution and considering the lifting capacity of crane. Taking into account the stress concentration and load distribution problems appearing in the analysis, the fitness function based on penalty function method is defined in [5].

The mathematical model is defined as follows:

$$\left\{ \begin{array}{l} \min \quad Energy(n, X) \\ s.t. \quad F_i(n, X) - [F_i] \leq 0 \\ \quad \quad n \in \{n_{\min}, n_{\min} + 1, \dots, n_{\max} - 1, n_{\max}\} \\ \quad \quad i = 1, 2, \dots, n \end{array} \right. \quad (1)$$

Where n is the number of lifting points; X is the position vector; $Energy(n, X)$ is the strain energy; $F_i(n, X)$ and $[F_i]$ are, respectively, actual load and allowable load of the i th point; n_{min} and n_{max} are, respectively, allowable minimum value and maximum value of the number of lifting points.

The fitness function based on penalty function method is given as follows:

$$E_val(n, X) = Energy(n, X) + Q * \sum_{i=1}^n \max \{0, F_i(n, X) - [F_i]\} \tag{2}$$

Where $E_val(n, X)$ is the fitness value when the number is n and the position vector is X ; Q is the penalty factor.

3 An Improved Evolutionary Programming

Evolutionary programming (EP) [6-8] is an artificial intelligence technology of adaptation, which imitates natural organisms' evolutionary process and mechanism to solve problems. It was originally proposed to solve the discrete optimization problem by L.J.Fogel etc, in 1966. Then it was developed to optimize the continuous functions by D.B.Fogel in 1992. From then on, EP is widely used and development.

As a kind of group search algorithm, EP faces premature convergence problem. There are lots of methods being proposed to prevent such problem [9-12]. As we hope, an efficient algorithm has such feature: at the beginning, most individuals search the wide solution space to approach the nearby area of the global optimal solution, and in the end, most individuals just search the nearby area to find out the global optimal solution. A double population evolutionary algorithm is proposed in [10]. In this algorithm, the population is divided equally into two subgroups. This algorithm can balance global optimization and local optimization, but it arouses redundant calculation: local optimization is redundant at the beginning and global optimization is redundant in the end. In order to make full use of balance ability of the double population algorithm and avoid redundancy, the parameter α is introduced.

Suppose that there are N_{max} individuals in the population, which is divided into two subgroups, group one contains N_1 individuals and group two contains N_2 individuals. Individual in group one changes all its parameters while individual in group two just changes one of the parameters every time. We in turn call such update methods “overall mutation” and “single-point mutation”. And we have N_1, N_2 as follows:

$$N_1 = round(\alpha * N_{max}) \tag{3}$$

$$N_2 = N_{max} - N_1 \tag{4}$$

To ensure the efficiency of the algorithm, the parameter α changes its value from 1 to 0, and the numerical change is slow for most of the time and sharply declining near the end.

$$\alpha(k) = \frac{G_{\max}^2 - k^2}{G_{\max}^2 - 1} \quad (5)$$

Where G_{\max} is maximum iterations; k is current iteration.

Individual in group one takes combination sequence number of lifting points as variant, and the update rule is:

$$y^k(N_{\max} + i) = y^k(i) + \text{round}(N(0, \sigma_1)), i = 1, 2, \dots, N_1 \quad (6)$$

$$\sigma_1 = L / 6 \quad (7)$$

Where $y^k(i)$ and $y^k(N_{\max} + i)$ are, respectively, the i th combination sequence number and the $(N_{\max} + i)$ th new combination sequence number; $N(0, \sigma_1)$ is normal distribution with mean value 0 and variance σ_1^2 ; L is the count of the current optional combinations and gradually decreases as iteration goes on.

Individual in group two just changes one parameter of the combination $X = (x_1, x_2, \dots, x_n)^T$, and the update rule is:

$$x_j^k(N_{\max} + N_1 + i) = x_j^k(i) + \text{round}(N(0, \sigma_2)), i = 1, 2, \dots, N_2 \quad (8)$$

Where j is a random number satisfied with $j \in \{1, 2, \dots, n\}$ and $\sigma_2 = 1.3$ in this paper.

In the EP, we get N_{\max} new individuals from N_{\max} old ones through mutation operation, then pick out N_{\max} individuals from these $2N_{\max}$ ones according to the fitness values. Run these processes again and again until we get the satisfactory results.

The flow chart shown in Fig. 1 is to show us that how to apply the improved EP to the optimization of lifting points of large-span steel structure.

4 Example

Fig.2 shows the offshore platform deck mentioned in [5]. This structure is about 41 meters long, 11 meters wide, and weighs about 91 tons. It consists of three rows, seven columns main girders, so there are 21 optional lifting points. Number these points as shown in Fig.2.

The structure is analyzed by a variety of methods to find out the best location of five lifting points in [5], and the conclusion says that PSO has the best optimization feature within the same time. In addition, an improved PSO guided by hiberarchy optimization idea is also applied to optimize the lifting points. In this paper, the improved EP guided respectively by hiberarchy and synchronic optimization idea are applied, and the results are compared with that gotten by improved PSO.

Before the optimization, we should deal with the structure to get super-element through substructure technique [13]. To do so, the optional lifting points can be

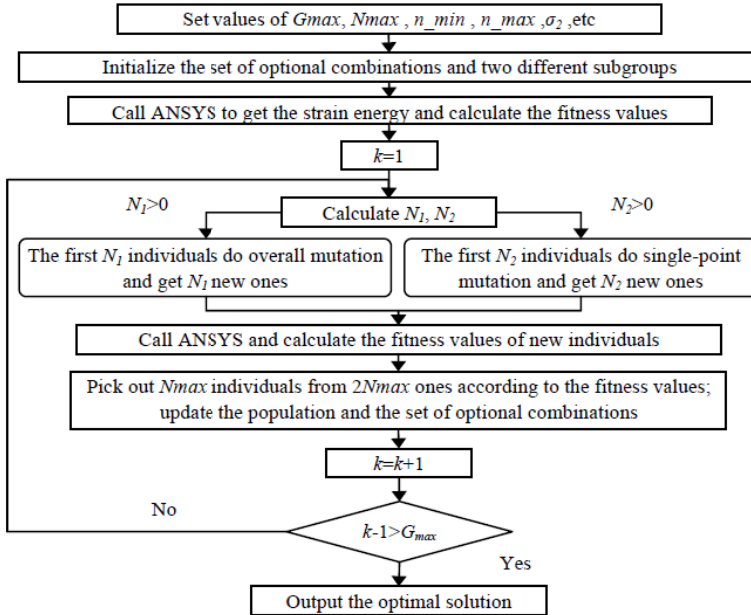


Fig. 1. Flow chart of lifting points optimization based on improved EP

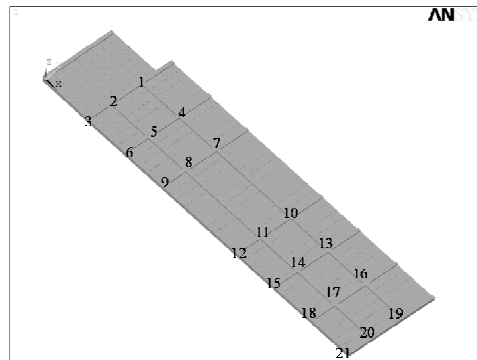


Fig. 2. The finite model of the offshore platform deck

treated as master DOF nodes. Using super-element to build finite element model can shorten the modeling time. And the structure calculation time can also be shortened because the internal DOFs of the structure are compressed.

In order to improve the data processing speed, all the optimization is operated in MATLAB. During the process, APDL [14] is used to compile the program files; and ANSYS is called in background to build and analyze the model, and output the results.

The flow chart of lifting points optimization design is presented as following:

- (1) ANSYS is called in background by MATLAB to build the model;
- (2) Initialize the lifting points information;
- (3) Call ANSYS, apply the load, solve it, get the strain energy and constraint reaction, and output the results to text files;
- (4) Load the data gotten in step (3), and calculate the fitness values;
- (5) Judge whether the results meet the conditions for the termination, if not, then turn step (6), otherwise, turn step (7);
- (6) Update the lifting points information, repeat step (3) through (5) ;
- (7) Stop the optimization and output the results.

The process of optimization guided by hiberarchy optimization idea contains two stages: to quantify the lifting points in the 1st stage; and to optimize the location in the 2nd stage. Same algorithm with different parameter values is applied in both stages to ensure the performance of EP and PSO in this paper.

In the improved PSO [15-17], linearly varying inertia weight is used:

$$w(k) = w_{\min} + \frac{w_{\max} - w_{\min}}{G_{\max}} (G_{\max} - k) \tag{9}$$

Where $w_{\min} = 0.4, w_{\max} = 0.9$; the cognitive and social behavioral factors $c_1 = c_2 = 2.05$; in the 1st stage $G_{\max} = 5$ and $N_{\max} = 60$; in the 2nd stage $G_{\max} = 25$ and $N_{\max} = 40$.

In the improved EP, we have the major parameters as Section 3 does; in the 1st stage $G_{\max} = 5$ and $N_{\max} = 60$; in the 2nd stage $G_{\max} = 25$ and $N_{\max} = 40$.

To do a comparative analysis, the fitness values of three, four and five lifting points are calculated, and the results are given in Table 1.

Table 1. Results of enumeration method

Number of lifting points	3	4	5
Count of combinations	$C_{21}^3 = 1330$	$C_{21}^4 = 5985$	$C_{21}^5 = 20349$
Optimal fitness value	4.28E4	2.12E4	1.19E4
Time-consuming/min	9.0	63.3	227.3

As can be seen from Table 1, a small increase in number of lifting points could sharply increase the count of combinations and the time-consuming. For comparative analysis, we set the number of lifting points between 3 and 5 in this paper.

The results of improved EP and improved PSO guided by hiberarchy optimization idea are shown in Table 2, and the curves of fitness value variation are shown in Fig.3.

Table 2. Results of hiberarchy optimization

Item	Time-consuming in the 1st stage /min	Time-consuming in the 2nd stage /min	Optimal number of lifting points	Optimal fitness value
PSO	2.1	9.1	5	1.55E4
EP	2.4	10.0	5	1.19E4

The algorithm guided by synchronic optimization idea takes the position vector as variant and, resultantly, the optimization of number is hidden.

The results of improved EP and improved PSO guided by synchronic optimization idea are shown in Table 3, and the curves of fitness value variation are shown in Fig.4.

Table 3. Results of synchronic optimization

Item	Time-consuming/min	Optimal number of lifting points	Optimal fitness value
PSO	7.6	5	2.39E4
EP	8.1	5	1.19E4

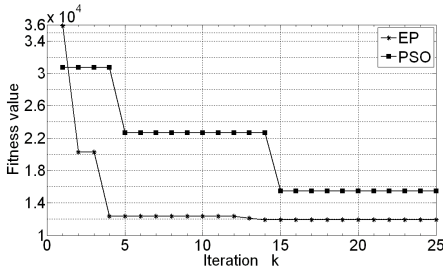


Fig. 3. Results of hiberarchy optimization

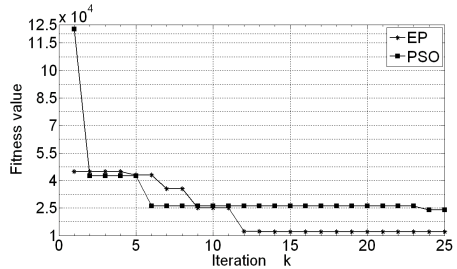


Fig. 4. Results of synchronic optimization

5 Conclusion

As can be seen from Table 2, in the hiberarchy optimization, the time-consuming of PSO in the 1st stage and the 2nd stage are almost the same as that of EP, respectively; both PSO and EP achieves the optimal number of lifting points, but EP obtains the global optimal solution while PSO doesn't. Table 3 shows that, in the synchronic optimization, the time-consuming of PSO is almost the same as that of EP; both PSO and EP achieves the optimal number of lifting points, but EP obtains the global optimal solution while PSO doesn't.

From the above comparative analysis, we can conclude that:

Comparing with the improved PSO, the improved EP has better performance in both the hiberarchy optimization and the synchronic optimization;

The improved EP is more efficient in the synchronic optimization than in the hiberarchy optimization.

In conclusion, The improved EP guided by the synchronic optimization idea is an efficient method in the optimization of lifting points for large-span steel structure.

References

1. Rui, Y.: Research on the Integral Hoisting Construction Method and the Controls Parameters. Ph.D. Dissertation of Tongji University, China (2008)
2. Rajasekaran, S., Annet, D., Sang Choo, Y.: Optimal Locations for Heavy Lifts for Offshore Platforms. *Asian Journal of Civil Engineering (Building and Housing)* 9(6), 605–627 (2008)
3. Bing, Z.: Research on Layout of Hoisting and Temporary Support in the Construction Process of Long-span Space Structures. Master Dissertation of Zhejiang University, China (2006)
4. Guo, Y.-L., Cui, X.-Q.: Key Technical Problems and Discussion in Construction Process of Large Span Steel Structures. *Industrial Construction* 34(12), 1–5 (2004)
5. Chen, B.-W.: The optimal research for the lift points of large span steel structure. Master Thesis, Dalian University of Technology (2011)
6. Fogel, D.B.: An Introduction to Simulated Evolutionary Optimization. *IEEE Trans. on Neural Networks* 5(1), 3–14 (1994)
7. Yun, Q.-X.: Evolutionary Algorithm. Metallurgical Industry Press, Beijing (2000)
8. Yan, X.: Actuality and Developmental Trend of the Evolutionary Programming. *Journal of Heze Teachers College* 25(4), 23–26 (2003)
9. Lin, D., Li, M.-Q., Kou, J.-S.: Two Methods to Prevent and Overcome Premature Convergence in Evolutionary Programming. *Journal of Systems Engineering* 16(3), 211–216 (2001)
10. Wang, X.-J., Xiang, D., Jiang, T., Lin, C.-S., Gong, S.-G., Fang, X.: A Novel Bi-Group Evolutionary Programming. *Chinese Journal of Computers* 29(5), 835–840 (2006)
11. Chellapilla, K.: Combining Mutation Operators in Evolutionary Programming. *IEEE Transactions on Evolutionary Computation* 2(3), 91–96 (1998)
12. Ji, M., Tang, H., Guo, J.: A Single-point Mutation Evolutionary Programming. *Information Processing Letters* (90), 293–299 (2004)
13. Jiang, S.X., et al.: Advanced ANSYS Finite Element Analysis Method and Application Examples. China Water Power Press, Beijing (2006)
14. Boyi Team: APDL Parametric Finite Element Analysis Technology and Application Examples. China Water Power Press, Beijing (2004)
15. Bo, L.: Particle Swarm Optimization Algorithm and its Engineering Application. Publishing House of Electronics Industry, Beijing (2010)
16. Poli, R., Kennedy, J., Blackwell, T.: Particle Swarm Optimization-An Overview. *Swarm Intell.* 1, 33–57 (2007)
17. Fang, G.: Research on Intelligent Particle Swarm Optimization Algorithm. Ph.D. Dissertation of Harbin Institute of Technology (2008)

Modifying Feasible SQP Method for Inequality Constrained Optimization

Zhijun Luo^{1,*}, Zhibin Zhu², and Guohua Chen³

¹ Department of Mathematics & Applied Mathematics, Hunan University of Humanities, Science and Technology, Loudi, 417000, P.R. China

² School of Mathematics and Computing Science, Guilin University of Electronic Technology, Guilin, 541004, P.R. China
ld1zj123@163.com

Abstract. This paper is concerned with an improved feasible sequential quadratic programming (FSQP) method which solves an inequality constrained nonlinear optimization problem. As compared with the existing SQP methods, at each iteration of our method, the base direction is only necessary to solve a equality constrained quadratic programming, the feasible direction and the high-order revised direction which avoids Maratos effect are obtained by explicit formulas. Furthermore, the global and superlinear convergence are proved under some suitable conditions.

Keywords: Nonlinear optimization, FSQP method, Global convergence, Superlinear convergence rate.

1 Introduction

Consider the nonlinear inequality constrained optimization problem:

$$\begin{aligned} \min f(x) \\ \text{s.t. } g_j(x) \leq 0, j \in I = \{1, 2, \dots, m\}, \end{aligned} \quad (1)$$

where $f, g_j : R^n \rightarrow R (j \in I)$ are continuously differentiable functions. Denote the feasible set for (1) by $X = \{x \in R^n \mid g_j(x) \leq 0, j \in I\}$.

SQP(Sequential Quadratic Programming[1]) method is one of the most effective methods for solving nonlinear programming. It generates iteratively the main search direction d_0 by solving the following quadratic programming(QP) subproblem:

$$\begin{aligned} \min \nabla f(x)^T d + \frac{1}{2} d^T H d \\ \text{s.t. } g_j(x) + \nabla g_j(x)^T d \leq 0, j = 1, 2, \dots, m, \end{aligned} \quad (2)$$

where $H \in R^{n \times n}$ is a symmetric positive definite matrix. However, such type SQP algorithms have two serious shortcomings: (1) SQP algorithms require that

* This work was supported in part by the National Natural Science Foundation (11061011) of China, and the Educational Reform Research Fund of Hunan University of Humanities, Science and Technology (NO.RKJGY1030).

the relate QP subproblem (2) must be consistency. (2) There exists Maratos effect. Many efforts have been made to overcome the shortcomings through modifying the quadratic subproblem (2) and the direction d [2]-[9]. Some algorithms solve the problem (1) by using the idea of filter method or trust-region [10]-[12].

For the problem (2), P. Spellucci [5] proposed a new method, the d_0 is obtained by solving QP subproblem with only equality constraints:

$$\begin{aligned} \min \quad & \nabla f(x)^T d + \frac{1}{2} d^T H d \\ \text{s.t.} \quad & g_j(x) + \nabla g_j(x)^T d = 0, j \in I. \end{aligned} \quad (3)$$

If $d_0 = 0$ and $\lambda \geq 0$ (λ is said to be the corresponding KKT multiplier vector.), the algorithm stops. The most advantage of these algorithms is merely necessary to solve QP subproblems with only equality constraints. However, if $d_0 = 0$, but $\lambda < 0$, the algorithm will not implement successfully. Recently, Z.B.Zhu [8] Consider the following QP subproblem:

$$\begin{aligned} \min \quad & \nabla f(x)^T d + \frac{1}{2} d^T H d \\ \text{s.t.} \quad & p_j(x) + \nabla g_j(x)^T d = 0, j \in L, \end{aligned} \quad (4)$$

where p_j is suitable vector, which guarantees to hold that if $d_0 = 0$, then x is a KKT point of (1), i.e. if $d_0 = 0$, then it holds that $\lambda_0 \geq 0$. Depended strictly on the strict complementarity, which is rather strong and difficult for testing, the superlinear convergence properties of the SQP algorithm is obtained. For avoiding the superlinear convergence depend strictly on the strict complementarity, another some SQP algorithms (see [13]) have been proposed, however it is regretful that these algorithms are infeasible SQP type and nonmonotone. In [14], a feasible SQP algorithm is proposed. Using generalized projection technique, the superlinear convergence properties are still obtained under weaker conditions without the strict complementarity.

We will develop an improved feasible SQP method for solving optimization problems based on the one in [8]. The traditional FSQP algorithms, in order to prevent iterates from leaving the feasible set, and avoid Maratos effect, it needs to solve two or three QP subproblems like (2). In our algorithm, per single iteration, it is only necessary to solve an equality constrained quadratic programming, which is very similar to (4). Obviously, it is simpler to solve the equality constrained QP problem than to solve the QP problem with inequality constraints. In order to void the Maratos effect, combined the generalized projection technique, a height-order correction direction is computed by an explicit formula, and it plays a important role in avoiding the strict complementarity. Furthermore, its global and superlinear convergence rate are obtained under some suitable conditions.

This paper is organized as follows: In Section 2, we state the algorithm; The well-defined of our approach is also discussed, the accountability of which allows us to present global convergence guarantees under common conditions in Section 3, while in Section 4 we deal with superlinear convergence.

2 Description of Algorithm

The active constraints set of (II) is denoted as follows:

$$I(x) = \{j \in I \mid g_j(x) = 0\}, I = \{1, 2, \dots, m\}. \tag{5}$$

Throughout this paper, following basic assumptions are assumed.

H 2.1. *The feasible set $X \neq \phi$, and functions $f, g_j (j \in I)$ are twice continuously differentiable.*

H 2.2. $\forall x \in X$, the vectors $\{\nabla g_j(x), j \in I(x)\}$ are linearly independent.

Firstly, for a given point $x^k \in X$, by using the pivoting operation, we obtain an approximate active $J_k = J(x^k)$, such that $I(x^k) \subseteq J_k \subseteq I$.

Sub-algorithm A:

Step 1 For the current point $x^k \in X$, set $i = 0, \epsilon_i(x^k) = \epsilon_0 \in (0, 1)$.

Step 2 If $\det(A_i(x^k)^T A_i(x^k)) \geq \epsilon_i(x^k)$, let $J_k = J_i(x^k), A_k = A_i(x^k), i(x^k) = i$, STOP. Otherwise goto Step 3, where

$$J_i(x^k) = \{j \in I \mid -\epsilon_i(x^k) \leq g_j(x^k) \leq 0\}, A_i(x^k) = (\nabla g_j(x^k), j \in J_i(x^k)). \tag{6}$$

Step 3 Let $i = i + 1, \epsilon_i(x^k) = \frac{1}{2}\epsilon_{i-1}(x^k)$, and goto Step 2.

Theorem 2.1. *For any iteration, there is no infinite cycle for above subalgorithm A. Moreover, if $\{x^k\}_{k \in K} \rightarrow x^*$, then there exists a constant $\bar{\epsilon} > 0$, such that $\epsilon_{k,i_k} \geq \bar{\epsilon}$, for $k \in K, k$ large enough.*

Now, the algorithm for the solution of the problem (II) can be stated as follows.

Algorithm A:

Step 0 Initialization:

Given a starting point $x^0 \in X$, and an initial symmetric positive definite matrix $H_0 \in R^{n \times n}$. Choose parameters $\epsilon_0 \in (0, 1), \alpha \in (0, \frac{1}{2}), \tau \in (2, 3)$. Set $k = 0$;

Step 1 For x^k , compute $J_k = J(x^k), A_k = A(x^k)$ by using Sub-algorithm A.

Step 2 Computation of the vector d_0^k :

2.1

$$B_k = (A_k^T A_k)^{-1} A_k^T, v^k = (v_j^k, j \in J_k) = -B_k \nabla f(x^k),$$

$$p_j^k = \begin{cases} -v_j^k, & v_j^k < 0, \\ g_j(x^k), & v_j^k \geq 0. \end{cases} \quad p^k = (p_j^k, j \in J_k). \tag{7}$$

2.2 Solve the following equality constrained QP subproblem at x^k :

$$\begin{aligned} \min & \nabla f(x^k)^T d + \frac{1}{2} d^T H_k d \\ \text{s.t.} & p_j^k + \nabla g_j(x^k)^T d = 0, j \in J_k. \end{aligned} \tag{8}$$

Let d_0^k be the KKT point of (8), and $b^k = (b_j^k, j \in J_k)$ be the corresponding multiplier vector. If $d_0^k = 0$, STOP. Otherwise, CONTINUE;

Step 3 Computation of the feasible direction with descent d^k :

$$d^k = d_0^k - \delta_k A_k (A_k^T A_k)^{-1} e_k. \tag{9}$$

Where $e_k = (1, \dots, 1)^T \in R^{|J_k|}$, and

$$\delta_k = \frac{\|d_0^k\| (d_0^k)^T H_k d_0^k}{2|(\mu^k)^T e_k| \cdot \|d_0^k\| + 1}, \quad \mu^k = -(A_k^T A_k)^{-1} A_k^T \nabla f(x^k). \tag{10}$$

Step 4 Computation of the high-order revised direction \tilde{d}^k :

$$\tilde{d}^k = -\delta_k A_k (A_k^T A_k)^{-1} (\|d_0^k\|^\tau e_k + \tilde{g}_{J_k}(x^k + d^k)), \tag{11}$$

where

$$\tilde{g}_{J_k}(x^k + d^k) = g_{J_k}(x^k + d^k) - g_{J_k}(x^k) - \nabla g_{J_k}(x^k)^T d^k. \tag{12}$$

Step 5 The line search:

Compute t_k , the first number t in the sequence $\{1, \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots\}$ satisfying

$$f(x^k + td^k + t^2 \tilde{d}^k) \leq f(x^k) + \alpha t \nabla f(x^k)^T d^k, \tag{13}$$

$$g_j(x^k + td^k + t^2 \tilde{d}^k) \leq 0, \quad j \in I. \tag{14}$$

Step 6 Update:

Obtain H_{k+1} by updating the positive definite matrix H_k using some quasi-Newton formulas. Set $x^{k+1} = x^k + t_k d^k + t^2 \tilde{d}^k$, and $k = k + 1$. Go back to step 1.

Lemma 2.2. For the QP subproblem (8) at x^k , if $d_0^k = 0$, then x^k is a KKT point of (1). If $d_0^k \neq 0$, then d^k computed in step 4 is a feasible direction with descent of (1) at x^k .

3 Global Convergence of Algorithm

In this section, firstly, it is shown that Algorithm A given in section 2 is well-defined, that is to say, for every k, that the line search at Step 5 is always successful.

Lemma 3.1. The line search in step 5 yields a stepsize $t_k = (\frac{1}{2})^i$ for some finite $i = i(k)$.

Proof. It is a well-known result according to Lemma 2.2. For (13),

$$\begin{aligned} s &\triangleq f(x^k + td^k + t^2 \tilde{d}^k) - f(x^k) - \alpha t \nabla f(x^k)^T d^k \\ &= \nabla f(x^k)^T (td^k + t^2 \tilde{d}^k) + o(t) - \alpha t \nabla f(x^k)^T d^k \\ &= (1 - \alpha)t \nabla f(x^k)^T d^k + o(t). \end{aligned}$$

For (14), if $j \notin I(x^k)$, $g_j(x^k) < 0$; $j \in I(x^k)$, $g_j(x^k) = 0$, $\nabla g_j(x^k)^T d^k < 0$, so we have

$$g_j(x^k + td^k + t^2 \tilde{d}^k) = \nabla f(x^k)^T (td^k + t^2 \tilde{d}^k) + o(t) = \alpha t \nabla g_j(x^k)^T d^k + o(t).$$

In the sequel, the global convergence of Algorithm A is shown. For this reason, we make the following additional assumption.

H 3.1. $\{x^k\}$ is bounded, which is the sequence generated by the algorithm, and there exist constants $b \geq a > 0$, such that $a\|y\|^2 \leq y^T H_k y \leq b\|y\|^2$, for all k and all $y \in R^n$.

Since there are only finitely many choices for sets $J_k \subseteq I$, and the sequence $\{d_0^k, d_1^k, \tilde{d}^k, v^k, b^k\}$ is bounded, we can assume without loss of generality that there exists a subsequence K , such that

$$x^k \rightarrow x^*, H_k \rightarrow H_*, d_0^k \rightarrow d_0^*, d^k \rightarrow d^*, \tilde{d}^k \rightarrow \tilde{d}^*, b^k \rightarrow b^*, v^k \rightarrow v^*, J_k \equiv J \neq \emptyset, k \in K, \tag{15}$$

where J is a constant set.

Theorem 3.2. *The algorithm either stops at the KKT point x^k of the problem (1) in finite number of steps, or generates an infinite sequence $\{x^k\}$ any accumulation point x^* of which is a KKT point of the problem (1).*

Proof. The first statement is easy to show, since the only stopping point is in step 3. Thus, assume that the algorithm generates an infinite sequence $\{x^k\}$, and (15) holds. According to Lemma 2.2, it is only necessary to prove that $d_0^* = 0$. Suppose by contradiction that $d_0^* \neq 0$. Then, from Lemma 2.2, it is obvious that d^* is well-defined, and it holds that

$$\nabla f(x^*)^T d^* < 0, \nabla g_j(x^*)^T d^* < 0, j \in I(x^*) \subseteq J. \tag{16}$$

Thus, from (16), it is easy to see that the step-size t_k obtained in step 5 are bounded away from zero on K , i.e.,

$$t_k \geq t_* = \inf\{t_k, k \in K\} > 0, k \in K. \tag{17}$$

In addition, from (13) and Lemma 2.2, it is obvious that $\{f(x^k)\}$ is monotonous decreasing. So, according to assumption H 2.1, the fact that $\{x^k\}_K \rightarrow x^*$ implies that

$$f(x^k) \rightarrow f(x^*), k \rightarrow \infty. \tag{18}$$

So, from (13), (16), (17), it holds that

$$0 = \lim_{k \in K} (f(x^{k+1}) - f(x^k)) \leq \lim_{k \in K} (\alpha t_k \nabla f(x^k)^T d^k) \leq \frac{1}{2} \alpha t_* f(x^*)^T d^* < 0, \tag{19}$$

which is a contradiction thus $\lim_{k \rightarrow \infty} d_0^k = 0$. Thus, x^* is a KKT point of (1).

4 The Rate of Convergence

Now we discuss the convergent rate of the algorithm, and prove that the sequence $\{x^k\}$ generated by the algorithm is one-step superlinearly convergent under some mild conditions without the strict complementarity. For this purpose, we add some regularity hypothesis.

H 4.1. *The sequence $\{x^k\}$ generated by Algorithm A is bounded, and possess an accumulation point x^* , such that the KKT pair (x^*, u^*) satisfies the strong second-order sufficiency conditions, i.e.,*

$$d^T \nabla_{xx}^2 L(x^*, u^*) d > 0, \forall d \in \Omega \triangleq \{d \in R^n : d \neq 0, \nabla g_{I^+}(x^*)^T d = 0\},$$

where, $L(x, u) = f(x) + \sum_{j \in I} u_j g_j(x)$, $I^+ = \{j \in I : u_j^* > 0\}$.

Lemma 4.1. *Let $H2.1 \sim H4.1$ holds, $\lim_{k \rightarrow \infty} \|x^{k+1} - x^k\| = 0$. Thereby, the entire sequence $\{x^k\}$ converges to x^* , i.e. $x^k \rightarrow x^*, k \rightarrow \infty$.*

Proof. From the Lemma 4.2, it is easy to see that

$$\lim_{k \rightarrow \infty} \|x^{k+1} - x^k\| = \lim_{k \rightarrow \infty} (\|t_k d^k + t_k^2 \tilde{d}^k\|) \leq \lim_{k \rightarrow \infty} (\|d^k\| + \|\tilde{d}^k\|) = 0.$$

Moreover, together with Theorem 1.1.5 in [4], it shows that $x^k \rightarrow x^*, k \rightarrow \infty$.

Lemma 4.2. *Suppose that assumptions $H 2.1-H 3.1$ hold, then,*

- 1) *There exists a constant $\zeta > 0$, such that $\|(A_k^T A_k)^{-1}\| \leq \zeta$;*
- 2) *$\lim_{k \rightarrow \infty} d_0^k = 0, \lim_{k \rightarrow \infty} d^k = 0, \lim_{k \rightarrow \infty} \tilde{d}^k = 0$;*
- 3) *$\|d^k\| \sim \|d_0^k\|, \|\tilde{d}^k\| = O(\|d^k\|^2)$, .*

Proof. 1) From Sub-algorithm A and Theorem 2.1, we have the following result.

$$\det(A_*^T A_*) = \lim_{k \in K} \det(A_k^T A_k) \geq \lim_{k \in K} \varepsilon_k \geq \bar{\varepsilon} > 0.$$

Thereby, the first conclusion 1) follows.

2) Prove $\lim_{k \rightarrow \infty} d_0^k = 0$.

According to the proof of Theorem 3.2, the fact that $x^k \rightarrow x^*(k \rightarrow \infty)$ implies that it is true.

The proof of $\lim_{k \rightarrow \infty} d^k = 0, \lim_{k \rightarrow \infty} \tilde{d}^k = 0$ are elementary from the result of 1) as well as formulas (9) and (11).

3) The proof of 3) is elementary from the formulas (9), (11) and assumption H 2.1

Lemma 4.3. *It holds, for k large enough, that*

- 1) $J_k \equiv I(x^*) \triangleq I_*, b^k \rightarrow u_{I_*} = (u_j^*, j \in I_*), v^k \rightarrow (u_j^*, j \in I_*)$.
- 2) $I^+ \subseteq L_k = \{j \in J_k : g_j(x^k) + \nabla g_j(x^k)^T d_0^k = 0\} \subseteq I(x^*)$.

Proof. 1) Prove $J_k \equiv I_*$.

On one hand, from Lemma 2.1, we know, for k large enough, that $I_* \subseteq J_k$. On the other hand, if it doesn't hold that $J_k \subseteq I_*$, then there exist constants j_0 and $\beta > 0$, such that

$$g_{j_0}(x^*) \leq -\beta < 0, j_0 \in J_k.$$

So, according to $d_0^k \rightarrow 0$ and assumption H2, it holds, for k large enough, that

$$p_{j_0}(x^k) + \nabla g_{j_0}(x^*)^T d_0^k = \begin{cases} -v_{j_0}^k + \nabla g_{j_0}(x^*)^T d_0^k \geq -\frac{1}{2}v_{j_0}^k > 0, b_{j_0}^k < 0, \\ g_{j_0}(x^k) + \nabla g_{j_0}(x^*)^T d_0^k \leq -\frac{1}{2}\beta < 0, v_{j_0}^k \geq 0. \end{cases}, \quad (20)$$

which is contradictory with (8) and the fact $j_0 \in J_k$. So, $J_k \equiv I_*$ (for k large enough).

Prove that $b^k \rightarrow u_{I_*} = (u_j^*, j \in I_*)$, $v^k \rightarrow (u_j^*, j \in I_*)$.

For the $v^k \rightarrow (u_j^*, j \in I_*)$ statement, we have the following results from the definition of v^k ,

$$v^k \rightarrow -B_* \nabla f(x^*) = -(A_*^T A_*)^{-1} A_*^T \nabla f(x^*)$$

In addition, since x^* is a KKT point of (II), it is evident that

$$\nabla f(x^*) + A_* u_{I_*} = 0, u_{I_*} = -B_* \nabla f(x^*). \quad (21)$$

i.e.

$$u_{I_*} = -(A_*^T A_*)^{-1} A_*^T \nabla f(x^*).$$

Otherwise, from (8), the fact that $d_0^k \rightarrow 0$ implies that

$$\nabla f(x^k) + H_k d_0^k + A_k b^k = 0, b^k \rightarrow -B_* \nabla f(x^*) = u_{I_*}.$$

The claim holds.

2) For $\lim_{k \rightarrow \infty} (x^k, d_0^k) = (x^*, 0)$, we have $L_k \subseteq I(x^*)$. Furthermore, it has $\lim_{k \rightarrow \infty} u_{I^+}^k = u_{I^+}^* > 0$, so the proof is finished.

In order to obtain superlinear convergence, a crucial requirement is that a unit step size is used in a neighborhood of the solution. This can be achieved if the following assumption is satisfied.

H 4.2. Let $\|(\nabla_{xx}^2 L(x^k, u_{J_k}^k) - H_k) d^k\| = o(\|d^k\|)$, where $L(x, u_{J_k}^k) = f(x) + \sum_{j \in J_k} u_j^k g_j(x)$.

According to Theorem 4.2 in [14], it is easy to obtain the following results.

Lemma 4.4. For k large enough, $t_k \equiv 1$.

Furthermore, In a way similar to the proof of Theorem 4.1 in [9], we may obtain the following theorem:

Theorem 4.5. Under all above-mentioned assumptions, the algorithm is super-linearly convergent, i.e., the sequence $\{x^k\}$ generated by the algorithm satisfies that

$$\|x^{k+1} - x^*\| = o(\|x^k - x^*\|).$$

References

1. Boggs, P.T., Tolle, J.W.: A Strategy for Global Convergence in a Sequential Quadratic Programming Algorithm. *SIAM J. Num. Anal.* 26, 600–623 (1989)
2. Han, S.P.: Superlinearly Convergent Variable Metric Algorithm for General Nonlinear Programming Problems. *Mathematical Programming* 11, 263–282 (1976)
3. Powell, M.J.D.: A Fast Algorithm for Nonlinearly Constrained Optimization Calculations. In: Waston, G.A. (ed.) *Numerical Analysis*, pp. 144–157. Springer, Berlin (1978)
4. Panier, E.R., Tits, A.L.: On Combining Feasibility, Descent and Superlinear Convergence in Inequality Constrained Optimization. *Mathematical Programming* 59, 261–276 (1993)
5. Spellucci, P.: An SQP Method for General Nonlinear Programs Using Only Equality Constrained Subproblems. *Mathematical Programming* 82, 413–448 (1998)
6. Lawrence, C.T., Tits, A.L.: A Computationally Efficient Feasible Sequential Quadratic Programming Algorithm. *SIAM J. Optim.* 11, 1092–1118 (2001)
7. Qi, L., Yang, Y.F.: Globally and Superlinearly Convergent QP-free Algorithm for Nonlinear Constrained Optimization. *JOTA* 113, 297–323 (2002)
8. Zhu, Z.B., Zhang, W.D., Geng, Z.J.: A feasible SQP method for nonlinear programming. *Applied Mathematics and Computation* 215, 3956–3969 (2010)
9. Luo, Z.J., Chen, G.H., Liang, J.L.: A variant of feasible descent SQP method for inequality constrained optimization. *International Journal of Pure and Applied Mathematics* 61(2), 161–168 (2010)
10. Gu, C., Zhu, D.T.: A non-monotone line search multidimensional filter-SQP method for general nonlinear programming. *Numerical Algorithms* 56(4), 537–559 (2011)
11. Hao, C.L., Liu, X.W.: A trust-region filter-SQP method for mathematical programs with linear complementarity constraints. *Journal of Industrial and Management Optimization (JIMO)* 7(4), 1041–1055 (2011)
12. Hao, C.L., Liu, X.W.: Global convergence of an SQP algorithm for nonlinear optimization with overdetermined constraints. *Numerical Algebra Control and Optimization* 2(1), 19–29 (2012)
13. Binnans, J.F., Launay, G.: Sequential quadratic programming with penalization the displacement. *SIAM J. Optimization* 54(4), 796–812 (1995)
14. Jian, J.B., Tang, C.M.: An SQP Feasible Descent Algorithm for Nonlinear Inequality Constrained Optimization Without Strict Complementarity. *An International Journal Computers and Mathematics with Application* 49, 223–238 (2005)

On the Solvable n -Lie Algebras

Liu Jianbo¹, Zhang Yanyan^{2,*}, Men Yafeng¹, and Chen Wenyong¹

¹ School of Mathematics and Statistics

² School of Computer Science and Telecommunication Engineering,
Northeastern University at Qinhuangdao, 066004, Qinhuangdao, China
zhangyy@mail.neuq.edu.cn

Abstract. The solvable n -Lie algebras are studied, we determined some properties of solvable n -Lie algebras, gave the definition of Borel n -subalgebra, and also got some results about Borel n -subalgebras.

Keywords: Lie algebra, n -Lie algebra, solvable.

1 Introduction

In the last few years, the theory of n -Lie algebras has attracted a lot of attention due to its close connection with the Nambu mechanics proposed by Y.Nambu in [1] as a generalization of the classical Hamiltonian mechanics. Historically, the first work dedicated to these topics was the paper of V. T. Filippov in [2]. In 1985, V. T. Filippov introduced the concept of an n -Lie algebra, which is a natural generalization of the concept of a Lie algebra to the case where the fundamental multiplication operation is n -ary ($n \geq 2$), When $n = 2$ the definition agrees with the usual definition of a Lie algebra [3][4]. In it, he considered n -ary multi-linear and skew symmetric operator $[x_1, \dots, x_n]$ satisfying:

$$[[x_1, \dots, x_n], y_2, \dots, y_n] = \sum_{k=1}^n [x_1, \dots, [x_k, y_2, \dots, y_n], \dots, x_n]$$

which was defined as the operator of an n -Lie algebra. There examples of n -Lie algebras and classification of $(n+1)$ -dimensional n -Lie algebras are given, structural notions such as simplicity, nilpotency are developed. In 1987, Sh. M. Kasymov, in [5], introduced and studied k -solvability, k -nilpotency ($1 \leq k \leq n$), Cartan subalgebras, Killing forms and representations of n -Lie algebras, he also proved an n -ary analog of the Engel's theorem. Recently, φ -free n -Lie algebra are studied in [6][7][8].

In this paper, we mainly discuss the solvable n -Lie algebras, give some properties of solvable n -Lie algebras, and also get some results about Borel n -Lie algebras.

* Corresponding author. This paper supported by "the Fundamental Research Funds for the Central Universities" China (N090323007).

2 Some Definitions on n -Lie Algebra

In this section, we give some definitions of n -Lie algebras, which will be used in next sections, more definition can be seen in [2] and [5].

Definition 1. An n -Lie algebra is a vector space V over a field F on which there is defined an n -ary multi-linear operation $[\cdot, \dots, \cdot]$ satisfying the identities:

$$[x_1, \dots, x_n] = (-1)^{\tau(\sigma)} [x_{\sigma(1)}, \dots, x_{\sigma(n)}]$$

and

$$[[x_1, \dots, x_n], y_2, \dots, y_n] = \sum_{k=1}^n [x_1, \dots, [x_k, y_2, \dots, y_n], \dots, x_n]$$

where σ runs over the symmetric group S_n and the number $\tau(\sigma)$ is equal to 0 or 1 depending on the parity of the permutation σ .

Definition 2. A subalgebra W of an n -Lie algebra V is a subspace of V , satisfying

$$[W, W, \dots, W] \subseteq W;$$

An ideal of an n -Lie algebra V is a subspace I of V , satisfying

$$[I, V, \dots, V] \subseteq I.$$

If I and J are ideals of V , then $I + J$ is an ideal of V ; If I and J are ideals of V , when $V = I + J$ and $I \cap J = \{0\}$, we denote $V = I \oplus J$.

Definition 3. If V is an n -Lie algebra, $[V, V, \dots, V]$ is called the derivation algebra of V , denoted by V^1 ; when $V^1 = \{0\}$, V is said to be an abelian n -Lie algebra.

Definition 4. An n -Lie algebra V is said to be simple, if $V^1 \neq 0$ and it has no ideals distinct from $\{0\}$ and V .

Theorem 1. ([9]) An $(n+1)$ -dimensional n -Lie algebra V is simple if and only if $\dim V = \dim V^1 = n + 1$.

Example 1. Let A be an $(n+1)$ -dimensional Euclidean space with an orthonormal basis e_1, \dots, e_{n+1} . We denote, by $[x_1, \dots, x_n]$, the vector product of the vectors $x_1, \dots, x_n \in A$, that is

$$[x_1, \dots, x_n] = \begin{vmatrix} x_{11} & x_{12} & \dots & x_{1n} & e_1 \\ x_{21} & x_{22} & \dots & x_{2n} & e_2 \\ \dots & \dots & \dots & \dots & \dots \\ x_{n+1,1} & x_{n+1,2} & \dots & x_{n+1,n} & e_{n+1} \end{vmatrix}$$

where $(x_{1i}, \dots, x_{n+1i})$ are the coordinates of the vectors $x_i, i = 1, \dots, n$, then $(A, [\cdot, \dots, \cdot])$ becomes an $(n+1)$ -dimensional simple n -Lie algebra, denoted by A_1 . It has the following multiplication table of the basis vectors:

$$[e_1, \dots, e_{i-1}, \hat{e}_i, e_{i+1}, \dots, e_{n+1}] = (-1)^{n+1+i} e_i,$$

where $i = 1, \dots, n + 1$ and the symbol \hat{e}_i means that e_i is omitted.

Definition 5. *The subspace*

$$Z(V) = \{x \in V \mid [x, V, \dots, V] = 0\}$$

is called the center of V . It is clear that $Z(V)$ is an ideal of V .

Definition 6. *An ideal I of V is said to be solvable if $I^{(r)} = 0$ for some $r \geq 0$, where $I^{(0)} = I$, and by induction, we define*

$$I^{(s+1)} = [I^{(s)}, \dots, I^{(s)}],$$

for $s \geq 0$. If V contains no nonzero solvable ideals, then V is called semi-simple. Simple algebras are semi-simple.

Theorem 2. ([2]) *If I and J are solvable ideals of V , then $I + J$ is a solvable ideal of V .*

The largest solvable ideal R is called the *radical* of V . It is clear that V is semi-simple if and only if $R = 0$.

Definition 7. *If I is an ideal of an n -Lie algebra V , the quotient algebra is the quotient space V/I defined by an n -ary multiplication*

$$[x_1 + I, \dots, x_n + I] = [x_1, \dots, x_n] + I,$$

for $x_i \in V, i = 1, \dots, n$. We still denote by V/I .

Definition 8. *If V is an n -Lie algebra ($n \geq 3$), a reduced $(n-1)$ -Lie algebra V_a from V is an $(n-1)$ -Lie algebra for any $a \in V$, with multiplication defined by*

$$[x_1, \dots, x_{n-1}]_a = [a, x_1, \dots, x_{n-1}].$$

Theorem 3. ([10]) *Let V be a 3-Lie algebra, $\{e_1, \dots, e_m\}$ be a basis of V . If V_i is the 2-reduced Lie algebra (briefly denoted reduced Lie algebra) V_{e_i} and $[\cdot, \cdot]_i$ is the corresponding multiplication of V_i , then every reduced Lie algebra from V is a Lie algebra and can be defined by means of the reduced Lie algebras V_i , with*

$$[x, y]_a = \sum_{j=1}^m \alpha_j [x, y]_j,$$

for every $a = \sum_{j=1}^m \alpha_j e_j \in V$.

3 Conclusion

In this section, we firstly give some properties of solvable n -Lie algebras and three equivalent conditions of the solvability of n -Lie algebras; Subsequently, we give the concept of a Borel n -subalgebra, and two propositions relative to it.

we firstly give three Lemmas, which can be seen in [2].

Lemma 1. *Let V be an n -Lie algebra, then $V^{(0)}, V^{(1)}, \dots, V^{(k)}$ are ideals of V .*

Lemma 2. *If V is a solvable n -Lie algebra, then each subalgebra and each homomorphic image of V is solvable.*

Lemma 3. *If an n -Lie algebra V contains a solvable ideal I such that the quotient algebra V/I is solvable, then the n -Lie algebra V is also solvable. This implies that the sum of two solvable ideals is again a solvable ideal.*

Theorem 4. *If V is a finite dimensional n -Lie algebra over an algebraically closed field F of characteristic zero, then the following conditions are equivalent:*

- (1) V is a solvable n -Lie algebra;
- (2) there exists a series of ideals of V :

$$V = V_0 \supseteq V_1 \supseteq \dots \supseteq V_r = \{0\},$$

such that V_i/V_{i+1} are abelian n -Lie algebras;

- (3) there exists a series of subalgebras of V :

$$V = H_0 \supseteq H_1 \supseteq \dots \supseteq H_k = \{0\},$$

such that H_{i+1} is an ideal of H_i and H_i/H_{i+1} are abelian n -Lie algebras;

- (4) there exists a series of subalgebras of V :

$$V = B_0 \supseteq B_1 \supseteq \dots \supseteq B_m = \{0\},$$

such that B_{i+1} is an ideal of B_i and $\dim(B_i/B_{i+1}) = 1$.

Proof. (1) \Rightarrow (2) Since V is a solvable n -Lie algebra, there is a positive integer number r , such that $V^{(r)} = 0$. Let $V_i = V^{(i)}$, for $i = 0, 1, \dots, r$. From Lemma 1, V_0, V_1, \dots, V_r are ideals of V , and

$$V_{i+1} = V^{(i+1)} = [V^{(i)}, \dots, V^{(i)}] \subseteq [V^{(i)}, V, \dots, V] \subseteq V^{(i)} = V_i,$$

for $i = 0, 1, \dots, r - 1$. Because of $V_{i+1} = [V_i, \dots, V_i]$, we have

$$[V_i/V_{i+1}, \dots, V_i/V_{i+1}] = [V_i, \dots, V_i] + V_{i+1} = \{0\},$$

that is, V_i/V_{i+1} are abelian n -Lie algebras.

(2) \Rightarrow (3) By (2), V_i are ideals of V , it is followed that V_i are subalgebras of V , let $H_i = A_i$, for $i = 0, 1, \dots, r$, (3) holds.

- (3) \Rightarrow (4) If there is a series of subalgebras of V :

$$V = H_0 \supseteq H_1 \supseteq \dots \supseteq H_k = \{0\},$$

H_{i+1} is an ideal of H_i , i.e.

$$[H_{i+1}, H_i, \dots, H_i] \subseteq H_{i+1},$$

and H_i/H_{i+1} are abelian n -Lie algebras, i.e.

$$[H_i/H_{i+1}, \dots, H_i/H_{i+1}] = \{0\}$$

or

$$[H_i, \dots, H_i] \subseteq H_{i+1}.$$

If $\dim(H_i/H_{i+1}) > 1$, let B be a proper subspace of H_i satisfying

$$H_i \supset B \supset H_{i+1},$$

so

$$[B, H_i, \dots, H_i] \subseteq [H_i, H_i, \dots, H_i] \subseteq H_{i+1} \subset B,$$

and

$$[H_{i+1}, B, \dots, B] \subseteq [H_i, H_i, \dots, H_i] \subseteq H_{i+1},$$

that is, B is an ideal of H_i and H_{i+1} is an ideal of B . In addition,

$$[H_i/B, \dots, H_i/B] = [H_i, \dots, H_i] + B \subseteq H_{i+1} + B = \{0\}$$

and

$$[B/H_{i+1}, \dots, B/H_{i+1}] = [B, \dots, B] + H_{i+1} \subseteq [H_i, \dots, H_i] + H_{i+1} = \{0\}$$

that is, H_i/B and B/H_{i+1} are abelian n -Lie algebras.

Iterating this processing several times, we can get a series of subalgebras of V satisfying the conditions in (4) since V is an finite dimensional n -Lie algebra.

(4) \Rightarrow (1) Since $\dim(B_i/B_{i+1}) = 1$, $B_m = 0$ and

$$B_{m-1}/B_m = B_{m-1}/0 = B_{m-1}, \dim B_{m-1} = \dim(B_i/B_{i+1}) = 1,$$

$i = 0, 1, \dots, m - 2$. In addition, one dimensional n -Lie algebra is solvable, then $B_{m-1}, B_i/B_{i+1}$ are solvable n -Lie algebras. By Lemma 3, all

$$B_{m-2}, B_{m-3}, \dots, B_1, B_0 = V$$

are solvable n -Lie algebras.

It is known that Borel subalgebras play an important role in the classification theory of Lie algebras. Now, we will give the concept of a Borel n -subalgebra of n -Lie algebra and prove two propositions related to it.

Definition 9. *The normalizer of a subalgebra H of an n -Lie algebra V is*

$$N_V(H) = \{x \in V \mid [x, H, \dots, H] \subseteq H\}.$$

H is called self-normalizing, if $H = N_V(H)$.

Definition 10. *A subalgebra B of an n -Lie algebra V is called a Borel n -subalgebra, if it is a maximal solvable subalgebra of V .*

Proposition 1. *Let B be a Borel n -subalgebra of an n -Lie algebra V over a field F , then $B = N_V(B)$, namely B is self-normalizing.*

Proof. Since B is a subalgebra of V , $[B, B, \dots, B] \subseteq B$, it is followed that $B \subseteq N_V(B)$. Conversely, if $x \in N_V(B)$, then $B_1 = B + Fx$ is a solvable subalgebra of V because of

$$[B_1, \dots, B_1] \subseteq [B, \dots, B] + [x, B, \dots, B] \subseteq B.$$

Hence $x \in B$ by the maximality of B . Therefore Proposition 1 holds.

Proposition 2. *If R is the solvable radical of an n -Lie algebra V , $R \neq V$, then the Borel n -subalgebras of V are in natural 1 – 1 correspondence with those of semi-simple n -Lie algebra V/R .*

Proof. If I is a solvable ideal of V , then $B + I$ is a solvable subalgebra of V , for any Borel n -subalgebra B of V , i.e. $I \subseteq B$. It is followed that $R \subseteq B$. In fact, $(B + I)^{(k)} \subseteq B^{(k)} + I$, for $k \geq 0$. by induction, when $k = 0$, it is trivial. Suppose $(B + I)^{(k)} \subseteq B^{(k)} + I$, then

$$\begin{aligned} (B + I)^{(k+1)} &\subseteq [(B + I)^{(k)}, \dots, (B + I)^{(k)}] \subseteq [B^{(k)} + I, \dots, B^{(k)} + I] \\ &\subseteq [B^{(k)}, \dots, B^{(k)}] + I = B^{(k+1)} + I \end{aligned}$$

Since B and I are solvable subalgebras, there exist r and s , such that $B^{(r)} = 0$ and $I^{(s)} = 0$, so $(B + I)^{(r+s)} = \{0\}$.

If φ is the canonical homomorphism from V into V/R , by Lemma 2, $\varphi(B) = B/R$ is a solvable subalgebra of V/R ; If B_1/R is a Borel n -subalgebra of V/R , then $\varphi^{-1}(B_1/R) = B_1 + R$ is solvable from the above discussion. Therefore Proposition 2 holds.

References

1. Nambu, Y.: Generalized Hamiltonian Mechanics. Phys. Rev. D7, 2405–2412 (1973)
2. Filippov, V.T.: N -Lie algebras. Sib. Mat. Zh. 26, 126–140 (1985)
3. Su, Y.C., Lu, C.H.: Introduction to the finite dimensional semisimple Lie algebra. Science Press, Beijing (2009)
4. Liu, J.B., Wang, X.M.: Linear algebra. Shanghai Jiaotong University Press, Shanghai (2012)
5. Kasymov, S.M.: On a Theory of n -Lie algebras. Algebra i Logika 26, 277–297 (1987)
6. Bai, R.P., Cheng, Y.: The Geometric Description of $(n-1)$ -Semisimple n -Lie Algebras. Acta Math. Appl. 33, 1087–1094 (2010)
7. Cheng, Y.: Decomposition of φ -free n -Lie algebra. J. of Baoding University 24, 8–9 (2011)
8. Cheng, Y., Meng, X.J.: A Criteria of φ -free n -Lie Algebras. Math. in Practice and Theory 40, 209–213 (2010)
9. Bai, R.P., Zhang, Z.X.: The inner derivation algebras of $(n+1)$ -dimensional n -Lie algebras. Comm. in Algebra 28, 2927–2934 (2000)
10. Saraiva, P.: Reduced n -Lie algebra. Comm. in Algebra 30, 2057–2074 (2002)

Modules of Lie Algebra $G(A)$

Zhang Yanyan¹, Liu Jianbo^{2,*}, Tao Wen², and Zhu Qin²

¹ School of Computer Science and Telecommunication Engineering

² School of Mathematics and Statistics,

Northeastern University at Qinhuangdao, 066004, Qinhuangdao, China

jbliu@mail.neuq.edu.cn

Abstract. Studied the structure and representations of Lie algebra $G(A)$, given a non-degenerate symmetric invariant bilinear form on $G(A)$, got the classification of the highest (lowest) weight modules on $G(A)$, and determined the maximal proper submodules when the highest (lowest) weight modules are reducible.

Keywords: Lie algebra, bilinear form, highest weight modules, lowest weight modules.

1 Introduction

It is well-known that the structure of non-symmetrizable Kac-Moody algebras are very complicated (see [1], [2], [3], [5], [4]). Even we do not know the multiplicities of imaginary roots of these Kac-Moody algebras (see [6], [7]). Our purpose of this paper is to expect to understand the structure of Kac-Moody algebras by studying related Lie algebras with simpler structure.

Recently, a class of interesting Lie algebras corresponding to symmetrizable Kac-Moody algebras was studied by Lu [8] and Zhang [9] where they studied finite-dimensional non-degenerate solvable Lie algebras. Generalizing their constructions, Liu, J. and Zhao, K., in [10], define the so-called deformed Kac-Moody algebras $G(A)$ associated to any generalized Cartan matrix A .

In this paper, we studied the structure and representations of $G(A)$ when $A = (2)_{1 \times 1}$, determined a non-degenerate symmetric invariant bilinear form on $G(A)$, and got the classification of the highest (lowest) weight modules.

2 Lie Algebra $G(A)$

Definition 1. An $n \times n$ integral matrix $A = (a_{ij})_{i,j=1}^n$ is called a generalized Cartan matrix (GCM) if

(C1). $a_{ii} = 2$, for all $i = 1, 2, \dots, n$;

(C2). $a_{ij} \leq 0$, for all $i \neq j$;

(C3). $a_{ij} = 0$ implies $a_{ji} = 0$.

* Corresponding author. This paper supported by "the Fundamental Research Funds for the Central Universities" China (N090323007).

In this paper we always assume that A is an $n \times n$ GCM, unless otherwise stated.

Let $g(A)$ be the Kac-Moody algebra associated to A , \mathfrak{h} the Cartan subalgebra of $g(A)$, $\Pi = \{\alpha_1, \alpha_2, \dots, \alpha_n\} \subseteq \mathfrak{h}^*$ the root basis, $\Pi^\vee = \{\alpha_1^\vee, \alpha_2^\vee, \dots, \alpha_n^\vee\} \subseteq \mathfrak{h}$ the coroot basis, and $e_1, e_2, \dots, e_n; f_1, f_2, \dots, f_n$ the Chevalley generators of $g(A)$. Note that $\langle \alpha_i^\vee, \alpha_j \rangle = a_{ij}$. Denote by Δ, Δ_+ and Δ_- the sets of all roots, positive roots and negative roots respectively. Set

$$Q = \sum_{i=1}^n \mathbb{Z}\alpha_i, \quad Q_\pm = \sum_{i=1}^n \mathbb{Z}_\pm \alpha_i$$

where \mathbb{Z}_+ (resp. \mathbb{Z}_-) is the set of all nonnegative (resp. nonpositive) integers. Then $\Delta = \Delta_+ \cup \Delta_-$ (a disjoint union), $\Delta_- = -\Delta_+$ and $\Delta_\pm = \Delta \cap Q_\pm$. The root space decomposition of $g(A)$ with respect to \mathfrak{h} is

$$g(A) = \sum_{\alpha \in \Delta_+} g_\alpha \oplus \mathfrak{h} \oplus \sum_{\alpha \in \Delta_+} g_{-\alpha}. \tag{1}$$

Let $\mathfrak{n}_+ = \sum_{\alpha \in \Delta_+} g_\alpha, \mathfrak{n}_- = \sum_{\alpha \in \Delta_+} g_{-\alpha}$ and $\mathfrak{b}_+ = \mathfrak{h} \oplus \mathfrak{b}_+$. Then \mathfrak{n}_+ (resp. \mathfrak{n}_-) is the subalgebra of $g(A)$ generated by e_1, e_2, \dots, e_n (resp. f_1, f_2, \dots, f_n), and g_α is the linear span of the elements of the form

$$[e_{i_1}, [e_{i_2}, [\dots [e_{i_{s-1}}, e_{i_s}] \dots]]] \text{ or } [f_{i_1}, [f_{i_2}, [\dots [f_{i_{s-1}}, f_{i_s}] \dots]]]$$

such that

$$\alpha_{i_1} + \alpha_{i_2} + \dots + \alpha_{i_s} = \alpha, \text{ (resp. } -\alpha)$$

In particular, $g_{\alpha_i} = \mathbb{C}e_i, g_{-\alpha_i} = \mathbb{C}f_i$, for $i = 1, 2, \dots, n$.

Let $(g(A), ad)$ be the adjoint representation of $g(A)$. Under this adjoint action, $g(A)$ can be regard as a \mathfrak{b}_+ -module. As a \mathfrak{b}_+ -module, $g(A)$ has a submodule \mathfrak{n}_+ . Hence we can obtain a \mathfrak{b}_+ -quotient module

$$\mathfrak{b}_- = g(A)/\mathfrak{n}_+ = \bar{\mathfrak{h}} \oplus \bar{\mathfrak{n}}_-. \tag{2}$$

It is clear that the set of weights of \mathfrak{b}_- with respect to \mathfrak{h} is $P(\mathfrak{b}_-) = \{0\} \cup \Delta_-$ and $\mathfrak{b}_+ \cdot \bar{\mathfrak{h}} = \{0\}$. For simplicity, we write the action of \mathfrak{b}_+ -module

$$x \cdot v, \forall x \in \mathfrak{b}_+, v \in \mathfrak{b}_-.$$

Now let us define the Lie algebra $G(A)$ associated to A as follows.

Set

$$G(A) = \mathfrak{b}_+ \oplus \mathfrak{b}_-. \tag{3}$$

Define the following bracket operator $[\cdot, \cdot]$ on $G(A)$:

$$[x, y] = [x, y]_0, \quad \forall x, y \in \mathfrak{b}_+; \tag{4}$$

$$[v_1, v_2] = 0, \quad \forall v_1, v_2 \in \mathfrak{b}_-; \tag{5}$$

$$[x, v] = -[v, x] = x \cdot v, \quad \forall x \in \mathfrak{b}_+, \forall v \in \mathfrak{b}_-, \tag{6}$$

where $[\cdot, \cdot]_0$ is the bracket operator on $g(A)$. It is easy to show that $G(A)$ becomes a Lie algebra under the above bracket operator $[\cdot, \cdot]$.

It is clear that, as vector spaces, $\mathfrak{h} \oplus \mathfrak{n}_- \cong \mathfrak{b}_-$. Let π be the canonical homomorphism from $g(A)$ onto \mathfrak{b}_- . Then $\sigma = \pi|_{\mathfrak{h} \oplus \mathfrak{n}_-}$ is an isomorphism between $\mathfrak{h} \oplus \mathfrak{n}_-$ and \mathfrak{b}_- , such that $\bar{\mathfrak{h}} = \sigma(\mathfrak{h})$, $\bar{\mathfrak{b}}_- = \sigma(\mathfrak{b}_-)$. We see that

$$[x, \sigma(y)] = \sigma([x, y]_0), \quad \text{for } x \in g_\alpha, y \in g_{-\beta}; \alpha \leq \beta, \alpha, \beta \in \Delta_+.$$

By the construction of $g(A)$, $(\mathfrak{h}, \Pi, \Pi^\vee)$ is a realization of A . We supplement $\alpha_{n+1}^\vee, \dots, \alpha_{2n-l}^\vee$ to $\Pi^\vee = \{\alpha_1^\vee, \alpha_2^\vee, \dots, \alpha_n^\vee\}$ to form a basis of \mathfrak{h} , where l is the rank of GCM A . Denote $z_i = \sigma(\alpha_i^\vee)$, $1 \leq i \leq 2n-l$. Thus $z_1, z_2, \dots, z_{2n-l}$ form a basis of $\bar{\mathfrak{h}}$. For any $y \in \mathfrak{b}_-$, we still write its image $\sigma(y)$ in $\bar{\mathfrak{b}}_-$ as y , i.e., write elements in $\bar{\mathfrak{b}}_-$ as elements in the Lie algebra $(\mathfrak{b}_-, [\cdot, \cdot]_0)$ by using f_1, f_2, \dots, f_n . Using these notations, we deduce that

$$[e_i, f_j] = \delta_{ij} z_i, \quad \text{for } 1 \leq i \leq n.$$

Let $\mathfrak{H} = \mathfrak{h} \oplus \bar{\mathfrak{h}}$. We can consider \mathfrak{H} as a Cartan subalgebra of $G(A)$. For any $\alpha \in \Delta_+$, define $\tilde{\alpha} \in \mathfrak{H}^*$ such that $\tilde{\alpha}|_{\mathfrak{h}} = \alpha$ and $\tilde{\alpha}|_{\bar{\mathfrak{h}}} = 0$. Thus $\tilde{\Delta} = \{\pm \tilde{\alpha} \in \mathfrak{H}^* \mid \alpha \in \Delta_+\}$ is the set of all roots of $G(A)$. In addition, for any $\alpha \in \Delta_+$, the root space attached to $\tilde{\alpha}$ (resp. $-\tilde{\alpha}$) is $G_{\tilde{\alpha}} = g_\alpha$ (resp. $G_{-\tilde{\alpha}} = g_{-\alpha}$). Hence, $\tilde{\alpha}$, $-\tilde{\alpha}$ and $\tilde{\Delta}$ can be identified with α , $-\alpha$ and Δ respectively. So we get the root space decomposition of $G(A)$ with respect to \mathfrak{H} :

$$G(A) = \sum_{\alpha \in \Delta_+} G_\alpha \oplus \mathfrak{h} \oplus \sum_{\alpha \in \Delta_+} G_{-\alpha}.$$

Denote $G_+ = \sum_{\alpha \in \Delta_+} G_\alpha$, $G_- = \sum_{\alpha \in \Delta_+} G_{-\alpha}$. Then we have the triangular decomposition of $G(A)$:

$$G(A) = G_+ \oplus \mathfrak{H} \oplus G_-.$$

Hence the universal enveloping algebra $U(G(A))$ of $G(A)$ can be factored as

$$U(G(A)) = U(G_+) \otimes U(\mathfrak{H}) \otimes U(G_-).$$

Now we collect some properties of $G(A)$ in the following lemma.

- Lemma 1.** 1) The set of all roots of $G(A)$ with respect to \mathfrak{H} is Δ ;
 2) As vector spaces, G_- is isomorphic to \mathfrak{b}_- ;
 3) $\bar{\mathfrak{h}} \oplus G_-$ is an abelian subalgebra of $G(A)$, $\bar{\mathfrak{h}}$ is in the center of $G(A)$;
 4) $\bar{\mathfrak{h}} \oplus G_-$ is a $\mathfrak{h} \oplus G_+$ -module, and $G_+ \cdot \bar{\mathfrak{h}} = \{0\}$.

Recall that a \mathbb{C} -valued symmetric bilinear form (\cdot, \cdot) on a complex Lie algebra g is said to be *invariant* if

$$([x, y], z) = (x, [y, z]), \quad \text{for all } x, y, z \in g. \tag{7}$$

Also recall that a complex $n \times n$ matrix A is called *symmetrizable* if there exists a non-singular $n \times n$ diagonal matrix D and an $n \times n$ symmetric matrix B such that

$$A = DB.$$

Lemma 2. ([3]) *Let $A = (a_{ij})_{i,j=1}^n$ be a symmetrizable GCM. Then the diagonal matrix*

$$D = \text{diag}(\epsilon_1, \epsilon_2, \dots, \epsilon_n)$$

in above equation can be chosen so that $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ are positive rational numbers. If, further, A is indecomposable, then the matrix D is unique up to a constant factor.

Theorem 1. ([10]) *The Lie algebra $G(A)$ has a non-degenerate symmetric invariant bilinear form (\cdot, \cdot) if and only if A is symmetrizable.*

3 Conclusion

Now we assume that A is one order GCM, i.e., $A = (2)_{1 \times 1}$. In this section, we will study the structure and the representations of the Lie algebra $G(A)$.

Let $A = (2)_{1 \times 1}$. Then we have

$$g(A) = sl_2 = \mathbb{C}y \oplus \mathbb{C}h \oplus \mathbb{C}x. \tag{8}$$

with defining relations:

$$[x, y] = h, [h, x] = 2x, [h, y] = -2y. \tag{9}$$

Thus $\mathfrak{b}_+ = \mathbb{C}h \oplus \mathbb{C}x$ and the \mathfrak{b}_+ -module $\mathfrak{b}_- = \mathbb{C}z \oplus \mathbb{C}y$ generated by the weight vector with the lowest weight -2. Hence

$$G(A) = \mathfrak{b}_+ \oplus \mathfrak{b}_- = \mathbb{C}h \oplus \mathbb{C}z \oplus \mathbb{C}x \oplus \mathbb{C}y. \tag{10}$$

and the bracket operator $[\cdot, \cdot]$ on $G(A)$ can be written as follows:

$$\begin{aligned} [h, x] &= 2x, [h, z] = 0, [h, y] = -2y, \\ [x, z] &= 0, [x, y] = z, [z, y] = 0. \end{aligned} \tag{11}$$

Theorem 2. *Let $A = (2)_{1 \times 1}$. Then there exists a non-degenerate symmetric invariant bilinear form (\cdot, \cdot) on $G(A)$, and more precisely we have a non-degenerate symmetric invariant bilinear form as follows:*

$$(x, y) = 1, \text{ and } (h, z) = 2,$$

and products of all other pairs of the basis elements are zero.

Proof. It is clear that A is symmetrizable. So, by Theorem 1, we know that there exists a non-degenerate symmetric invariant bilinear form (\cdot, \cdot) on $G(A)$. Now we suppose that (\cdot, \cdot) is a non-degenerate symmetric invariant bilinear form on $G(A)$. It is easy to see that

$$(G_\alpha, G_\beta) = \{0\}, \quad \forall \alpha + \beta \neq 0.$$

Hence we only need to determine the value of (x, y) , (h, h) , (z, z) and (h, z) .

Denote by $(\cdot, \cdot)_0$ the bracket operator on $g(A)$. We know, by the proof of Theorem 1, that we can get a non-degenerate symmetric invariant bilinear form (\cdot, \cdot) on $G(A)$ via a non-degenerate symmetric invariant bilinear form $(\cdot, \cdot)_0$ on $g(A)$. Since the standard non-degenerate symmetric invariant bilinear form on $g(A) = sl_2$ is the Killing form, i.e.,

$$(x, y)_0 = 1, \text{ and } (h, h)_0 = 2,$$

Hence we have

$$(x, y) = (x, y)_0 = 1 \text{ and } (h, z) = (h, h)_0 = 2$$

and products of all other pairs of the basis elements are zero.

Also, we know that $\mathfrak{H} = \mathbb{C}h \oplus \mathbb{C}z$. For $\lambda \in \mathfrak{H}^*$, λ is uniquely determined by $\lambda(h)$ and $\lambda(z)$. Let $\lambda(h) = m$ and $\lambda(z) = n$. We define the 1-dimensional $\mathfrak{H} \oplus \mathbb{C}y$ -module $\mathbb{C}\phi_\lambda$ via

$$\begin{cases} (ah + bz) \cdot \phi_\lambda = (am + bn)\phi_\lambda, & \forall a, b \in \mathbb{C} \\ y \cdot \phi_\lambda = 0. \end{cases}$$

Thus the lowest weight module on $G(A)$ is

$$\bar{V}(\lambda) = \text{Ind}_{\mathfrak{H} \oplus \mathbb{C}y}^{G(A)} \mathbb{C}\phi_\lambda = U(G(A)) \otimes_{U(\mathfrak{H} \oplus \mathbb{C}y)} \mathbb{C}\phi_\lambda.$$

And $\bar{V}(\lambda) \simeq U(\mathbb{C}x) \simeq \mathbb{C}[x]$ as vector spaces, where $\mathbb{C}[x]$ is the polynomial algebra of one indeterminate. Let $\alpha \in \mathfrak{H}^*$ with $\alpha(h) = 2$ and $\alpha(z) = 0$. Then we have the weight space decomposition of $\bar{V}(\lambda)$:

$$\bar{V}(\lambda) = \sum_{i=0}^{\infty} \bar{V}(\lambda)_{\lambda+i\alpha},$$

where

$$\bar{V}(\lambda)_{\lambda+i\alpha} = \mathbb{C}x^r \phi_\lambda$$

and

$$\dim \bar{V}(\lambda)_{\lambda+i\alpha} = 1.$$

Since $\bar{V}(\lambda)$ is spanned by $x^r \phi_\lambda$, $r \in \mathbb{Z}$, we can write the action of basis elements of $G(A)$ on $\bar{V}(\lambda)$ as follows:

$$\begin{cases} h \cdot x^r \phi_\lambda = (m + 2r)x^r \phi_\lambda, \\ z \cdot x^r \phi_\lambda = nx^r \phi_\lambda, \\ x \cdot x^r \phi_\lambda = x^{r+1} \phi_\lambda, \\ y \cdot x^r \phi_\lambda = -rn x^{r-1} \phi_\lambda. \end{cases}$$

From the action, we know that $\bar{V}(\lambda)$ is an irreducible module over $G(A)$ if $n = \lambda(z) \neq 0$.

If $n = \lambda(z) \neq 0$, then $\bar{V}(\lambda)$ is reducible. It has submodules

$$V_i = \text{Span}\{x^{i+r}\phi_\lambda | r \in \mathbb{Z}_+\}.$$

It is clear that $J = V_1$ is the maximal proper submodule of $\bar{V}(\lambda)$. The quotient module $L(\lambda) = \bar{V}(\lambda)/J$ is the 1-dimensional module with action as follows:

$$h \cdot \phi_\lambda = m \cdot \phi_\lambda, \quad z \cdot \phi_\lambda = 0, \quad x \cdot \phi_\lambda = 0, \quad y \cdot \phi_\lambda = 0.$$

From the action and the discussions above, we know that

Theorem 3. *Let $A = (2)_{1 \times 1}$. Then the lowest weight module $\bar{V}(\lambda)$ over $G(A)$ is an irreducible module if and only if $n = \alpha(z) \neq 0$.*

Similarly, for any $\lambda \in \mathfrak{H}^*$, we define the 1-dimensional $\mathfrak{H} \oplus \mathbb{C}x$ -module $\mathbb{C}\tilde{\phi}_\lambda$ via

$$\begin{cases} (ah + bz) \cdot \tilde{\phi}_\lambda = (am + bn)\tilde{\phi}_\lambda, & \forall a, b \in \mathbb{C} \\ x \cdot \tilde{\phi}_\lambda = 0. \end{cases}$$

Thus the highest weight modules over $G(A)$ is

$$\tilde{V}(\lambda) = \text{Ind}_{\mathfrak{H} \oplus \mathbb{C}x}^{G(A)} \mathbb{C}\tilde{\phi}_\lambda = U(G(A)) \otimes_{U(\mathfrak{H} \oplus \mathbb{C}x)} \mathbb{C}\tilde{\phi}_\lambda.$$

As the discussion of the lowest weight module, we can deduce:

Theorem 4. *Let $A = (2)_{1 \times 1}$. Then the highest weight module $\tilde{V}(\lambda)$ over $G(A)$ is an irreducible module if and only if $n = \lambda(z) \neq 0$. Moreover,*

$$J = \text{Span}\{y^{1+r}\tilde{\phi}_\lambda | r \in \mathbb{Z}_+\}.$$

is the maximal proper submodule of $\tilde{V}(\lambda)$ when $n = \lambda(z) \neq 0$.

References

1. Kac, V.G.: Infinite dimensional Lie algebras, 3rd edn. Cambridge University Press, New York (2006)
2. Moody, R.V., Pianzola, A.: Lie algebras with triangular decompositions. Canadian Mathematical Society Series of Monographs and Advanced Texts. A Wiley-Interscience Publication, John Wiley Sons, Inc., New York (1995)
3. Wan, Z.: Introduction to Kac-Moody algebra. World Scientific (1991)
4. Su, Y.C., Lu, C.H.: Introduction to the finite dimensional semisimple Lie algebra. Science Press, Beijing (2009)
5. Liu, J.B., Wang, X.M.: Linear algebra. Shanghai Jiaotong University Press, Shanghai (2012)
6. Benkart, G., Kang, S.J., Misra, K.C.: Indefinite Kac-Moody algebras of classical type. Adv. Math. 105, 76–110 (1994)
7. Benkart, G., Kang, S.J., Misra, K.C.: Weight multiplicity polynomials for affine Kac-Moody algebras of type $A_r^{(1)}$. Compositio Math. 104, 153–187 (1996)
8. Zhang, H., Lu, C.: The isomorphic realization of non-degenerate solvable Lie algebras of maximal rank. Algebra Colloq. 15, 347–360 (2008)
9. Zhang, H.: A class of non-degenerate solvable Lie algebras and their derivations. Acta Math. Sin. (Engl. Ser.) 24, 7–16 (2008)
10. Liu, J., Zhao, K.: Deformed Kac-Moody algebras and their representations. J. Algebra 319, 4692–4711 (2008)

Detection Performance Analysis of Tests for Spread Targets in Compound-Gaussian Clutter

Xiandong Meng, Zhiming He, Xiaowei Niu, and Ganzhong Feng

School of Electronic Engineering, University of Electronic Science and Technology of China,
Chengdu 611731, China
xdmeng_uestc@163.com

Abstract. The problem of adaptive detection of spatially distributed targets or targets embedded in compound-Gaussian clutter with unknown covariance matrix is studied. At first, the test decision statistic of the generalized likelihood ratio test (GLRT), Rao test and Wald test which have been widely applied to the distributed targets detection of modern Wideband radar is derived. Next, the numerical results are presented by means of Monte Carlo simulation strategy. Assume that cells of signal components are available. Those secondary data are supposed to possess either the same covariance matrix or the same structure of the covariance matrix of the cells under test. In this context, the simulation results highlight that the asymptotic properties of the three tests in different coherent train pulses, and that the performance loss of the real target length mismatches the setup in receiver.

Keywords: GLRT, Rao test, Wald test, distributed targets, compound-Gaussian clutter.

1 Introduction

The problem of detecting spread targets has received great attention recently. It naturally arises that the detecting ability of high resolution radars (HRRs). The clutter is very complex and it can not be considered as Gaussian distributed in HRRs. Much work has been directed so far towards the clutter model in HRRs. And the spikiness clutter is usually modeled as a compound-Gaussian vector [2, 6, 8-11].

Various adaptive detection algorithms in non-Gaussian background have been studied [2, 3, 6-11]. To avoid complicated integral computation of the optimum NP detector, these papers resort to the so called two-step GLRT detection scheme which is proposed by Robey [5], and is adopted in many HRR GLRT-based detection algorithms [1-4, 9, 12]. About the test design, Conte et al. [1] Derived the two-step adaptive detector-Rao test and Wald test, and gave the formula of the test decision statistic, which are widely applied in modern wideband radar. Conte et al. [2] also investigated the NP detector versus the GLRT detector. De Maio et al. [3] derived the coincidence of the Rao Test, Wald Test, and GLRT in Partially Homogeneous Environment. N. Bon et al. [4] derived GLRT subspace detection for range and Doppler distributed targets.

In this paper, we mainly focus on the problem of adaptive detection for spatially distributed targets in compound-Gaussian background. Precisely, we focus on the detection performance of GLRT, Rao test and Wald test. At first, based on the works above cited we derive the test decision statistic of the three tests. Then, the performance analyses of the tests are carried out via Monte Carlo simulations. And we mainly consider the asymptotic properties of the three tests in different coherent train pulses; then, the performance loss of the real target length mismatches the setup in receiver is simulated.

2 Problem Statement

Assume that a coherent train of N pulses is transmitted by the radar and that the incoming waveform of the receiver is properly demodulated, filtered and sampled. The binary hypothesis can be written as:

$$\begin{cases} H_0 : z_t = c_t & t = 1, 2, \dots, L, L+1 \dots L+K \\ H_1 : \begin{cases} z_t = \alpha_i p + c_t & t = 1, 2, \dots, L \\ z_t = c_t & t = L+1, \dots, L+K \end{cases} \end{cases} \quad (1)$$

Where $z = [z(0) \dots z(N-1)]^T$ is an N -dimensional vector and the N are complex samples, T denotes the transpose operator, p denotes the steering vector, and the α_i is unknown deterministic parameter which account for the channel propagation effects and the target reflectivity. $z_1, z_2 \dots z_L$ are collected from cells under test that are referred as primary data, and $z_{L+1}, z_{L+2} \dots z_{L+K}$ are secondary data which not contain any useful target echo and exhibit the same structure of the covariance matrix as the primary data. Here the received data vectors are assumed to be independent between each range cell.

As to the clutter $c_1, c_2 \dots c_{L+K}$ in (1), they are modeled as compound-Gaussian vectors; $c_t = \sqrt{\tau_t} g_t$, here the speckle g_t is modeled as a zero mean complex Gaussian vector with covariance matrix

$$E\{xx^H\} = M \quad (2)$$

Where $E\{\cdot\}$ denotes statistical expectation operator. The texture τ_t is a positive random variable with an unknown PDF.

3 The Test Design

3.1 GLRT Design

The GLRT design has been derived by many authors [1, 13], which can be written as

$$T_G = -N \sum_{t=1}^L \ln \left(1 - \frac{|p^H M^{-1} z_t|^2}{(p^H M^{-1} p)(z_t^H M^{-1} z_t)} \right) \quad (3)$$

3.2 Rao Test and Wald Test Design

The Rao test decision statistic can be written as [3, 13]

$$T_R = \sum_{i=1}^L \frac{|p^H M^{-1} z_i|^2}{p^H M^{-1} p \cdot z_i^H M^{-1} z_i} \tag{4}$$

And the Wald test decision statistic can be written as [3, 13]

$$T_W = \sum_{i=1}^L \frac{|p^H M^{-1} z_i|^2}{p^H M^{-1} p \left(z_i^H M^{-1} z_i - \frac{|p^H M^{-1} z_i|^2}{p^H M^{-1} p} \right)}. \tag{5}$$

4 Detection Performance Assessment

The performance analysis of the test is carried out via Monte Carlo simulations. Assume the clutter is K-distributed. In the simulation below, we mainly review the detection preference of tests in different parameters, and assuming $p_{fa} = 10^{-4}$ throughout the section. L is related to the range extent of the target and the range resolution of HRRs. We consider small values of $L(L \leq 10)$ and $N=8$ to save simulation time.

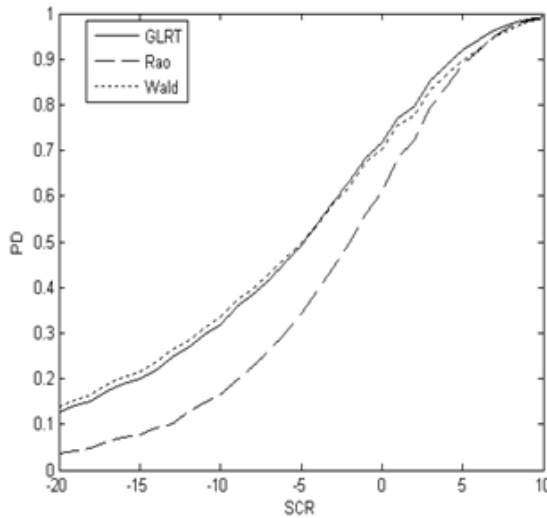


Fig. 1. Detection performance of GLRT, Rao and Wald test for the numbers of pulses $N=8$

4.1 The Asymptotic Properties of GLRT, Rao Test and Wald Test

We mainly analyze the detection performance of the three detectors in different N , and give the condition when the detection performance of the tests tends to coincidence.

The Fig. 1 gives the detection performance analysis of the three detectors in different sample numbers $N=8$. The detection performance of the Rao test is obviously appreciably worse than the other two tests when $N=8$, and the detection performance of the Wald test and GLRT is closely. When SCR is low, the detection performance of Wald test is appreciably better than GLRT. While, when SCR is high, the detection performance of Wald test is worse than GLRT.

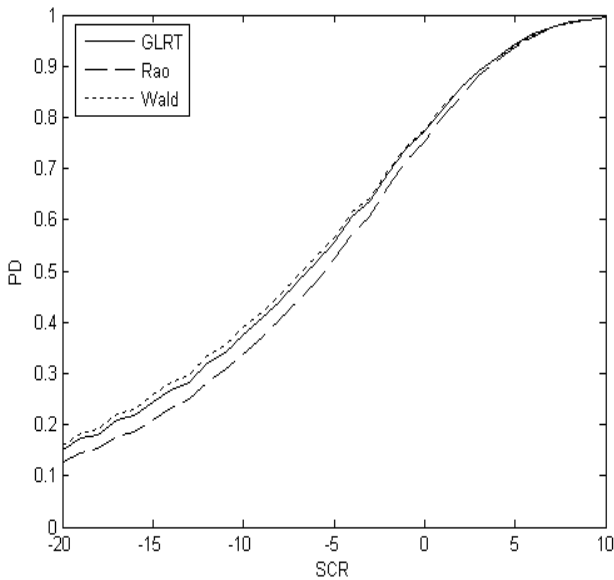


Fig. 2. Detection performance of GLRT, Rao and Wald test for the numbers of pulses $N=16$

In Fig. 2, the detection performance of the three detectors are very close when $N=16$, and only a little difference when SCR is low. According to the results above, we can find that the detection performance of Rao test is correspondingly worse, while the three tests have the same asymptotic properties.

4.2 The Performance Loss with Different L

Assume that setup the target occupy 10 range cells in the receiver, while, in fact, the target occupy $10(M1)$, $6(M2)$ and $2(M3)$ respectively, i.e. the target has the three models in the TABLE 1, assume the total energy of every model is normalized.

Table 1. The energy of every spread range target cell

	Range cell $ \alpha, l^2$									
	1	2	3	4	5	6	7	8	9	10
M1	1/10	1/10	1/10	1/10	1/10	1/10	1/10	1/10	1/10	1/10
M2	1/6	1/6	1/6	1/6	1/6	1/6	0	0	0	0
M3	1/2	1/2	0	0	0	0	0	0	0	0

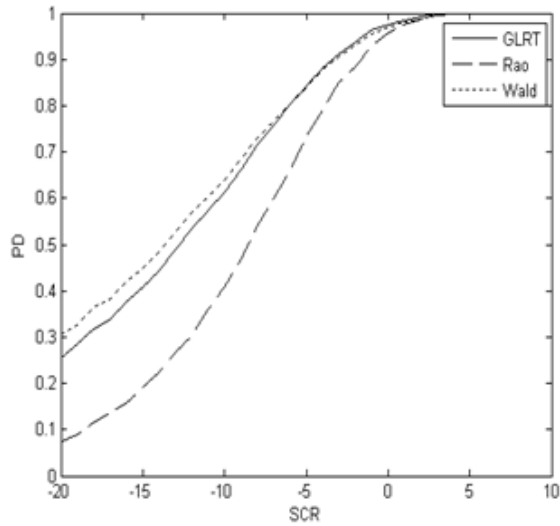


Fig. 3. The detection performance of GLRT, Rao test and Wald test shows in model M1

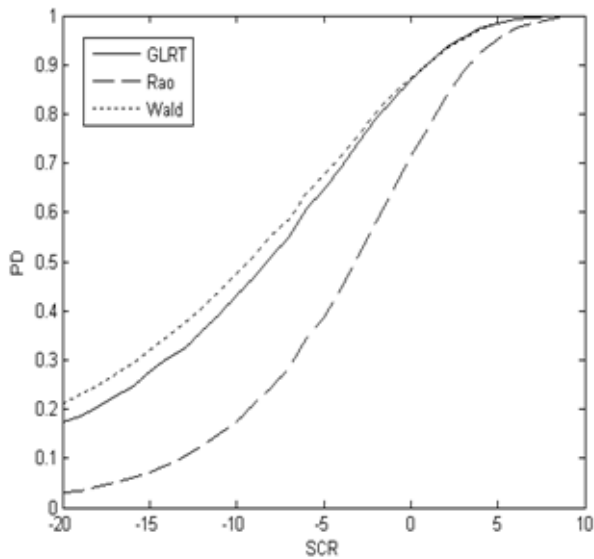


Fig. 4. The detection performance of GLRT, Rao test and Wald test shows in model M2

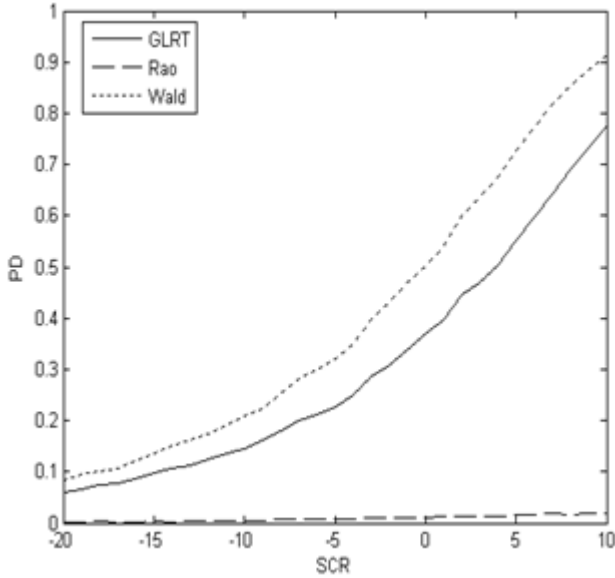


Fig. 5. The detection performance of GLRT, Rao test and Wald test shows in model M3

Fig. 3, Fig. 4 and Fig. 5 respectively give the detection performance of target model in M1, M2 and M3, and the parameters reference in the simulation as: $\nu=0.5$, $N=8$ and $L=10$. In the Fig. 3 display the result of target model in M1, and the real length of target matches the value L , then the receiver correspondingly shows a super detection performance. In the Fig. 4, the real length of target mismatches the value L , so all the detection performance of the three tests in Fig. 4 are worse than in the Fig. 3. While, when the real target length is seriously mismatches the value L , then the performance of the receiver badly drop, especially the Rao test already cannot work normally, and the performance of Wald test appreciably better than GLRT.

5 Conclusions

In this paper, we have addressed the problem of adaptive detection of range-spread targets in compound Gaussian clutter. Specifically, we have shown that: the asymptotic properties of GLRT, Rao test and Wald test; the performance loss with different L ; Performance analysis of the tests when the target amplitudes are fluctuant.

In conclusion, we can state that: the asymptotic properties of GLRT, Rao test and Wald test are coincidence; with limit sample numbers, the detection performance of GRLT and Wald test are appreciably better than Rao test; the design of Rao test is simpler than the Wald test and GLRT. All the three tests are free of the apriority information of the texture, so the three tests are suit to all compound Gaussian clutter environment, but to some special compound Gaussian clutter, the three tests are not Uniformly Most Powerful (UMP), and not the optimal.

Acknowledgments. The authors are very grateful to the anonymous referees for their careful reading, helpful comments. This work was supported by the Fundamental Research Funds for the Central Universities (ZYGX2011J013, ZYGX2009J017). National Nature Science Foundation of China and China Academy of Engineering Physics (No.10876006).

References

1. Conte, E., De Maio, A., Ricci, G.: GLRT-based adaptive detection algorithms for range-spread targets. *IEEE Transactions on Signal Processing* 49(7), 1336–1348 (2001)
2. Conte, E., De Maio, A., Galdi, C.: Signal detection in compound-Gaussian noise: Neyman–Pearson and CFAR detectors. *IEEE Transactions on Signal Processing* 48(2), 419–428 (2000)
3. De Maio, A., Iommelli, S.: Coincidence of the Rao Test, Wald Test, and GLRT in Partially Homogeneous Environment. *IEEE Signal Processing Letters* 15 (2008)
4. Bon, N., Khenchaf, A., Garello, R.: GLRT subspace detection for range and doppler distributed targets. *IEEE Transactions on Aerospace and Electronic Systems* 44(2), 678–695 (2008)
5. Robey, F.C., Fuhrmann, D.R., Nitzberg, R., Kelly, E.J.: A CFAR adaptive matched filter detector. *IEEE Transactions on Aerospace and Electronic Systems* 28(1), 208–216 (1992)
6. Conte, E., De Maio, A., Ricci, G.: Asymptotically optimum radar detection in compound Gaussian clutter. *IEEE Transactions on Aerospace and Electronic Systems* 31(2), 617–625 (1995)
7. Bandiera, F., De Maio, A., Ricci, G.: Adaptive Radar Detection of Distributed Targets in Homogeneous and Partially Homogeneous Noise Plus Subspace Interference. *IEEE Transactions on Signal Processing* 55(4) (April 2007)
8. Gini, F., Farina, A.: Vector subspace detection in compound-Gaussian clutter, part I: survey and new results. *IEEE Transactions on Aerospace and Electronic Systems* 38, 1295–1311 (2002)
9. Shuai, X., Kong, L., Yang, J.: Performance analysis of GLRT-based adaptive detector for distributed targets in compound-Gaussian clutter. *Signal Processing* 90, 16–23 (2010)
10. Shuai, X., Kong, L., Yang, J.: AR-model-based adaptive detection of range-spread targets in compound Gaussian clutter. *Signal Processing* 91, 750–758 (2011)
11. Pascal, F., Chitour, Y., Forster, P., Larzabal, P.: Covariance structure maximum-likelihood estimates in compound Gaussian noise existence and algorithm analysis. *IEEE Transactions on Signal Processing* 56(1), 34–48 (2008)
12. CFAR adaptive subspace detector is a scale-invariant GLRT. *IEEE Trans. on Sign. Process* 47, 2538–2541 (1999)
13. Shuai, X.: Research on Detection Algorithms of Range-Targets. Doctoral Dissertation, 68–77 (2011)

Rough Differential Equations in Rough Function Model

Yun Wang¹, Xiaojing Xu², and Zhiqin Huang¹

¹ School of Mathematics Sciences, University of Jinan, 250022 Jinan, P.R. China

² School of Computer Science and Technology, Shandong Architecture University, 250001

Jinan, P.R. China

ss_wangy1@ujn.edu.cn

Abstract. To describe rough differential equations in rough function model, a series of new concepts of rough differential equations are proposed by generalizing theory of difference equations and ordinary differential equations into rough function model. Based on these basic concepts, fundamental properties of linear rough differential equations solutions are discussed and four principles of rough superposition are put forward, which lay foundations for solving methods of rough differential equations. Several insufficiencies in original solving methods of rough differential equations are pointed out and improved respectively. Three types of typical rough differential equations are given. According to characteristics of different types of rough differential equations, corresponding solving methods and solving steps as well as different forms of solutions are discussed. Finally, some general remarks are made on other types of rough differential equations.

Keywords: rough set, rough function, rough differential equation.

1 Introduction

The theory of rough sets [1] is a powerful mathematical approach to deal with imprecision, uncertainty and vagueness in data analysis. The development and achievement of rough set theory are outstandingly manifested in various deformation and generalization of rough set model from the initial rough set model [2].

Although many theoretical and practical problems relating to data analysis have been solved successfully by rough set theory, due to the limitation of rough set theory being based on set theory, a large number of theoretical and applicable problems having strong connections with function theory, such as synthesis and analysis of rough controllers, generation and optimization of discrete dynamic systems [3-6] and so on, could not be described and solved only by the lower and upper approximate sets in rough set theory. Therefore, Pawlak [6-8] generalized the concepts of rough sets into real numbers domain and provided descriptions of rough functions in real numbers domain. The researches on rough function model supplement and develop the new theory based on rough set theory, and lay dependable foundations for continuing researches.

On rough differential equations in rough function model, [8] gave the general form of first-order rough differential equations. Then from an example of solving a simple rough differential equation, it pointed out that there were two methods of solving rough differential equations. The direct basis of the both methods was a proposition in

the part of rough integrals given by [8]. However, there exist several insufficiencies in applications of the proposition. In addition, essential theories were not given by [8], including basic concepts necessary for intensively study of rough differential equations, fundamental properties of solutions, and so on. Moreover, general forms of solutions with respect to different types of rough differential equations of simple common were not developed yet. In the light of the above problems, the investigations on rough differential equations to be discussed in this paper are as follows. (1) A series of essential concepts of rough differential equations are given in section 2, such as rough differential equations of first order and n -th order, their initial value problems, linear homogeneous (non-homogeneous) rough differential equations, groups of rough differential equations, etc. (2) In section 3, rough linear superposition principles of linear rough differential equations are proposed, which lay theoretical foundations for discussing solving methods of rough differential equations and groups of equations. (3) In section 4, insufficiencies in original solving methods of rough differential equations are pointed out. Corresponding solving methods of different types of common rough differential equations are analyzed, and general forms of the solutions are given. And finally, section 5 concludes the whole paper.

Some preliminary concepts necessary to this paper such as rough functions and rough derivatives, etc. can be seen in references

2 Basic Concepts of Rough Differential Equations

A rough differential equation [8] is a equation which contains single-variable unknown rough functions discretely valuing in $[n]$ and their rough derivatives.

Definition 1. Let $f^{(k)}$ ($k \in Z^+$) be the k -th order rough derivative of the one-variable rough function f , then the rough differential equation

$$\phi(i, f(i), f'(i), \dots, f^{(n)}(i)) = 0 \tag{1}$$

is called the rough differential equation of n -th order, where Φ is a rough function of several variables.

When $n=1$, the rough differential equation of n -th order is referred to the first-order rough differential equation [8]:

$$f'(i) = \phi(i, f(i)) \tag{2}$$

Definition 2. The problem of solving the n -th order rough differential equation (1) satisfying initial conditions $f(0)=j_0, f'(0)=j_0', \dots, f^{(n-1)}(0)=j_0^{(n-1)}$ ($j_0^{(i)} \in [m], i=0, 1, \dots, n-1$) is called the initial values problem of the n -th order rough differential equation, which is also called the Cauchy problem. The $f(i)$, satisfying the above conditions, is referred to as the solution of the Cauchy problem.

Rough differential equations are divided into two types of the linear and the non-linear, according to the rough unknown function items and their rough derivative items appearing linearly or non-linearly in rough differential equations.

Definition 3. The general form of a linear rough differential equation of n -th order is defined as

$$f^{(n)}(i) + a_1(i)f^{(n-1)}(i) + \dots + a_{n-1}(i)f'(i) + a_n(i)f(i) = g(i) \tag{3}$$

where $a_l(i)$ ($l=1, 2, \dots, n$), $g(i)$ are rough functions in $[n]$.

When $g(i) \equiv 0$, the expression (3) is called the linear homogeneous rough differential equation of n -th order, denoted as (3)'. When $g(i) \neq 0$, (3) is called the linear non-homogeneous rough differential equation of n -th order, denoted as (3)''.

In the theory of differential equations, the so-called ‘‘integrating’’ a differential expression and ‘‘differentiating’’ an integral expression is one of the most important skills, and is a fundamental tool of study various properties of solutions¹⁸. The use of indefinite integrals causes that the result of the solution contains a random constant, so that there is the concept of the ordinary solution correspondingly. While in the theory of rough differential equations, rough functions are defined in $[n]$. When solving rough differential equations, we only take rough integrations with the lower limitation 0 and the upper limitation i . Therefore, the constant needs to be determined in the result of the solution is $f(0)$, and the significance of the ‘‘ordinary’’ solution is not obvious. Hence, the solution containing the constant $f(0)$ and the solution of the Cauchy problem with determined $f(0)$ are not to be delineated here.

Definition 4. Let $\vec{f}(i)$ and $\vec{g}(i)$ be n -dimensional rough function vectors, $A(i)$ be a $n \times n$ matrix, we call the group of equations

$$\vec{f}'(i) = A(i)\vec{f}(i) \tag{4}$$

is a linear homogeneous rough differential equations group.

The group of equations

$$\vec{f}'(i) = A(i)\vec{f}(i) + \vec{g}(i) \tag{5}$$

is defined as a linear non-homogeneous rough differential equations group, where $\vec{g}(i) \neq \vec{0}$.

3 Properties of Linear Rough Differential Equation Solutions

Linear rough differential equations are bases of studying non-linear rough differential equations. Therefore, the theory of linear rough differential equations is quite important in both theory and applications. This paper only discusses linear rough differential equations and their applications. Nonlinear rough equations will be discussed in the future.

Theorem 1 (Principle 1 of Rough Superposition). If both $f_1(i)$ and $f_2(i)$ are solutions of (3)', C_1 and C_2 are random integer constants, then $C_1f_1(i) + C_2f_2(i)$ is also a solution of (3)'. And if $\vec{\psi}_1(i)$ and $\vec{\psi}_2(i)$ are solution vectors of the equations group (4), then $C_1\vec{\psi}_1(i) + C_2\vec{\psi}_2(i)$ is also a solution vector of (4).

Proof. Substituting $f(i) = C_1f_1(i) + C_2f_2(i)$ in the left side of (3)', we have $[C_1f_1(i) + C_2f_2(i)]^{(n)} + a_1(i)[C_1f_1(i) + C_2f_2(i)]^{(n-1)} + \dots + a_n(i)[C_1f_1(i) + C_2f_2(i)]$. By rough derivatives operating rules [17], the above expression can be turned into $C_1[f_1^{(n)}(i) + a_1(i)f_1^{(n-1)}(i) +$

$\dots+a_n(i)f_1(i)]+C_2[f_2^{(n)}(i)+a_1(i)f_2^{(n-1)}(i)+\dots+a_n(i)f_2(i)]$. It follows that the above expression is identically equal to zero. Therefore, we can conclude that $C_1f_1(i)+C_2f_2(i)$ is a solution of (3)′.

Then we prove the latter part of the theorem. Substituting $\bar{f}(i)=C_1 \bar{\psi}_1(i)+C_2 \bar{\psi}_2(i)$ in the right side of (4), we have $A(i)[C_1 \bar{\psi}_1(i)+C_2 \bar{\psi}_2(i)]=C_1A(i) \bar{\psi}_1(i)+C_2A(i) \bar{\psi}_2(i)$. So we can conclude that $C_1 \bar{\psi}_1(i)+C_2 \bar{\psi}_2(i)$ is a solution vector of (4).

Theorem 2 (Principle 2 of Rough Superposition). If both $f_1(i)$ and $f_2(i)$ are solutions of (3)′′, then $f_1(i)+f_2(i)$ is a solution of (3)′. If $\bar{\psi}_1(i)$ and $\bar{\psi}_2(i)$ are solution vectors of (5), then $\bar{\psi}_1(i)-\bar{\psi}_2(i)$ is a solution vector of (4).

Theorem 3 (Principle 3 of Rough Superposition). If $f_1(i)$ is a solution of (3)′, $f_2(i)$ is a solution of (3)′′, $f_1(i)+f_2(i)$ is a solution of (3)′′. If $\bar{\psi}_1(i)$ is a solution vector of (4), $\bar{\psi}_2(i)$ is a solution vector of (5), $\bar{\psi}_1(i)+\bar{\psi}_2(i)$ is a solution vector of (5).

Theorem 4 (Principle 4 of Rough Superposition). Assume $g(i)=g_1(i)+g_2(i)+\dots+g_p(i)$ and $f_q(i)$ ($p \in Z^+, q=1, 2, \dots, p$) satisfy

$$f_q^{(n)}(i) + a_1(i)f_q^{(n-1)}(i) + \dots + a_{n-1}(i)f_q'(i) + a_n(i)f_q(i) = g_q(i) \tag{6}$$

where, at least one of $g_q(i)$ is not identically equal to 0, then $f_1(i)+f_2(i)+\dots+f_p(i)$ satisfies (3)′′. And assume $\bar{g}(i) = \bar{g}_1(i) + \bar{g}_2(i) + \dots + \bar{g}_p(i)$ and $\bar{\psi}_q(i)$ satisfy

$$\bar{f}_q'(i) = A(i)\bar{f}_q(i) + \bar{g}_q(i) \tag{7}$$

where, at least one of $\bar{g}_q(i)$ is not identically equal to $\bar{0}$, then $\bar{\psi}_1(i) + \bar{\psi}_2(i) + \dots + \bar{\psi}_p(i)$ satisfies (5).

The proofs of theorem 2~4 are analogous with that of theorem 1, omitted.

4 Solving Methods of Rough Differential Equations

In [8], Pawlak gave two methods of solving rough differential equations which are (1) solving the analytic expression of the solution; (2) obtaining every result of the solution by the recursive expression of relation. The foundation of the two methods is the proposition as follows.

Proposition 1. [8] $\int_0^i f'(j)\Delta(j) = f(i)+k$, where k is an integer constant.

The applications of this proposition have the following insufficiencies yet, which are (1) in proposition 1, it is only pointed out that k is an integer constant. In fact, k can be determined and is solvable. It can be easily checked that $k=-f(0)$ by definitions of rough integrals and rough derivatives. (2) Methods and formulas, which can be applied directly to obtaining solutions of rough differential equations in analytic forms, are not given by [8]. In the light of this point, this paper will discuss ordinary common rough differential equations according to their types. Solving methods and general forms of solutions will be given correspondingly. (3) The recursive formula of

solving rough differential equations, $f(i+1)=f(i)+f'(i)$, is just the defining expression of the first order rough derivations. It has no direct relationship with proposition 1. Moreover, the initial condition of the recursive relation expression given by [8] is $f(0)=k$. The k among it is not the integer constant k in proposition 1.

In general, the overwhelming majority of the rough differential equations are linear rough differential equations with constant coefficients, which can be solved accurately in analytic forms. In addition, what is given is the expression of $f(i)$ depending on i and the initial condition $f(0)$. In the following, several kinds of common rough differential equations are discussed according to their types.

4.1 First-Order Rough Differential Equations with Reducible Order

Definition 5. The first-order rough differential equation with reducible order is defined as

$$f'(i) = Q(i) \tag{8}$$

where $Q(i)$ is a rough function in $[n]$.

The left side of the equation is the first-order rough derivative of the unknown rough function $f(i)$. The right is expression of the independent variable, and the unknown rough function is not obviously included.

Theorem 5 $f(i) = \int_0^i f'(j)\Delta(j) + f(0)$.

Proof. From the definition of rough integrals, $\int_0^i f'(j)\Delta(j) = \sum_{j=0}^{i-1} f'(j)\Delta(j)$.

Therefore, $\int_0^i f'(j)\Delta(j) = f(1)-f(0)+f(2)-f(1)+\dots+f(i)-f(i-1)=f(i)-f(0)$. It follows that

$$f(i) = \int_0^i f'(j)\Delta(j) + f(0)$$

Theorem 5 gives the formula of directly solving the analytic solution of (8). This method is called the formula method. The applications of the formula method need combining computing methods of rough integrals and fundamental formulas of rough integrals. The Cauchy solution of (8) can be immediately obtained by the value of the initial condition $f(0)$.

Example 1. Solve $f'(i)=il^i+2$, where l is an integer constant and $l \neq 1$.

According to theorem 5, $f(i) = \int_0^i f'(j)\Delta(j) + f(0)$. From $f'(i) = il^i+2$, we have $f(i) =$

$$\int_0^i jl^j\Delta(j) + 2 \int_0^i \Delta(j) + f(0)$$

By the fundamental formula of rough integrals $\int_a^b k^i\Delta(i) = [\frac{1}{k-1}k^i]_a^b, (k \neq 1)$, $\int_0^i jl^j\Delta(j) = \frac{1}{l-1} \int_0^i j(l^j)' \Delta(j)$. Then by the rough integration by

$$\text{parts } \int_a^b f(i)g'(i)\Delta(i) = [f(i)g(i)]_a^b - \int_a^b f'(i)g(i)\Delta(i) - \int_a^b f(i)g'(i)\Delta(i),$$

$$\int_0^i jl^j\Delta(j) = \frac{il^i}{l-1} - \frac{l^{i+1}-l}{(l-1)^2}$$

By $\int_a^b k\Delta(i) = [ki]_a^b, \int_0^i \Delta(j) = i$. So $f(i) = \frac{il^i}{l-1} - \frac{l^{i+1}-l}{(l-1)^2} + 2i + f(0)$ is the solution sought.

Definition 6. The solution of a rough differential equation derived by the formula method is called the formula solution. The solution of a rough differential equation derived by the recursive expression $f(i+1)=f(i)+f'(i)$ is called the recursive solution.

Theorem 6. The recursive solution of (8) is $f(i)=f(0)+\sum_{j=0}^{i-1}Q(j)$.

Proof. According to $f(i+1)=f(i)+Q(i)$, when $i=0$, we have $f(1)=f(0)+Q(0)$, so the conclusion holds. Assume the conclusion holds when $i=u \in [n]$, that is $f(u)=f(0)+\sum_{j=0}^{u-1}Q(j)$. When $i=u+1$, we can conclude that $f(u+1)=f(u)+Q(u)=f(0)+\sum_{j=0}^{u-1}Q(j)+Q(u)=f(0)+\sum_{j=0}^uQ(j)$. Therefore, when $i=u+1$ the conclusion holds. By the mathematical induction, theorem 6 can be proved.

In theory, the recursive method can be used to solve all of the first-order and higher order rough differential equations, especially when rough functions in rough differential equations are presented in tabular form.

Proposition 2. The formula solution and the recursive solution of the first-order rough differential equation with reducible order (8) are equivalent.

Proof. From $f'(i)=Q(i)$ and the definition of rough integrals, $f(i)=\int_0^i f'(j)\Delta(j)+f(0)=\int_0^i Q(j)\Delta(j)+f(0)=\sum_{j=0}^{i-1}Q(j)+f(0)$. On the other hand, $f(i)=f(0)+\sum_{j=0}^{i-1}Q(j)=f(0)+\sum_{j=0}^{i-1} \int_0^i f'(j)\Delta(j)$. So the theorem is deduced.

4.2 Linear Rough Differential Equations of First Order

Definition 7. The linear rough differential equations of first order is defined as

$$f'(i)=P(i)f(i)+Q(i),$$

where both $P(i)$ and $Q(i)$ are rough functions in $[n]$.

When $Q(i)=0$, it is called a linear homogeneous rough differential equation of first order and is denoted by (9)'. When $Q(i) \neq 0$, called a linear non-homogeneous rough differential equation of first order, denoted by (9)''.

Theorem 7. The recursive solution of the first-order linear non-homogeneous rough differential equation (9)'' is $f(i)=f(0)\prod_{j=0}^{i-1} [1+P(j)]+\sum_{j=0}^{i-1} Q(j)\prod_{l=j+1}^{i-1} [1+P(l)]$.

The proof is analogous with that of theorem 6, omitted.

From theorem 7, the following corollaries can be immediately derived.

Corollary 1. The recursive solution of (9)' is $f(i)=f(0)\prod_{j=0}^{i-1} [1+P(j)]$.

Corollary 2. The Cauchy problem solution of (9)'' when $f(0)=0$ is $f(i)=\sum_{j=0}^{i-1} Q(j) \prod_{l=j+1}^{i-1} [1+P(l)]$.

Example 2. Solve $f'(i)=2i f(i)+i$.

This equation belongs to a linear non-homogeneous rough differential equation of first order, where $P(i)=2i$, $Q(i)=i$. According to theorem 7, it can be directly gained that the solution of this equation $f(i)=f(0) \prod_{j=0}^{i-1} (1+2j)+\sum_{j=0}^{i-1} j \prod_{l=j+1}^{i-1} (1+2l)$. It yields $f(1)=f(0)$, $f(2)=3f(0)+1$, $f(3)=15f(0)+7$, It follows that substituting different values of i into the above expression, the solution needed can be obtained immediately.

4.3 Linear Rough Differential Equations Group with Constant Coefficients

Solving linear rough differential equations group with constant coefficients is equal to solving rough differential equations group of (5) in which $A(i)$ is in the form of a constant matrix. General steps will be given as follows.

Step 1 Eliminate some unknown rough functions and their various order rough derivatives from the rough differential equations group. A higher order linear rough equation with constant coefficients, which containing only one unknown rough function, can be thus obtained.

Step 2 Seek the unknown rough function satisfying this higher order linear rough differential equation with constant coefficients.

Step 3 Substitute the obtained rough function into the original equations group to seek the rest unknown rough functions.

4.4 Illustrating Other Types of Rough Differential Equations

In the discussion of the above several types of rough differential equations, the formula method and the recursive method are applied to solve. For other types, some of them can be solved by being turned into the known type of rough differential equations after proper deformation according to the characteristics of the equations. The digital curvature function is illustrated as follows.

Example 3. Solve the digital curvature function satisfying the condition that its curvature is a constant, that is solve $f'(i+1)-f'(i)=C$.

Notice that $f'(i)=f(i+1)-f(i)$ and $[f(i+1)-f(i)]'=f'(i+1)-f'(i)$. The formerly equation can be rewritten as $f''(i)=C$, that is $f'(i)=Ci+C_0$ ($C_0=f'(0)$). Thus it can be turned into the type of (8). Solutions of the equation can be simultaneously obtained by the formula method and the recursive method. Since the solving of the formula method has been applied in Example 1, it will not be repeated here. By theorem 6, $f(i)=f(0)+\sum_{j=0}^{i-1} Q(j)$.

Here, $Q(j)=cj+f'(0)$. So $f(i)=f(0)+c \sum_{j=0}^{i-1} j + \sum_{j=0}^{i-1} f'(0)=f(0)+c \frac{i(i-1)}{2} +if'(0)$. Hence the solution of $f'(i+1)-f'(i)=C$ is $f(i)=\frac{C}{2}i(i-1)+f'(0)i+f(0)$.

5 Conclusion

On the basis of difference and differential equation theory, a theoretical system of rough differential equations is established. (1) Solutions in analytic form of the first-order rough differential equations with reducible order can be obtained by the formula method and recursive method. (2) For non-homogeneous equations in linear rough differential equations of first order, their analytic solutions can be obtained by the recursive method. Thus, analytic solutions of homogeneous rough differential equations can be given in a corollary form. (3) General solving steps for linear rough differential equations groups with constant coefficients are given. Finally, (4) we make some general remarks on rough differential equations besides the above typical types, and an example is illustrated.

Acknowledgments. This work was supported by NSFC (61070241), DFSC (BS2009SF021), and NSFSC (ZR2010FM035).

References

1. Pawlak, Z.: Rough sets. *International Journal of Computer and Information Sciences* 11, 341–356 (1982)
2. Yang, Y.J., Hinde, C.: A new extension of fuzzy sets using rough sets: R-fuzzy sets. *Information Sciences* 180, 354–365 (2010)
3. Wang, S.L., Zhang, L., Zhang, Q.: Numerical Investigation of Entropy Generation and Optimization on a Centrifugal Fan. *Advanced Science Letters* 4, 2240–2245 (2011)
4. Pawlak, Z.: Rough calculus. In: *Proceedings of 2nd Annual Joint Conference on Information Science*, vol. 1, pp. 344–345 (1995)
5. Shah, C.P., Singh, K., Dwivedi, C., Kumar, M., Bajaj, P.N.: Synthesis and Characterization of Sodium Selenosulphate Induced PVA-Capped Gold Nanoparticles. *Advanced Science Letters* 3, 288–294 (2010)
6. Pawlak, Z.: Rough real functions and rough controllers. In: *Proceedings of the Workshop on Rough Set and Data Mining at 23rd Annual Computer Science Conference*, pp. 58–64 (1995)
7. Pawlak, Z.: Rough functions. *Bull PAS Tech Ser.* 355, 249–251 (1997)
8. Pawlak, Z.: Rough sets, rough function and rough calculus. In: Pal, S.K., Skowron, A. (eds.) *Rough-Fuzzy Hybridization, A New Trend in Decision-making*, pp. 99–109. Springer, Singapore (1999)

Cloud Computing Infrastructure and Application Study

Ming Ye and ZeHui Qu

School of Computer and Information Science, Southwest University, 400715,
Chongqing, China
yeming@163.com, zmxym@swu.edu.cn

Abstract. With the significant advances in Information and Communication Technology and the utilization of cloud platforms grows over the last half century, users are realizing that the implicit promise of clouds (leveraging them from the tasks related with infrastructure management) is not fulfilled. This paper proposes a cloud-based infrastructure that is optimized so as to support large-scale agriculture information computing. This cloud infrastructure mainly consists of virtualization platform for agriculture information cloud computing and management. At the same time, this paper also provides insights on market-based resource management strategies that encompass both customer-drive service and management. Furthermore, the paper evaluates the performances of CPU and Internet-based services workload in the environment of proposed cloud computing platform infrastructure and management service. Experiments shows that the proposed cloud compute infrastructure and management service is effective and essential for large-scale agriculture information computing.

Keywords: Cloud computing, Virtualization computing, Management service, Cloud platform architecture.

1 Introduction

Cloud computing and distributed computing, grid computing is the same strain. Although the concept of cloud computing is emerging a few years ago, but the accumulation of technology is concerned, already up to two to three decades [1-5]. In essence, cloud computing is to a large number of distributed and high-cost computer, networking equipment and storage resource management to support a variety of applications over the Internet to provide customers with high-quality low-cost service [6-9].

Cloud computing has the following characteristics: (1) with significant scale: Google cloud computing already has more than 100 million servers, Amazon, IBM, Microsoft, Yahoo, etc. "cloud" are hundreds of thousands of servers. Enterprises typically have hundreds of thousands of private cloud servers. "Cloud" gives users an unprecedented computing power [10-13]. (2) Virtualization: Cloud computing allows users at any location, using a variety of terminal access to application services. The requested is resources from the "cloud" rather than fixed physical entity. Application in the "cloud" somewhere in the running, but in fact users do not know or worry about

the specific location of running applications. Only need a laptop or a cell phone, you can achieve through the network services to all our needs, even supercomputing to the task. (3) High reliability: "Cloud" of data using multiple copies of tolerance, computing nodes are interchangeable with the structure and other measures to ensure the service reliability, the use of cloud computing and reliable than using the local computer. (4) Versatility: Cloud computing is not for a particular application, in the "cloud" can be constructed under the support of the application of ever-changing, with a "cloud" can support different applications running simultaneously [14-17]. (5) High scalability: "Cloud" the size of dynamically scalable to meet the growing size of applications and user needs. (6) on-demand services: "Cloud" is a huge resource pool, you demand to buy; cloud can be like running water, electricity, gas as billing. (7) Extremely cheap: As the "cloud" of special measures can be used very low-cost fault-tolerant nodes to form the cloud, "cloud" of automated centralized management to a large number of companies without the burden of increasingly high cost of data center management, "cloud" of the universal utilization of resources significant than the traditional system upgrade, so users can fully enjoy the "cloud" of low-cost advantages, often spend hundreds of dollars just a few days' time will be able to complete the previously required tens of thousands of dollars, several months to complete the task. Cloud computing can completely change people's future lives, but also with attention to environmental issues, so as to truly contribute to human progress, rather than simply upgrade their technology. (8) And potentially dangerous: cloud computing services in addition to providing computing services, but also must provide a storage service. However, the current monopoly of cloud computing services in the private sector (enterprises), whereas they only able to provide commercial credit. For government agencies, commercial organizations (in particular, banks that hold sensitive data like business organizations) cloud computing service for the selection should maintain adequate vigilance. Once the business users to use the private sector to provide large-scale cloud computing services, no matter how strong its technical advantages, inevitably these private institutions to "data (information)," The importance of the hijack the entire community. For the information society, "information" is essential [17-19].

On the other hand, the cloud data for data users other than the owner of cloud computing users is confidential, but to provide the business sector in terms of cloud computing really no secret at all. It is like an ordinary person cannot monitor other people's phones, but the telecommunications company, they can always listen to any phone. All of these potential dangers are business and government organizations to choose cloud computing services [20-22].

2 Infrastructure for Cloud Computing

The underlying cloud infrastructure and environment must be designed and implemented to be flexible and scalable. Unfortunately, the history of designing, delivering, and managing very large scale federally-developed systems does not offer many success stories to build upon. If not implemented properly, the government risks significant challenges and costs in migrating information to different technologies as the third-party

vendor upgrades its processing and storage environment. If this type of upgrade is managed in-house, resident IT professionals can more readily manage migration and harmonization of data, users, and processes. But the procedures that a cloud vendor executes in scaling its environment is managed without the input of its customers, and may change services the customer requires. Customers require the ability to increase bandwidth, speed, and response time. In some cases, the cost to move data to a cloud infrastructure has proven quite costly in terms of time (bandwidth) and money. Some cloud users have resorted to using physical media to send data in order to expedite changes in their business needs. All IT systems are subject to regular considerations of their life spans and durability. The question arises as to how long certain technologies will exist. The need for interoperability, the ability to switch providers, compatibility between vendors, and avoidance of migration issues will all be demanded by customers. As the government approaches the cloud, this will likely be very problematic as there are no universal, ratified standards within the industry or through NIST that would govern these issues. So it is important for users how to design the integrated solving strategy for cloud computing platform.

In this paper, we present infrastructure and management strategy for cloud computing to meet users' demand. Cloud computing platform integrated in together by certain relatively independent subsystems; include data server, computing server, busiess intelligence application. Uers can make use of cloud computing infrastructures for their serve, for example, data storage, information sharing and remote computing and so on. See Fig. 1.

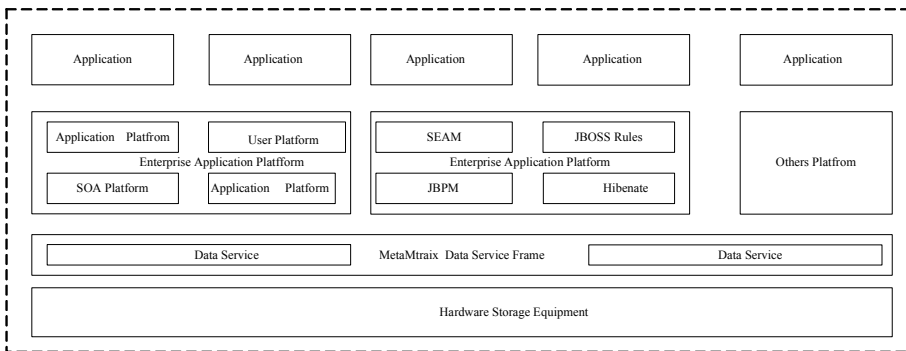


Fig. 1. Cloud Platform Architecture

3 Management Service Platform

Customer satisfaction is the crucial success factor to excel in the service industry, cloud computing service providers have to meet them in order to achieve customer satisfaction. Hence, we need to design management system platform for data centers and cloud computing that provide customers, such as enabling communication to keep customers informed and obtain feedback, monitor cloud computing platform runtime, convenience

customers to use and so on. These management strategies can also encourage trust and confidence in customers by emphasizing on the security measures undertaken against risks and doubts.

3.1 Cloud Platform Management

Firstly, management service platform furnish a overall management module. Cloud management platform can reflect working environment of cloud system, include environment of software and hardware. At the same time, cloud management platform can manage all resource of cloud computing platform and effectively assign it, include virtual machine.

Secondly, management service platform furnish a monitoring server. Fig.2 shows monitoring platform can monitor amount of total CPU, Hosts up and Hosts down. At the time, monitoring server can display the using instance of CPU, Memory and Nodes, and so on. We also can calculate average workload of CPU by this monitoring-platform.

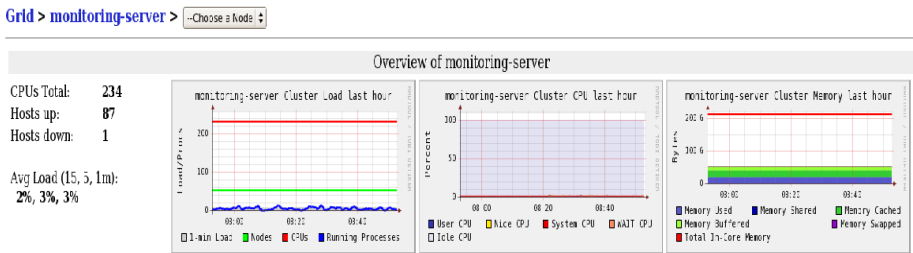


Fig. 2. Service monitoring platform

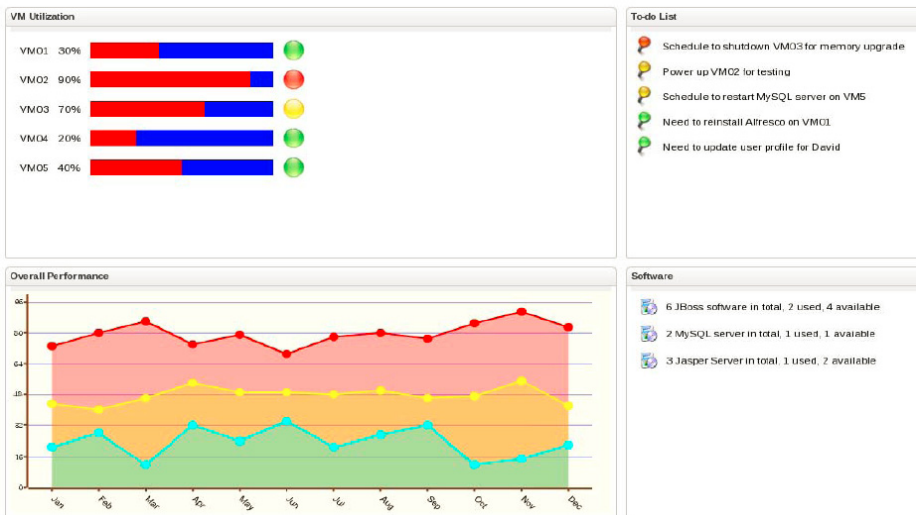


Fig. 3. Cloud platform performance monitoring

3.2 Cloud Platform Performance Monitoring

Thirdly, management service platform furnish a cloud computing center performance monitoring module. Fig.6 shows overall performance of the agriculture information cloud can be monitored. This monitoring platform can monitor overall workload of agriculture information cloud and realtimely display in graphics interface. Once system manager monitor system workload is very scale, they can immediately deal with. At the time, we can improve the throughput of the agriculture information cloud.

4 Performance Evaluation

We conduct extensive experiments on the performance of our proposed cloud Infrastructure. Fig.4shows the Agriculture information Cloud environment setup used for performance evaluation. The Agriculture information Cloud contains 30 personal computers (PCs) with 1 master node and 32 execution nodes located across 3 student computer laboratories in the Department of Computer Science and Software Engineering. This setup demonstrates that the Agriculture information Cloud is able to present a unified resource to the users/brokers by harnessing computing resources located physically apart in the 3 laboratories. Synthetic workloads are created by utilizing trace data of HPC applications. The experiments utilize 238 reservation requests.

Experiment result shows that CPU workload and internet workload of this work model of the agriculture information cloud are acceptant. But this work model in the standard distribution is unable to perform well in heterogeneous loud computing infrastructure. Though agriculture information cloud is successful in homogeneous computing environments, experimental observations reveal that the homogeneity assumptions of reduce can cause wrong and often unnecessary speculative execution in heterogeneous environments, sometimes resulting even worse performance than with speculation disabled. This valuation and performance results demonstrate that Cloud execution management systems need to be designed to handle heterogeneity that is present n workloads, applications, and computing infrastructure.

5 Conclusions

Cloud computing is a new and promising paradigm delivering IT services as computing utilities. As Clouds are designed to provide services to external users, providers need to be compensated for sharing their resources and capabilities. In this paper, we have proposed architecture for market-oriented allocation of resources within Clouds. We have also presented a management service strategy to manage cloud computing platform. Moreover, we have evaluated this cloud computing center based on this infrastructure and system arctitecture. The result proves CPU workload and internet workload are acceptant.

In particular, we have presented various Cloud efforts in practice from the market-oriented perspective to reveal its emerging potential for the creation of third-party services to enable the successful adoption of Cloud computing. In addition, we need programming environments and tools that allow rapid creation of Cloud applications. Data Centers are known to be expensive to operate and they consume huge amounts of electric power. As Clouds are emerging as next-generation data centers and aim to support ubiquitous service-oriented applications, it is important that they are designed to be energy efficient to reduce both their power bill and carbon footprint on the environment. To achieve this software systems level, we need to investigate new techniques for allocation of resources to applications depending on quality of service expectations of users and service contracts established between consumers and providers. Finally, we need to address regulatory and legal issues, which go beyond technical issues. Some of these issues are explored in related paradigms such as Grids and service-oriented computing systems. Hence, rather than competing, these past developments need to be leveraged for advancing Cloud computing. Also, Cloud computing and other related paradigms need to converge so as to produce unified and interoperable platforms for delivering IT services as the 5th utility to individuals, organizations, and corporations.

Acknowledgements. The authors thank the editors and the anonymous reviewers for their helpful comments and suggestions. This work is supported by “the Fundamental Research Funds for the Central Universities (XDJK2012C023).

References

1. Aubert, B.A., Patry, M., Rivard, S.: A Framework for information technology outsourcing risk management. *The Database for Advances in Information Systems* 36(4), 9–28 (2005)
2. Baker, S.: Google and the wisdom of the clouds, <http://www.msnbc.msn.com/id/22261846/> (retrieved February 27, 2009)
3. Bertot, J., Jaeger, P.T., Shuler, J.A., Simmons, S.N., Grimes, J.M.: Reconciling government documents and e-Government: Government information in policy, librarianship, and education. *Government Information Quarterly* 26, 433–436 (2010)
4. Brodtkin, J.: Loss of customer data spurs closure of online storage service “the link up” (2008)
5. Burroughs, J.M.: What users want: assessing government information preferences to drive information services. *Government Information Quarterly* 26, 203–218 (2009)
6. Jaeger, P.T., Lin, J., Grimes, J.M.: Cloud computing and information policy: computing in the policy cloud. *Journal of Information Technology & Politics*, 269–283 (2008)
7. Jaeger, P.T., Lin, J., Grimes, J.M., Simmons, S.N.: Where is the cloud Geography, economics, environment, and jurisdiction in cloud computing. *First Monday* 14(5) (2009)
8. Jaeger, P.T., McClure, C.R., Bertot, J.: The e-rate program and libraries and library consortia. *Information Technology and Libraries* 24(2), 57–67 (2005)
9. Wen, G.Y., Marshak, A., Cahalan, R.F., et al.: 3-D aerosol-cloud radiative interaction observed in collocated MODIS and ASTER images of cumulus cloud fields. *Journal of Geophysical Research-Atmospheres* 112(D13) (2007)

10. Grossman, R.L., Gu, Y.H., Sabala, M., Zhang, W.Z.: Compute and storage clouds using wide area high performance networks. *Future Generation Computer Systems-The International Journal of Grid Computing Theory Methods and Application* 25(2), 179–183 (2009)
11. Buyya, R., Abramson, D., Venugopal, S.: The grid economy. *Proceedings of the IEEE* 93(3), 698–714 (2005)
12. Stuer, G., Vanmechelena, K., Broeckhovea, J.: A commodity market algorithm for pricing substitutable grid resources. *Future Generation Computer Systems* 23(5), 688–701 (2007)
13. Venugopal, S., Chu, X., Buyya, R.: A negotiation mechanism for advance resource reservation using the alternate offers protocol. In: *Proc. 16th Int. Workshop on Quality of Service, IWQoS 2008*, Twente, The Netherlands (June 2008)
14. Van Looy, B., Gemmel, P., Van Dierdonck, R. (eds.) *Services Management: An Integrated Approach*: Financial Times. Prentice Hall, Harlow (2003)
15. Schneider, B., White, S.S.: *Service Quality: Research Perspectives*. Sage Publications, Thousand Oaks (2004)
16. Yeo, C.S., Buyya, R.: Integrated risk analysis for a commercial computing service. In: *Proc. 21st IEEE Int. Parallel and Distributed Processing Symposium, IPDPS 2007*, Long Beach, USA (March 2007)
17. Crouhy, M., Galai, D., Mark, R.: *The Essentials of Risk Management*. McGraw-Hill, New York (2006)
18. Moeller, R.R.: *COSO Enterprise Risk Management: Understanding the New Integrated ERM Framework*. John Wiley and Sons, Hoboken (2007)
19. Barham, P., Dragovic, B., Fraser, K., Hand, S., Harris, T., Ho, A., Neugebauer, R., Pratt, I., Warfield, A.: Xen and the art of virtualization. In: *Proc. 19th ACM Symposium on Operating Systems Principles, SOSP 2003*, Bolton Landing, USA (October 2003); *Integrated Solving Strategy for Cloud Computing. Applied Mechanics and Materials* 44-47 (2011)
20. Ye, M., et al.: Integrated real-time scheduling strategy based on Small-scale wireless sensor networks. *Sensor Letters* 9, 1–5 (2011)
21. Ye, M.: Novel Proctocl Model Design of Wireless Sensor Real-time Control Network. *Key Engineering Materials*, 460–461 (2011)
22. Ye, M., et al.: *Integrated Solving Strategy for Cloud Computing. Applied Mechanics and Materials* 44-47, 3299–3303 (2011)

An Ant Colony Algorithm for Solving the Sky Luminance Model Parameters

Ping Guo¹, Lin Zhu¹, Zhujin Liu¹, and Ying He²

¹ School of Computer Science, Chongqing University, Chongqing, 400044, China

² College of Architecture and Urban Planning, Chongqing University, Chongqing, China
guoping@cqu.edu.cn

Abstract. The new concept of sky luminance distributions which is modeling skies under a wide range of occurrences from the overcast sky to cloudless situations without or with sunlight respectively was proposed by CIE. The numerical expressions of this concept contain five adjustable parameters (a , b , c , d , e). Each type of sky proposed by CIE represents one combination of the parameters. In this paper, according to the research on the characteristics of the numerical expressions, for a measured sky type, a heuristic algorithm for solving complex optimization problems —ant colony optimization will be used to analyze and optimize the influencing factors of the sky luminance and finally get its parameters value, the experiment results show that it has high accuracy and good effect.

Keywords: Sky luminance distribution, Swarm Intelligence, Ant Colony Optimization, CIE.

1 Introduction

The first non uniform CIE standard for the luminance distribution on the overcast sky was suggested by Moon and Spencer (1942) [1], the changes of luminance from horizon to zenith in ratio 1:3. The luminance distribution on the clear sky was derived by Kittler (1967) and together with the CIE overcast sky published as ISO/CIE standard in 1996. In addition, in 1983, the International Daylight Measurement Programme (IDMP) was set up. On the basis of the programme, fifteen sky types of relative luminance distributions in the SSLD model by Kittler et al. (1998) are based on scan measured luminance data at Tokyo, Berkeley and Sydney and were proposed at the same time. Five overcast, five clear and five transitional skies are modeled by the combination of gradation and indicatrix functions [2]. This solution was proposed as a CIE code draft CIE (2001) which is under review by CIE National Committees at present.

The fifteen standard relative luminance distributions which are based on six groups of a and b values for the gradation function and six groups of c , d and e values for the indicatrix function [7]. Supposed that we have real sky measurements, how can we obtain those five coefficients value? It is difficult for traditional mathematical analysis methods to solve this problem. However, the ant colony algorithm is suitable for solving complex optimization problems.

Ant colony optimization (ACO) algorithm is a novel random global searching algorithm firstly proposed by Italian scholars Dorigo M [11-13]. It is inspired by the ant foraging behavior. Until now, ant colony optimization algorithm has being successfully applied to the traveling salesman problem [3-4], quadratic assignment problem [5], vehicle routing problem and such combinatorial optimization problems [14-15]. The research in recent years shows that it has powerful advantage to solve discrete space optimization problems [6-8]. As we know, it is quite difference from research on continuous space optimization problems. Here we reference the thought of ACO in discrete space and present a new ant colony optimization algorithm which can solve common function optimization problems.

In this paper, ant colony algorithm will be used to get the coefficients value of sky relative luminance distribution numerical expressions. It is organized as follows. In section 2, we present the sky relative luminance distribution numerical expressions. The solution method to obtain the coefficients value is presented in section 3. The experiments are shown in section 4. Finally, the conclusions are summarized in the last section 5.

2 CIE Sky Relative Luminance Distribution Model

The position of the sun and of the arbitrary sky element as well as parameters a, b, c, d, e which describe atmospheric conditions have to be taken as input calculation quantities. The position of the arbitrary sky element is defined by the zenith angle Z and the azimuth difference A_z between the element and the solar meridian, then its distance from the sun is defined by equation (1).

$$\chi = \arccos(\cos Z_s \cos Z + \sin Z_s \sin Z \cos A_z) \tag{1}$$

Where $A_z = |\alpha - \alpha_s|$. α and α_s are azimuthal angles of the vertical plane of the sky element and sun position respectively.

According to the above definition, the ratio of zenith luminance L_z to diffuse horizontal illuminance D_v , is expressed in an functional formula following the current CIE Standard :

$$\frac{L_z}{D_v} = \frac{\varphi(0^0)f(Z_s)}{\int_{Z=0}^{\pi/2} \int_{\alpha=0}^{2\pi} [\varphi(Z)f(\chi) \sin Z \cos Z] dZ d\alpha} \tag{2}$$

The luminance gradation function φ relates the luminance of a sky element to its zenith angle:

$$\varphi(Z) = \begin{cases} 1 + a \exp(b / \cos Z) & 0 \leq Z < \pi / 2 \\ 1 & Z = \pi / 2 \end{cases} \tag{3}$$

Equation (3) applies also to its value at the zenith:

$$\varphi(0^0) = 1 + a \exp b \tag{4}$$

The function f expresses the scattering indicatrix which relates the relative luminance of a sky element to its angular distance from the sun:

$$f(\chi) = 1 + c(\exp(d\chi) - \exp(d\pi / 2)) + e \cos^2 \chi \tag{5}$$

Its value at the zenith is expressed in eq. (6):

$$f(Z_s) = 1 + c(\exp(dZ_s) - \exp(d\pi / 2)) + e \cos^2 Z_s \tag{6}$$

In equation (2), L_z , D_v , Z_s , α_s can be measured by instrument. So the left side of equation (2) is a constant, and the right side is formed by the gradation function φ which contains two parameters a , b and indicatrix function which contains three parameters c , d , e . In order to obtain feasible solutions about equation (2), defined $g(a,b,c,d,e)$ like that:

$$g(a,b,c,d,e) = \left| \frac{L_z}{D_v} - \frac{\varphi(0^0)f(Z_s)}{\int_{Z=0}^{\pi/2} \int_{\alpha=0}^{2\pi} [\varphi(Z)f(\chi) \sin Z \cos Z] dZ d\alpha} \right| \tag{7}$$

So the problem becomes a function optimization problem that to obtain a solution $(a_0, b_0, c_0, d_0, e_0)$ to meet

$$g(a_0, b_0, c_0, d_0, e_0) = \min\{g(a, b, c, d, e)\} \tag{8}$$

3 Solving Methods

Eq. (8) is a function optimization problem; there are several methods to solve it, such as Simulated Annealing Algorithm, Neural Network Algorithm, Tabu Search Algorithm and so on. In this section, ant colony algorithm will be used to solve it.

3.1 Basic Ant Colony Algorithm

Taking TSP problem for example, we explain the structure of basic ant colony algorithm. Suppose that the number of the city is n . Initially, each ant is assigned to a city randomly and constructs a valid tour. d_{ij} ($i, j=1,2,\dots,n$) can be interpreted as the distance between city i and j . $\tau_{ij}(t)$ ($i, j=1,2,\dots,n$) is the residual pheromone on the line between city i and j at t moment. At the beginning ($t=0$), the pheromone on each line between city i and j is equal, $\tau_{ij}(0)=C$.

The transfer direction of ant k ($k=1,2, \dots,m$) is based on the pheromone on each trail. To keep track of the cities already visited, every ant maintains a tabu list $tabu_k$ ($k=1,2, \dots,m$), in which its actual partial tour is stored. In search process, ants depend on the transition probability $p_{ij}^k(t)$ to choose next selectable city. At moment t the definition of $p_{ij}^k(t)$ is:

$$P_{ij}^k = \begin{cases} \frac{[\tau_{ij}(t)]^\alpha [\eta_{ij}(t)]^\beta}{\sum_{s \in allowed_k} [\tau_{is}(t)]^\alpha [\eta_{is}(t)]^\beta} & j \in allowed_k \\ 0 & Otherwise \end{cases} \tag{9}$$

Where $allowed_k$ is the selectable city of ant k in the next step which is defined as $\{City_tabu_k\}$, $\eta_{ij}=1/d_{ij}$ is the inverse of the distance d_{ij} ($i, j=1,2,\dots,n$), Furthermore, α and β are positive parameters, whose values determine the relation between pheromone information and heuristic information.

In Eq. (9), we defined $\eta_{ij}=1/d_{ij}$. In this way we favor the choice of edges which are shorter and which have a greater amount of pheromone.

In ant system, the global updating rule is implemented as follows. Once all ants have built their tours, pheromone will be updated on all edges according to:

$$\tau_{ij}(t+n) = (1-\rho)\tau_{ij}(t) + \Delta\tau_{ij}(t) \tag{10}$$

$$\Delta\tau_{ij}(t) = \sum_{k=1}^m \Delta\tau_{ij}^k(t) \tag{11}$$

Where $\rho \in (0, 1]$ is the evaporation rate. Pheromone evaporation is needed to avoid too rapid convergence of the algorithm. It implements a useful form of forgetting, favoring the exploration of new areas in the search space. $\Delta\tau_{ij}(t)$ is the increment pheromone increment on trail (i, j) in this cycle. While $\Delta\tau_{ij}^k(t)$ is the pheromone increment of the k -th ant on trail(i, j) at the moment t .

Ant Circle System model is usually used the calculation model as (12). The Q is the pheromone which ants carry, and L_k is the length of the path ant k had passed.

$$\Delta_{ij}^k(t) = \begin{cases} \frac{Q}{L_k} & \text{if } k\text{th ant passed trail } (i, j) \\ 0 & \text{otherwise} \end{cases} \tag{12}$$

3.2 Improved Ant Colony Algorithm for Continuous Function Optimization

Inspired by the ACO for solving TSP problem, we should find which path can let the objective function minimization for continuous function optimization.

Division of the Solution Space. Obviously, in the TSP problem, the ants have a definite direction in the route selection, while in continuous function optimization problems, it can not offer an alternative route and direction for ants directly. Because there are countless options on continuous space, in order to solve continuous space optimization problems, one of feasible and commonly used method is dividing solution space into several discrete intervals called nodes. Supposed that the decision variables of the objective function is an n -dimensional vector $X=(x_1, x_2, \dots, x_n)$, $\min_i \leq x_i \leq \max_i, i=1,2, \dots, n$. So if x_i is divided into n discrete intervals, as the Fig. 1 shows, x_i is also called layer x_i , the function optimization problems will become an

n -dimensional decision problem, the solution space is n^n large. Obviously, the more accurate solution want to obtain the larger solution space will be. Just as Fig. 1 shows, an ant has chosen a node on each layer, the path is built, when we choose a random value of each node of the path, the solution $(a_i, b_i, c_i, d_i, e_i)$ of this ant will be obtained.

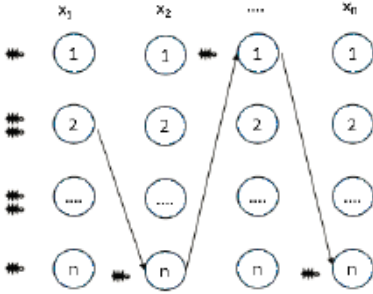


Fig. 1. Division of solution space

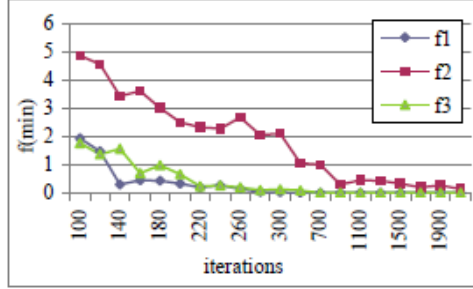


Fig. 2. The curve of f_1, f_2, f_3

The Structure of Improved Algorithm. Prior to this, continuous space had divided into several nodes. Now we explain the structure of the improved ant colony algorithm. Suppose that the number of ant is m and n is the number of nodes of each layer. Initially, each ant will be assigned to one node of first layer x_j ($Locate(ant^k)=rand()\%n$) and begin to find a feasible solution. $\tau_{x_nij}(t)$ is the residual pheromone on the line between layer x_n node i and layer x_{n+1} node j at t moment. At the beginning ($t=0$), $\tau_{x_nij}(0) = C$. Ants depend on the transition probability $p_{x_nij}^k(t)$ to choose the node on layer x_{n+1} . At moment t the definition of $p_{x_nij}^k(t)$ is :

$$p_{x_nij}^k(t) = \begin{cases} \arg \max \{ \tau_{x_nij}(t) \} & q \leq q_0, i, j \in [1, m] \\ J & otherwise \end{cases} \tag{13}$$

$$J = \begin{cases} \tau_{x_nij}(t) / \sum_{s=1}^m \tau_{x_nis}(t) & j, s \in [1, m] \\ 0 & otherwise \end{cases} \tag{14}$$

Where $q \in [0,1]$ is a random number which is used to determine the probability of random selection, constant $q_0 \in [0,1]$ which is usually set to 0.8.

In the improved ant colony algorithm, the update of the pheromone divided into two steps. First is local updating rule. In ants establishing the paths, the residual information should be weakened on the passed paths constantly by eq. (13). So the probability of the next ant choosing the same path can be reduced, expect it has already be determined the best path by many times of circulation.

$$\tau_{x_nij}(t+1) = (1-\rho)\tau_{x_nij}(t) + \rho\Delta\tau_0 \tag{15}$$

Where $\Delta\tau_0 = \Delta\tau_0^1 + \Delta\tau_0^2 + \dots + \Delta\tau_0^m$ is the residual pheromone after ants through the path. And $\rho \in (0, 1]$ is the evaporation rate.

Another is global updating rules. When the ant has chosen a node on each layer, the path is built as Fig. 1 shows. For this path we can obtain a random value of each node of the path $(a_i, b_i, c_i, d_i, e_i)$ and use to compute the function value. After all of the ants have finished its' tour construction, we can choose an ant whose function value is minimum. This ant is called iteration optimal ant. If it is smaller than the function value of the global optimal ant, the previous global optimal ant will be replaced. When all of the ants finished path construction, the nodes of the iteration optimal ant passed will be updated. This updating rule is for the optimal path in order to increase algorithm convergence speed.

$$\tau_{x_nij}(t+n) = \tau_{x_nij}(t) + \Delta\tau \tag{16}$$

$$\Delta\tau = 1 / (g_{min} + c) \tag{17}$$

Where g_{min} is the min value of the function in this iteration, C is an constant, the reason why need to add C is to prevent g_{min} is a negative number or a number close to 0. General C is always a big positive number.

Algorithm Description and Test. The improved ant colony algorithm for solving continuous function optimization problems in this paper is described as follows:

Initialization. Set the initial value of each parameter. Discretizing the continuous domain into several regions, set every ant at layer x_l . Set $\tau_{x_nij}(0) = C$.

Ants state transition. Depend on $p_{x_nij}^k(t)$ ants choose a node at x_n .

Update local pheromone according to formula.15

According to the path each ant obtained, compute the objective function value.

Update global pheromone according to formula.16. Ants return to the layer x_l .

Iterative loop. Repeat steps (2) to (5) until the number of iterations achieve maximum n_0 . Save the global optimal path and cut the other nodes and repeat steps (1) to (5) until the width of the node smaller than ε .

To prove the effectiveness of the algorithm introduced in the paper, we use it to solve three continuous function optimization problems $f_1[9]$, $f_2[9]$ and $f_3[10]$ as follows:

$$f_1 = -20 \exp(-0.2 \sqrt{1/n \sum_{i=1}^n x_i^2}) - \exp(1/n \sum_{i=1}^n \cos(2\pi x_i)) + 20 + e$$

Where, $x_i \in [-32, 32]^n$, min point: $x_i=0$, and $\min(f_1)=0$.

$$f_2 = \sum_{i=1}^n (x_i^2 - 10 \cos(2\pi x_i) + 10)$$

Where, $x_i \in [-5.12, 5.12]^n$, min point: $x_i=0$, and $\min(f_2)=0$.

$$f_3 = \sum_{i=1}^n x_i^2 + (\sum_{i=1}^n \frac{ix_i}{2})^2 + (\sum_{i=1}^n \frac{ix_i}{2})^4$$

Where, $x_i \in [-5, 10]^n$, min point: $x_i=0$, and $\min(f_3)=0$.

In order to reduce the occasional affect, we compute 200 times for each function. Set the algorithm parameters value like that: $\rho=0.7$, $\tau_0=0.01$, $q_0=0.8$, $\epsilon=0.0000001$. For f_1 , $m=8$, $n=4$, $n_0=300$, for f_2 , set $m=16$, $n=4$, $n_0=400$ and for function 3 set $m=40$, $n=8$, $n_0=1000$. Following are the solving rate of each function show on the Table 1.

Fig. 2 are the curves which shows the average value changes of f_1, f_2, f_3 , when the iterations increase, where Abscissa is iterations, ordinate is $f_{avg}(\min)$.

As the Fig. 2 shows, with the increase of the iterations, the algorithm always can obtain optimal solution.

Table 1. The Solving Rate of the each function

Function	Solving Rate	Parameters
f_1	97.5%	M=8, n=3, $n_0=300$
f_2	73.5%	M=40, n=8, $n_0=1000$
f_3	85%	M=16, n=4, $n_0=400$

4 Experimental and Results

The algorithm has proved have good effect, now we will use it to get the coefficients value of sky relative luminance distribution numerical expressions.

The change areas and trends of the coefficients which were proposed in section 2 from sunny to cloudy are as Fig. 3.

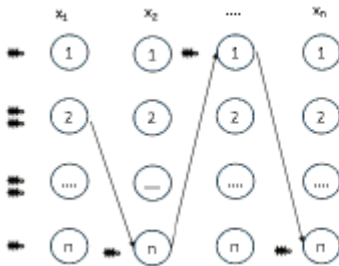


Fig. 3. The trends of the coefficients

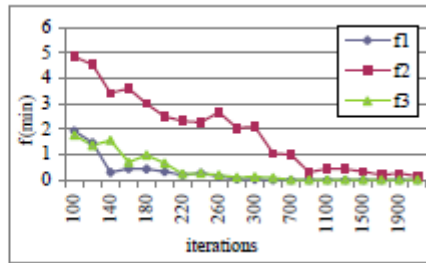


Fig. 4. The result of the coefficients

According to the test of typical overcast sky, the value of L_z, D_v, Z_s, α_s had actual measured are $\alpha_s=51.3(\text{deg})$, $Z_s=48.616(\text{deg})$, $L_z=1732(\text{cd/m}^2)$, D_v as=3750(lx), and Sky Type is cloudy.

Basing on this, use the improved ant colony algorithm to solve the problem formula.(8), Set the algorithm parameters value like that : $\rho=0.7$, $\tau_0=0.01$, $q_0=0.8$, $\epsilon=0.0000001$, $m=8$, $n=4$, $n_0=300$. Following are the result after solve 200 times.

As the Fig. 4 shows, the solutions (a, b, c, d, e) are fluctuating around (4, -0.7, 0, -1,0) which is the standard value that CIE proposed for cloudy. It shows that the algorithm have high accuracy and good effect to solve CIE sky relative luminance distribution.

5 Conclusion

In this paper, we studied the sky luminance distribution model and its numerical expression and according to the characteristics of the numerical expressions, we use an improved ant colony optimization to analyze and optimize the influencing factors of the sky luminance. For every measured sky type, the algorithm can get the coefficients value of sky relative luminance distribution numerical expressions. The experimental results show that the algorithm has good efficiency and accuracy.

Acknowledgments. This work was supported by the National Natural Science Foundation of China-Youth Fund (Grant No. 1010200220090070).

References

1. Moon, P., Spencer, D.E.: Illumination from a non-uniform sky. *Illum. Eng.* 37, 707–726 (1942)
2. Darula, S., Kittler, R.: A catalogue of fifteen sky luminance patterns between the CIE standard skies. In: Proc. 24th of the CIE Session, vol. 1, part 2, pp. 7–9. CIE Publ. 133, Warsaw (1999)
3. Dorigo, M., Gambardella, L.M.: Ant Colony System: A Cooperative Learning Approach to the Traveling Salesman Problem. *IEEE Trans. on Evolutionary Computation* 1(1), 53–66 (1997)
4. Xing, J.Q., Zhu, Q.S., Guo, P.: Research on ant colony clustering combination method. *Computer Engineering and Application* 45, 146–148 (2009)
5. Gambardella, L.M., Taillard, E.D., Dorigo, M.: Ant Colonies for the Quadratic Assignment Problem. *J. Oper. Res. Soci.* 50, 167–176 (1999)
6. Zhou, J.X., Yang, W.D., Li, Q.: Improved Ant Colony Algorithm and Simulation for Continuous Function Optimization. *Journal of System Simulation* 21, 1685–1688 (2009)
7. Wittkopf, S.K.: Analysing sky luminance scans and predicting frequent sky patterns in Singapore. *Lighting Res. Technol.* 39, 31–51 (2007)
8. Dorigo, M., Caro, G.D., Gambardella, L.M.: Ant algorithms for discrete optimization. *Artificial Life* 5, 137–172 (1999)
9. Xiao, J., Li, L.P.: A hybrid ant colony optimization for continuous domains. *Expert Systems with Applications* 38, 11072–11077 (2011)
10. Socha, K., Dorigo, M.: Ant colony optimization for continuous domains. *European Journal of Operational Research* 185, 1155–1173 (2008)
11. Yang, L., Fu, Z.Q., Wang, D., Li, H.L., Xia, J.B.: An improved ant colony algorithm for continuous space optimization. In: Proceedings of the Ninth International Conference on Machine Learning and Cybernetics, pp. 1829–1843. IEEE Press, Qingdao (2010)
12. Wang, J., Xiao, Q., Zhang, J.: Improved ant colony optimization for solving constrained continuous function optimization problems. *J. Computer Engineering and Design* 31, 1027–1031 (2010)
13. Zhao, H.Y., Li, G.C., Cui, J.: Continuous and colony optimization based on normal distribution model of pheromone. *Computer Engineering and Applications* 46, 54–56 (2010)
14. Sun, H.Y., Chen, L.: A new approach for solving continuous optimization using ant colony optimization. *Journal of ShanDong University* 39, 24–30 (2009)

Tugboat Scheduling Problem Based on Trust-Based Ant Colony Optimization

Su Wang, Min Zhu, Jun Zheng, and Kai Zheng

Computer Center, East China Normal University, North Zhongshan Road No.3663, Shanghai,
200062, China

{swang, mzhu, jzheng, kzhen} @cc.ecnu.edu.cn

Abstract. Tugboat scheduling is an important decision problem in container terminals. This paper proposes a mathematic model of tugboat scheduling and applies an improved Trust-based Ant Colony Optimization method to get the best scheduling plan. The concept of trust value is considered as heuristic information to affect path selecting and the pheromone influence information is introduced to avoid the local optima. The results of the simulation experiment suggests that Trust-based ACO is well suited for tugboat scheduling problem in container terminals and can get better performance than basic ACO.

Keywords: Ant Colony Optimization, Tugboat Scheduling, Trust Value, Container Terminal.

1 Introduction

Nowadays, container terminals are continuously facing more and more challenges and rising competition, it is necessary to reduce the turnaround time of ships which are mainly determined by the efficiency of equipment in container terminal. Tugboat is one kind of important equipment because the port lane is narrow and the water is shallow so that ships can not be sailed directly to berths and should dock with the help of tugboat towing. Efficient tugboat scheduling operation can reduce the turnaround time of ships and improve the utilization of tugboat. It is an important decision problem because the number and the service ability of tugboats are finite.

Several studies have been conducted to improve the efficiency of ship operations in container terminals including: berth and quay crane scheduling [1]-[3], yard crane scheduling [4]-[6], AGV scheduling [7], storage space allocation [8]. But there are few researches on tugboat scheduling in container terminals. Takayuki [9] introduced tugboat development situation and analyzed tugboat business in Asian countries, especially in Japan. But it did not introduce how to assign and schedule tugboats to finish ship docking. Liu [10] studied the tugboat operation scheduling problem and employed a hybrid evolutionary algorithm to solve this problem. Su Wang [11] studied the tugboat scheduling problem based on genetic algorithm and ant colony optimization, but only basic ACO algorithm was applied and not discussed deeply. This paper attempts to give the model of tugboat scheduling problem and present the Trust-based ACO to resolve the tugboat scheduling problem.

2 The Model of Tugboat Scheduling Problem

When a ship arrives at a container terminal, it usually needs the specific equipment called tugboat to help the ship berth. Moreover, the moving between two berths and the department of ships also need tugboats. Tugboat is a kind of small ship but with large engine. Tugboat's ability is stated by engine's horsepower (ps). The tugboat with higher horsepower has higher service ability and can tug bigger ship. Because the number and the engine of tugboat are limited, it is important to schedule tugboats at an optimum level for reducing scheduling cost and improving tugboat utilization.

The tugboat scheduling problem is to consider a matching relationship of multiple tugboats and ships. Tugboats are assigned to ships based on some scheduling rules. Each ship needs to schedule one or more tugboats for docking. But the tugboat can only work for one ship at the same time. The scheduling rules are shown as Table 1.

Table 1. Scheduling rules between ships and tugboats

Length of Ship (m)	Scheduled Horsepower of Tugboat	Scheduled Number of Tugboat
0-100	≥2600 ps	≥1
100-200	≥5200 ps	≥2
200-250	≥6400 ps	≥2
250-300	≥6800 ps	≥2
≥300	≥8000 ps	≥2

The goal is to determine a match between ships and tugboats so that the total scheduling cost is minimized on the premise of all of the ships finish berthing.

This paper describes tugboat scheduling problem with the following notations:

n : number of ships.

t_i : the least need horsepower of ship i .

m : number of tugboats.

c_j : horsepower of tugboat j .

x_{ij} : decision variable represents that tugboat j is to be scheduled by ship i or not.

Minimize:

$$\sum_{i=1}^n \sum_{j=1}^m x_{ij} c_j \tag{1}$$

The constraints of the problem are shown below:

$$x_{ij} = \begin{cases} 1 & \text{Tugboat } j \text{ is scheduled by ship } i \\ 0 & \text{otherwise} \end{cases} \tag{2}$$

$$\sum_{j=1}^m x_{ij} c_j \geq t_i \tag{3}$$

$$\sum_{i=1}^n \sum_{j=1}^m x_{ij} \leq m \tag{4}$$

Constraint (2) states that tugboat j is scheduled by ship i or not. Constraint (3) means that the total horsepower of the scheduled tugboats must meet the least need of ship. Constraint (4) describes that the number of the scheduled tugboat should be fewer than the total number of tugboats.

3 ACO Method for Tugboat Scheduling Problem

3.1 Ant Colony Optimization

ACO was proposed by Dorigo.M as a multi-agent approach to difficult combinatorial optimization problems [12]. The basic idea of ACO is inspired by the way ants explore their environment in search of a food source. However, basic ACO algorithm has some weaknesses such as low convergence speed and local optima. It is required to improve ACO algorithm to enhance the global search ability, convergence speed and adjust the algorithm to resolve tugboat scheduling problem in container terminals.

For tugboat scheduling problem, a ship is represented by an ant and a tugboat is represented by a note. Ship i schedules tugboats means that ant i deposits pheromone on paths related to the scheduled tugboats. For example, there are ten available tugboats $(n_1, n_2, n_3, \dots, n_{10})$ scheduled by four ships. All ants start from the virtual note n_0 . In tugboat scheduling problem, ants seek route for forming some tugboat unions for ships. In Fig.1, ants establish four tugboat unions $\{R_1, R_3\}$, $\{R_2, R_4, R_8\}$, $\{R_6, R_9\}$, $\{R_7\}$ for four ships. R_5 is not scheduled. The union $\{R_2, R_4, R_8\}$ schedules tugboat R_2 , R_8 and R_4 . Ants deposits pheromone on l_{02}, l_{28}, l_{84} .

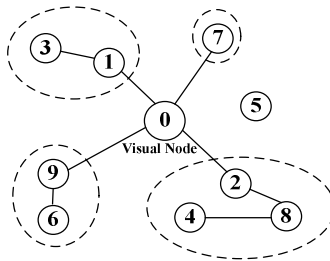


Fig. 1. Nine tugboats scheduled by four ships

The node transition rule is as follows. An ant k in R_i chooses R_j to move to following the rule:

$$s = \begin{cases} \arg \max_{j \notin tabu_k} \{ \tau_{ij}^\alpha \cdot \eta_{ij}^\beta \} & q \leq q_0 \\ S & q > q_0 \end{cases} \tag{5}$$

q is a random variable between $[0,1]$. q_0 is a constant parameter. The pheromone information is denoted by τ_{ij} and the heuristic information is denoted by η_{ij} . τ_{ij} is the

cost of scheduling tugboat R_j after tugboat R_i . With probability q_0 , the ant chooses the unscheduled tugboat which maximizes $\tau_{ij}^\alpha \cdot \eta_{ij}^\beta$. α and β determine the relative influence of the pheromone and the heuristic values. With probability $1 - q_0$, the next tugboat is chosen according to the probability S determined by p_{ij}^k .

$$p_{ij}^k(t) = \begin{cases} \frac{\tau_{ij}(t)^\alpha \cdot \eta_{ij}^\beta}{\sum_{j \in \text{tabu}_k} \tau_{ij}(t)^\alpha \cdot \eta_{ij}^\beta} & j \in \text{allow } k \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

In ACO, ants are allowed to deposit pheromone on the path. Because all tugboats in one tugboat union are combined to serve one ship, all paths in one tugboat union should be updated by the same pheromone. When one tugboat union finishes, the pheromone of paths in the union are updated as formula (7) and (8). When all tugboat unions finish, the pheromone of the best union paths is updated as formula (7) and (9).

$$\tau_{ij}^{\text{new}} = \rho \tau_{ij}^{\text{old}} + (1 - \rho) \Delta \tau_{ij} \quad (7)$$

$$\Delta \tau_{ij} = \begin{cases} \frac{Q_{\min}}{\sum M} & \text{ant } k \text{ select note } i \text{ and note } j \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

$$\Delta \tau_{ij} = \begin{cases} \frac{Q_{\min}}{\sum_k \sum M} & (i, j) \in \text{the best tour} \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

Where (i, j) is one path of the tour, ρ is a parameter governing pheromone decay. Q_{\min} is the minimal cost for ship k , $\sum M$ is the real cost of completing ship k . With the decrease of scheduling cost $\sum M$, the pheromone τ_{ij} increases so that ants have the more probability of selecting tugboat j after tugboat i .

3.2 Trust-Based ACO for Tugboat Scheduling Problem

Recently, many researchers only consider the effect of path length in ACO. But in fact, ants not only consider path length, but also path condition. For example, there are two paths shown as Fig. 2.

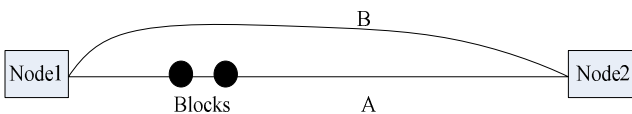


Fig. 2. Two Paths between Node1 and Node 2

Path A is shorter than path B, but with two blocks. When A blocks up, ants will take more time even if A is shorter. In that situation, ants may select the longer path B but with no block. Based on the above analysis, the trust mechanism is introduced in the Trust-based ACO and the trust value is defined as heuristic information to represent the trust degree of path. The selecting process is affected by two parameters, one is the pheromone represented by path length and the other one is the heuristic information represented by path condition. The trust value between [0,1] can be learned from the experience of ants. The path with more blocks has lower trust value. If the path has no block, trust value is 1.

In tugboat scheduling problem, not all of the tugboats can complete berthing job smoothly because of the machine depreciation and sudden faults which have effect on the scheduling operation. The new tugboat and less broken tugboat has good service ability and can compete berthing job smoothly. It means that the tugboat has higher trust value. Otherwise, the old tugboat and more broken tugboat has lower service ability and trust value.

The trust value related to the tugboat service ability is defined as the formula (10):

$$v_i = v_0 \times \frac{m}{n} \times k \quad (10)$$

v_i is the trust value of tugboat i . v_0 is the initial trust value of tugboat i . m is the number of completed job of tugboats i . n is the number of assigned job to tugboat i . k is the depreciation rate of tugboat i .

To avoid local minimum, the concept of Pheromone Influence (PI) is introduced. In ACO, ants select path according to the pheromone at the beginning. But in the real searching model, the pheromones on paths are few and have little influence in the initial stages. So ants always display the characteristic of autonomy. They select the path randomly at the beginning so that they can pass more paths and escape from the local minimum. With the pheromone accumulating, the communication among ants becomes more and more. Then ants move at the collective level and they will select the path with more pheromone. PI is a parameter to measure the level of the pheromones. In the Trust-based ACO algorithm, when the average pheromone of all the paths is fewer than PI, ants select path randomly; otherwise, ants select path by pheromone. The formula of PI is as follows:

$$PI = \frac{\sum_{i=1}^m (1/L_i)}{m} \times \lambda \quad (11)$$

n is the number of ants. (L_1, L_2, \dots, L_m) is the set of paths built by n ants randomly. λ is a integer parameter, usually 3 or 4.

Using Trust-based ACO Algorithm to solve tugboat scheduling problem in container terminals, the pseudo code is as follows:

Step 1: Initialize parameters.
 Step 2: while (available tugboat) and (uncompleted ship)
 for ($i=1; i \leq n; i++$)
 if (ship i has completed) continue;
 j = the present node of ant i ;
 if (the average pheromone of all paths from node $i < PI$)
 Ants select next paths randomly;
 else
 Ants select next paths according to pheromone by formula (5) and (6);
 Step3: All ants complete their searching process;
 Update pheromone according to formula (7) and (8);
 if (the cost of the current solution $<$ the cost of the best solution)
 Update pheromone according to formula (7) and (9);
 Update the best solution;
 Step4: $NC++$;
 if ($NC < NC_{max}$) or (no evolving solution)
 Goto Step2;
 else
 Output the best solution.

4 Simulation Optimization Results

In this section the simulation experiment of tugboat scheduling problem using the Trust-based ACO algorithm is presented. This paper uses two parameters to measure the performance of tugboat scheduling problem. One is the utilization rate of tugboat $\bar{\rho}$ and the other one is the average waiting time of ships \bar{t} . The two parameters are calculated by formula (12) and (13):

$$\bar{\rho} = \frac{\sum_{i=1}^m P_i \cdot T_i}{\sum_{i=1}^m P_i \cdot T} \quad (12)$$

P_i is the horsepower of tugboat i . T_i is the operation time of tugboat i during the whole simulation. T is the total simulation time. m is the number of tugboats.

$$\bar{t} = \frac{\sum_{j=1}^n (T_j - A_j)}{n} \quad (13)$$

T_j is the start docking time of ship j . A_j is the arriving time of ship j . n is the number of ships.

Table. 2 shows the allocation of tugboats in two different container terminals. Table.3 gives one case of arriving ships at the same time in container terminal 1.

Table 2. Distribution of tugboat at two ports

Horsepower(PS)	1200	2600	3200	3400	4000
Container Terminal 1	1	7	3	3	3
Container Terminal 2	0	7	3	3	4

Table 3. Six arriving ships at the same time

Ship No.	1	2	3	4	5	6
Ship Length (m)	145	98	223	262	187	163

In tugboat scheduling problem, the number of ants is a variable parameter which can be changed according to the number of ships. Other parameters are: $\tau_0=0.5, \alpha=1, \beta=3, \rho=0.8, NC_{max}=5000$. For the ships in Table 3, the comparison results between ACO and the Trust-based ACO are showed in Fig.3. The best solution of ACO shown in (a) is 33200ps converging at the generation 4160 and the best solution of the Trust-based ACO shown in (b) is 32800ps converging at the generation 2600. It is known that the Trust-based ACO can get better solution and faster convergence speed.

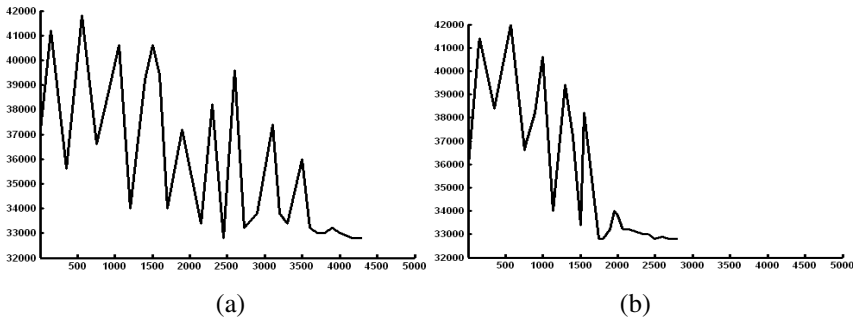


Fig. 3. Comparison of the best solution of ACO and Trust-based ACO

Table 4. Comparison of simulation optimization results

Parameters	Utilization Ratio of Tugboats (%)		Average Waiting Time of Ships (m)	
Container Terminal	1	2	1	2
Only Simulation	20.5	25.4	7.1	5.4
Basic ACO	32.3	38.6	3.2	2.4
Trust-based ACO	36.5	42.5	1.9	1.3

For the ships in 30 days, the experiment results are shown in Table.4. From the results, it is clear that the Trust-based ACO can get the higher utilization ratio of tugboats and less average waiting time of ships than basic ACO and simulation without ACO. Moreover, the performance of container terminal 1 is better than container terminal 2. Although the number of the two terminals is identical, container

terminal 2 has one more tugboat with 4000ps and it can get higher operation efficiency and serve more ships.

5 Conclusion and Future Work

Tugboat scheduling problem is an important decision problem in container terminals. It is necessary to study it using optimization method to get the best scheduling solution so that it can improve the operation efficiency of container terminal. In this paper, a mathematic model for tugboat scheduling problem based on match mechanism is defined. The improved Trust-based ACO algorithm is presented to solve this problem. In the algorithm, ants are incline to select the tugboats with less scheduling cost and better service ability. The results of simulation show that the Trust-based ACO can get better solution and performance than basic ACO.

In the future, it is decided to study the Trust-based ACO deeply to resolve other scheduling problem in container terminal such as berth scheduling problem, quay crane scheduling problem and etc. On the other hand, the heuristic learning algorithm should be introduced to obtain the trust valve.

Acknowledgments. This work is supported by the Natural Science Foundation of Shanghai (No.10ZR1410400).

References

1. Han, X.-L., Lu, Z.-Q., Xi, L.-F.: A proactive approach for simultaneous berth and quay crane scheduling problem with stochastic arrival and handling time. *European Journal of Operational Research* 207, 1327–1340 (2010)
2. Bierwirth, C., Meisel, F.: A survey of berth allocation and quay crane scheduling problems in container terminals. *European Journal of Operational Research* 202, 615–627 (2010)
3. Kim, K.H., Park, Y.-M.: A crane scheduling method for port container terminals. *European Journal of Operational Research* 156(3), 752–768 (2004)
4. Guo, X., Huang, S.Y., Hsu, W.J., Low, M.Y.H.: Dynamic yard crane dispatching in container terminals with predicted vehicle arrival information. *Advanced Engineering Informatics* 25, 472–484 (2011)
5. Chen, L., Bostel, N., Dejax, P., Cai, J., Xi, L.: A tabu search algorithm for the integrated scheduling problem of container handling systems in a maritime terminal. *European Journal of Operational Research* 181, 40–58 (2007)
6. Zeng, Q., Yang, Z.: Integrating simulation and optimization to schedule loading operations in container terminals. *Computers & Operations Research* 36, 1935–1944 (2009)
7. Barcos, L., Rodríguez, V., Álvarez, M.J., Robusté, F.: Routing design for less-than-truck load motor carriers using Ant Colony Optimization. *Transportation Research Part E: Logistics and Transportation Review* 46, 367–383 (2010)
8. Bazzazi, M., Safaei, N., Javadian, N.: A genetic algorithm to solve the storage space allocation problem in a container terminal. *Computers & Industrial Engineering* 56, 44–52 (2009)
9. Mori, T.: The present situation and the issues on tugboat business in Japan, <http://www.h2.dion.ne.jp/~t-mori/ronbun.html>

A Cloud Architecture with an Efficient Scheduling Technique

Nawsher Khan, A. Noraziah, and Tutut Herawan

Faculty of Computer Systems and Software Engineering University Malaysia Pahang
Lebuh Raya Tun Razak, 26300 Gambang, Kuantan, Pahang, Malaysia
nawsher@gmail.com, {noraziah,tutut}@ump.edu.my

Abstract. Reliability, efficiency (in term of time consumption) and effectiveness in resources utilization are the desired quality attributes of Cloud scheduling system, the main purpose of which is to execute jobs optimally, i.e. with minimum average waiting, turnaround and response time. Replication provides improved availability, decreased bandwidth use, increased fault tolerance, and improved scalability. To speed up access, file can be replicated so a user can access a nearby replica. In this paper, we propose an architecture to convert Globally One Cloud to Locally Many Clouds. By combining replication and scheduling, this architecture will improve efficiency, accessibility, reliability, availability and scalability. In the case of failure of one sub cloud or one cloud service, clients can start using another cloud under “failover” techniques. As a result, no one cloud service will go down.

Keywords: Cloud Computing, Sub-Cloud, Replication, Scheduling.

1 Introduction

Reliability, accessibility, efficiency and scalability are the key factors and needs to be the key features of cloud computing. Nowadays, computing, data storage and data transferring requirements from end-users are growing, demanding more capacity, more reliability and the capability to access information from anywhere in the world. Cloud services (computing, storage and transferring) meet this demand by providing transparent, easy and reliable solutions. Since late 2007, the concept of cloud computing was proposed [1] and it has been utilized in many areas with many achievements [2,3]. Cloud computing is deemed as the next generation of IT platforms that can deliver computing as a kind of utility [4]. Foster *et al.* made a comprehensive comparison of grid computing and cloud computing [5].

By a cloud, we mean an infrastructure that provides resources and/or services over the Internet. A storage cloud provides storage services (block or file based services); a data cloud provides data management services (record-based, column-based or object-based services); and a compute cloud provides computational services. Often these are layered (compute services over data services over storage service) to create a stack of cloud services that serves as a computing platform for developing cloud-based applications, for example; Google’s Google File System (GFS), BigTable and MapReduce infrastructure [6,7], Amazon’s S3 storage cloud, SimpleDB data cloud,

EC2 compute cloud [8]; and the open source Hadoop system [9,10]. For the majority of applications, databases are the preferred infrastructure for managing and archiving data sets, but as the size of the data set begins to grow larger than a few hundred terabytes, current databases become less competitive with more specialized solutions, such as the storage services e.g., [9,10] that are parts of data clouds. For example, Google's GFS manages Petabytes of data [11].

Cloud architectures are middleware services for different purposes i.e. resource allocation management, job scheduling, security [24], authorization and data management etc. When a user requests a file, a large amount of bandwidth could be spent to send the file from the server to the client and the delay or response time involved could be high [12]. Besides that, maintaining local copies of data on each accessing site are cost prohibitive while storing all data in a centralized manner is impractical due to remote access latency [13]. This may lead the Internet turns to be the bottleneck in accessing the files in the Cloud Computing. Due to the high latency of the Wide Area Network (WAN), the main issue is to design the strategy for efficient data access and share around the world with considerably low time complexity in data Grid and Cloud research [14]. Furthermore, in order to manage the data there are another several problems must be considered such as failures or malicious attacks during execution, fault tolerance, scalability of data and etc. These problems can be solving by using the replication techniques [15].

User of cloud computing can and have created a virtual server room on one desktop. 100 million virtual machines are being created per year or 273,972 per day or 11,375 per hour. The number of physical servers in the world today is about 50 million. By 2013, for approximately 60 percent of server workloads will be virtualized means will convert to virtual cloud [16,17]. For the better management of this day by day increasing heavy storage data, thus we need an efficient scheduling technique.

In this paper, we propose an architecture to convert Globally One Cloud to Locally Many Clouds. By combining replication and scheduling, this architecture will improve efficiency, accessibility, reliability, availability and scalability. In the case of failure of one sub cloud or one cloud service, clients can start using another cloud under "failover" techniques. As a result, no one cloud service will go down.

The rest of the paper is organized as follows. Section 2 describes related works that address the problem of scheduling. Section 3 describes scheduling architecture and expected results. Finally, the conclusion of this work is presented the Section 4.

2 Related Works

There are some recent and related works that address the problem of scheduling. Ranganathan and Foster [18] proposed the realization of importance of data locality in job scheduling problem. The authors presented a Data Grid architecture base on three main components i.e. External Scheduler (ES), Local Scheduler (LS) and Dataset Scheduler (DS). ES receives submitted jobs from user, then depend on ES's scheduler policies, it decide which job to send to which remote site. How to schedule all the

jobs, LS of each site decide on its local resources. Keeping track of popularity for each dataset currently available and making data replicating decision, this is the responsibility of Tang, et al. [19] and Nguyen, et al. [20] improved the older Ranganathan and Foster works [18] by integrating the scheduling and replication strategy to improve the scheduling performance.

Analyzing the works of Ranganathan and Foster [18], Tang, et al. [19] and Nguyen, et al. [20], the authors have proposed new scheduling architecture. For the integration of scheduling and replication, in this stage, the author has focused on Total Completion Time (TCT) for a job. In above works the authors uses the following formula for Total Time Completion (TT) for a job as given as follow

$$TT_{k,i} = \max\{QT_{(i)}, DT(f(k), i)\} + ET_{k,i} \tag{1}$$

$$ETTC_{ji} = \max\{DT(f(j), i), QT_{(i)}\} + EET_{ji} \tag{2}$$

In $\max\{DT(f(j), i), QT_{(i)}\}$, one value DT either QT is ignoring, because only one of them will be maximum, the other minimum value is ignoring, Even both are two different parameters, and have their importance separately.

In initial stage, the author separated these two parameters and mathematically proved by using Little’s law [21].

New architecture presents a new cloud model, in which the failure possibility of a service is not acceptable as previously discussed in cloud service proposed by Khan *et.al.* [22].

3 Proposed Method

3.1 Globally One Cloud to Locally Many Clouds

For the improvement of data efficiency, easy accessibility, strong reliability and always availability, the current research proposes an approach in order to overcome these qualities in cloud services. We need to make six (let say) local sub clouds on the basis of six sub-continent instead of one Global cloud. Each local sub cloud will have same anytime updated replica (copy) of each other. Due to the local cloud, data in a wisely manner will offer a faster access to files required by cloud client, hence increase the job execution’s performance.

Due to local cloud, reliability and accessibility will increase. If one sub cloud goes down, any client from anywhere in the world can use another sub cloud under the “failover” techniques as described in Figure 1.

Dealing with large amount of data makes the requirement for efficiency in data access more critical. A good scheduling strategy will allow shortest access to the required data; therefore reduce the data access time. Vice versa, replication strategy that allows place data in a wisely manner will offer a faster access to required files. The goal of replication is to shorten the data access not only for user accesses but enhancing the job execution performance. For this approach we need an architecture, in which scheduling is compatible with replication.



Fig. 1. Simple Sub-Cloud Structure for the Globe

3.2 Scheduling Architecture

The scheduling architecture is encapsulated in three distinct modules as shown in Figure 2.

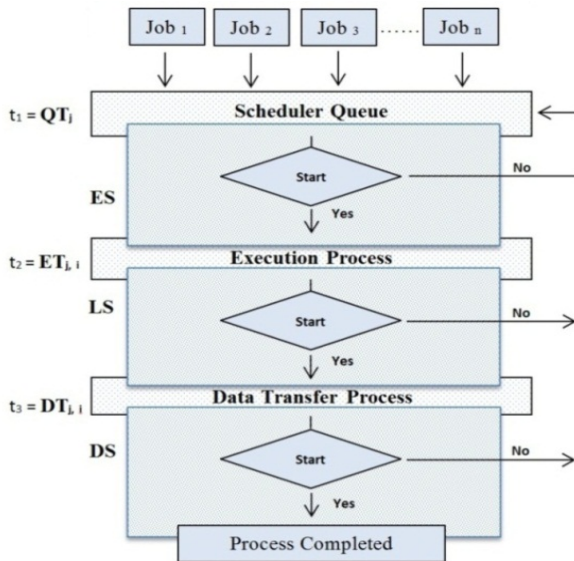


Fig. 2. Scheduling Architecture

External Scheduler (ES): ES decides the remote site to which to send the job to depending on some scheduling algorithm. ES uses external information for taking decision as input such as load at a remote site or the location of a dataset.

Local Scheduler (LS): Assigned jobs are managing by local scheduler to run at a particular site. The allocation of allocated jobs is also responsibility of the LS. LS decide about the priority and refusing to run jobs submitted by a certain user.

Dataset Scheduler (DS): At each site DS keeps track the popularity of each locally available data set. By using external information such as whether the data already exist at the concern site and loading to target remote site.

With encapsulation of three schedulers ES, LS and DS, the time for completion of a job is also encapsulation of three Times. QT_i or t_i is that time, which a job passing in waiting after entrance and before starting execution. Time denoted by $ET_{j,i}$ is the job execution time after entrance from queue and before starting transferring process. $DT_{j,i}$ is the time after completion execution time and before completion data transferring process. Means the Total Completion Time (TCT) for a job is the sum of all these three times.

3.3 Scheduling Strategy

The resource scheduling strategy is based on the estimation of cost (time) of executing a job in each site. It is possible to assume that job is submitted to LS (Local Scheduler) one by one. When receiving a job submission, the LS will estimate the time for completing executing (ETTC) a job in a site i , as mentioned by Nguyen et al. as below

$$ETTC_{ji} = \max \{DT(f(j), i), QT_{(i)}\} + EET_{ji}, \tag{3}$$

where DT is the Data Transfer Time for job; QT is the Queuing Time in site; and EET is the Estimate Execution Time for a job.

For Turnaround Time simply we use, Total Completion Time (TCT) for a job which is the sum of all these three Times as follow

$$TCT_{j,i} = QT_{(i)} + ET_{j,i} + DT(f(j)), \tag{4}$$

where TCT is Total Completion Time, QT is Queuing Time, ET is the job Execution Time and DT is Data Transfer time.

For assumptions, let T_w mean waiting time, T_s mean service (execution) time for each arrival, T_r mean time an item spends in system. Hence, According to Little's Law [21], we get

$$T_r = T_w + T_s. \tag{5}$$

In our equation (4) above, the residence time is given as

$$T_r = QT_{(i)} + ET_{j,i}, \tag{6}$$

where QT is the Time pass in Queue for a job and ET is the Service time which is Execution time in real as described in Figure 3 below.

Putting equation (6) in (4), the Total Completion Time for a job (TCT) is given by

$$TCT_{j,i} = QT_w + ET_s + DT(f(j)), \tag{7}$$

Where, QT_w is waiting time in a queue for a job, ET_s or QT_s is the passing time by processor during execution and is the Transfer time. In other words, we can say that QT_w is the *Queuing Time*, ET_s is the *Execution Time* and DT is the *Data Transfer Time* for a job. For waiting time in a queue, by Little’s Law [21], for single server, finite population formula given in [25], W_q will be as

$$QT_w = W_q = \frac{L_q}{\lambda(M - L)} \tag{8}$$

Where W_q is waiting time in queue, L_q is the number of jobs in queue, λ is jobs arrival rate, M is total population of jobs and L is the number of jobs in system. According to formula given in [25], ET will be as

$$QT_s(j) = ET(j) = \frac{L}{\lambda(M - L)} \tag{9}$$

Thus, the Data Transfer Time (DT) according to [18,19,20] can be formed as

$$DT(f(j)) = size f(k) / BW_{(i,j)}, \tag{10}$$

Where, Size $f(k)$ is the File Size in bytes and BW is the available Bandwidth between computing sites. Putting equations (8), (9), and (10) in (4) the result will be as

$$TCT_{ji} = \frac{L_q}{\lambda(M - L)} + \frac{L}{\lambda(M - L)} + size f(k) / BW_{(i,j)} \tag{11}$$

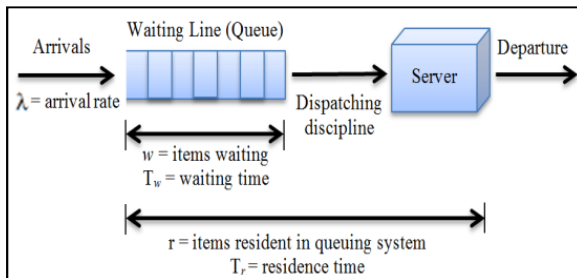


Fig. 3. Queuing system structure and parameters for single-server queue

3.4 Expected Results

In order to evaluate the performance of proposed architecture the author using CloudSim. Even the comparison of result may be will not so good, but with adding new one parameter, there will be great impact in term of accuracy during calculating the Total Time of Completion for a job, which will bring improvement in efficiency.

Dealing with large amount of data makes the requirement for efficiency in data access more critical. For the improvement of data efficiency, easy accessibility, strong reliability and always availability, the user of the world need more copy (replica) of local sub clouds on the basis of six sub-continent. Each local sub cloud will have same anytime updated copy of each other. Due to the local cloud, data in a wisely manner will offer a faster access to files required by cloud client, hence increase the job execution's performance. Due to local cloud, reliability and accessibility will increase. In this case if one sub cloud goes down, any client from anywhere in the world can use another sub cloud under the "failover" techniques. Cloud service failure is no acceptable in this architecture.

4 Conclusion

This paper presents architecture for cloud computing to support efficient data access for the job. Cloud should be a "true service provider" which is the demand of each user. The propose architecture will fulfill this property, because there is no chance to go down all local clouds. It is the ability of applications that are delivered in a Cloud vector, to provide speed and agility to the business and make it more competitive. For better scheduling techniques, this research gives importance to each parameter while calculating Total Time of Completion for a job. We need to evaluate our new model by using M/M/C, M/M/Inf, M/M/C/K and M/M/C/*M queuing models. There will be a great impact on accuracy by taking each parameter separately.

In future work, we plan to investigate more realistic scenarios and real user access patterns. Additionally, we plan to develop a complete real time model by using CloudSim, with a combination of replication [23] and scheduling.

Acknowledgement. This paper is supported by Postgraduate Research Grant Scheme (PRGS) no. vote GRS 090314 from Universiti Malaysia Pahang.

References

1. Weiss: Computing in the Cloud. ACM Networker 11, 18–25 (2007)
2. Brantner, M., Florescu, D., Graf, D., Kossmann, D., Kraska, T.: Building a Database on S3. In: The Conference Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data, pp. 251–263 (2008)
3. Moretti, Bulosan, J.: An abstraction for data-intensive cloud computing. In: Proceeding of IEEE International Symposium Parallel & Distributed Systems (IPDPS 2008), pp. 1–11 (2008)
4. Buyya, R., Yeo, C.S.: Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility. Future Generation Computer Systems 25, 599–616 (2009)
5. Foster, Yong, Z., Raicu, I., Lu, S.: Cloud computing and grid computing 360-degree compared. In: Proceeding of Grid Computing Environments Workshop (GCE 2008), pp. 1–10 (2008)

6. Dean, J., Ghemawat, S.: Map Reduce: Simplified data processing on large clusters. In: OSDI 2004: Proceedings of the 6th Conference on Symposium on Operating Systems Design & Implementation, pp. 137–149 (2004)
7. Ghemawat, S., Gobiuff, H., Leung, S.-T.: The Google file system. In: Proceedings of the Nineteenth ACM Symposium on Operating Systems Principles (SOSP 2003), pp. 29–43 (2003)
8. Amazon.com, Amazon S3 (2009), <http://aws.amazon.com/s3>
9. Borthaku, D.: The hadoop distributed file system: Architecture and design (2007), <http://lucene.apache.org/hadoop> (retrieved)
10. Hbase Development Team. Hbase: Bigtable-like structured storage for hadoopbase (2007), <http://wiki.apache.org/lucene-hadoop/Hbase>
11. Dean, J., Ghemawat, S.: MapReduce: Simplified data processing on large clusters. Communications of the ACM 51(1), 107–113 (2008)
12. Charrada, F.B., Ounelli, H., Chettaoui, H.: An Efficient Replication Strategy for Dynamic Data Grids. In: Proceedings of Conference on Grid, Cloud and Internet Computing (3PGCIC), pp. 50–54 (2010)
13. Shorfuzzaman: Distributed Popularity Based Replica Placement in Data Grid Environments. In: Proceedings of the 2010 ACM International Conference on Parallel and Distributed Computing, Applications and Technologies (PDCAT), pp. 66–77 (2010)
14. Zhao, W., Xu, X., Xiong, N.: A Weight-Based Dynamic Replica Replacement Strategy in Data Grids. In: Proceedings of IEEE Asia-Pacific Services Computing Conference (APSCC 2008), pp. 1544–1549 (2008)
15. Bsoul, M., Al-Khasawneh, A.: Enhanced Fast Spread Replication strategy for Data Grid. Journal of Network and Computer Applications 34(2), 575–580 (2011)
16. http://www.idc.com/prodserv/idc_cloud.jsp (November 1, 2011)
17. <http://www.yourdigitalspace.com/2010/10/the-amount-of-data-generated-and-consumed-on-the-internet/> (retrieved on November 27, 2011)
18. Ranganathan, K., Foster, I.: Simulation Studies of Computation and Data Scheduling Algorithms for Data Grids. Journal of Grid Computing 1(1), 53–62 (2003)
19. Tang, M., Lee, B.-S., Tang, Z., Yeo, C.-K.: The Impact of data replication on job scheduling performance in the Data Grid. Future Generation Computer Systems 22(3), 254–268 (2006)
20. Nguyen, D.N., Lim, S.B.: Combination of Replication and Scheduling in Data Grids. International Journal of Computer Science and Network Security 7(3), 304–308 (2007)
21. Stalling, W.: Queuing Analysis, <http://WilliamStallings.com/StudentSupport.html> (accessed on May 12, 2012)
22. Khan, N., Noraziah, A., Deris, M.M., Ismail, E.I.: *CLOUD COMPUTING: Comparison of Various Features*. In: Ariwa, E., El-Qawasmeh, E. (eds.) DEIS 2011. CCIS, vol. 194, pp. 243–254. Springer, Heidelberg (2011)
23. Zin, N.M., Ahmad, N., Fauzi, A.A.C., Herawan, T.: Replication Techniques in Data Grid Environments. In: Pan, J.-S., Chen, S.-M., Nguyen, N.T. (eds.) ACIIDS 2012, Part II. LNCS, vol. 7197, pp. 549–559. Springer, Heidelberg (2012)
24. Fauzi, A.A.C., Ahmad, N., Herawan, T., Zin, N.M.: On Cloud Computing Security Issues. In: Pan, J.-S., Chen, S.-M., Nguyen, N.T. (eds.) ACIIDS 2012, Part II. LNCS, vol. 7197, pp. 560–569. Springer, Heidelberg (2012)
25. D. I. Mag. Christian. Stationary Queuing Models with Aspects of Customer Impatience and Retrial Behaviour. Deutsch-Wagram, Vienna, Austria (2009), http://www.telecomm.at/documents/Stationary_QM.pdf

A Local Search Particle Swarm Optimization with Dual Species Conservation for Multimodal Optimization

Dingcai Shen^{1,2,*} and Xuewen Xia¹

¹ School of Computer and Information Science
Hubei Engineering University, Xiaogan, China

² State Key Laboratory of Software Engineering
Wuhan University, Wuhan, China
d.c.shen@163.com

Abstract. This work presents a new optimization technique called dual species conservation particle swarm optimization (DSPSO) for finding multiple optima (global or local) of multimodal functions. The basis of the proposed algorithm is repeatedly using species conservation and hill-valley detecting mechanism to refine the species set. To improve the balance between exploration and exploitation of the standard Particle Swarm Optimization (PSO), a local search around found optima strategy is adopted in PSO. The performance of DSPSO is validated on a set of widely used multimodal benchmark functions. Numerical results show that the proposed technique is effective and efficient in finding multiple solutions of selected benchmark.

Keywords: Particle swarm optimization, Species conservation, Multimodal optimization, Local search strategy.

1 Introduction

Particle Swarm Optimization (PSO) is a population-based optimization technique proposed by Kennedy and Eberhart [5]. PSO has been shown to successfully optimize a wide range of continuous functions. As other Evolutionary Algorithms (EAs) did, most of the application of PSO are often deal with single global optimum. However, in many real world application, there often exist multiple optima and worth to locate multiple global and local optima of a given objective function [3]. For example, in the field of engineering design, due to physical constraints, the best solutions may be hard to manufacture, or with easy access to assembly and maintenance, or for the reliability, etc., that we may choose the less fit solutions as our final choice.

* This work is supported by Youth Foundation of Hubei Engineering University (No. Z2007037) and Natural Science Foundation of Hubei Province of China (No. 2011CDC161).

Multimodal problems are often shown to be hard to solve by canonical EAs because of the presence of many local optima. There are some widely adopted techniques that extend the canonical EAs to enable them to identify and maintain multiple optima among the whole search space. These techniques can be considered from two aspects categories [8]: (i) iterative methods [14,11], which applying the same optimization algorithm iteratively or hybrid with other algorithms to locate multiple optima of a multimodal function; (ii)parallel (implicit or explicit) subpopulation models, which acquire multiple solutions for a multimodal optimization problems by dividing a population into several sub-population, niches or species that evolve in parallel, such as AFMDE [10], MNGA [13], Species based algorithm [2,7], Crowding [9,11],fitness sharing [4], and so on.

In this paper, a new multimodal optimization algorithm named dual species conservation particle swarm optimization (DSPSO) is presented. This DSPSO algorithm is designed by integrating species conservation and hill-valley detecting mechanism, to enhance the exploitation of canonical PSO, a local search strategy is adopted. The advantage of DSPSO is reduce the dependency of niche radius, which can be hard to choose without a prior. The second round species conservation in DSPSO result in the refinement of species seeds. Thus reduce the number of seeds to be conserved in the next generation. Numerical experiments on a set of widely used multimodal benchmark functions show that the DSPSO is able to effectively identify and maintain multiple optima of all tested functions.

The paper is organized as follows: section 2 describes related work on multimodal problems, and presents a brief of particle swarm optimization. DSPSO is described in detail in section 3. In section 4, the proposed algorithm is compared with canonical SCGA [7] on commonly used multimodal optimization benchmark functions. Finally, section 5 provides conclusion and discussion of the future work.

2 Background

2.1 Particle Swarm Optimization

Particle swarm optimization (PSO) [5] is a population-based stochastic optimization technique developed by Eberhart and Kennedy in 1995. It is inspired from the metaphor of social interaction observed among insects or animals. PSO can be easily implemented and has been proved to be both effective and fast when applied to many global optimization function.

PSO is similar to other population-based Evolution Algorithms (EAs). At the initial step, a population of potential solutions are randomly generated in a d -dimensional search space. Each particle (individual) of the population “flies” through the d -dimensional search space, adjusting its flying according to its own experience and that of neighboring particles. Thus, the movement of the particles have a tendency towards its best previously visited position and towards the position of the global best individual.

In a d -dimensional problem space, the i th particle of the swarm can be represented by a d -dimensional vector, $\mathbf{x}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,d})$ and $\mathbf{v}_i = (v_{i,1}, v_{i,2}, \dots, v_{i,d})$, respectively. The best previously visited position of the i th particle is denoted as $\mathbf{p}_i = (p_{i,1}, p_{i,2}, \dots, p_{i,d})$. For generation t of PSO algorithm, v_i and x_i are updated as follows:

$$v_i(t + 1) = \chi(v_i(t) + \varphi_1(p_i(t) - x_i(t)) + \varphi_2(p_g(t) - x_i(t))) \tag{1}$$

$$x_i(t + 1) = x_i(t) + v_i(t + 1) \tag{2}$$

where:

$$\varphi_1 = c_1r_1, \varphi_2 = c_2r_2,$$

$$\chi = \frac{2\kappa}{|2 - c - \sqrt{c^2 - 4c}|} \tag{3}$$

and c_1 and c_2 are the acceleration constants, typically set to 2.05. $c = c_1 + c_2$. κ is also a constant, usually set at 1. r_1, r_2 are uniformly distributed random number confined in the $[0, 1]$, and p_g is the best position visited so far by the considered particle's topological neighborhood; and χ is a constriction factor used to limit velocity.

The PSO algorithm performs repeated applications of the update process until a specified number of iterations has been reached, or the number of maximum evaluations has been exceeded.

2.2 Species Conservation

Literally in ecology, species are distinctly different of organisms. Competition for resources among organisms and the ever-changing environment impels some species to extinction, while others survive to maintain the balance in nature. In order to locate multiple optima in multimodal function optimization problems, many species based techniques have been introduced to existing EAs. Li et al. [7] introduced a recent technique called species conservation for multimodal optimization, which evolving parallel sub-populations by exploiting the detected species. The species conservation algorithm concentrates on two ways: The determination of species according to their similarity, each of these species is gathered around a dominating individual called the species seed; and the preservation of these species seeds found in the current generation to the next generation of the evolutionary cycle.

To define a species, a parameter known as species distance, which we denoted by σ_s , should be defined. The species distance specifies the maximum distance between two individuals for which they are considered to be similar and are considered to belong to the same species. At the beginning of the species determination, the population were sorted in decreasing order of fitness. The species seeds set denoted by X_s was initially set to empty. All individuals were checked successively against the species seeds found so far. A individual will be added to the species set if none of the found species seeds is within the half of the species distance ($\sigma_s/2$) to the individual considered.

The procedure that is used for conserving the species seeds is conducted after the new population was constructed. During the process of evolution, some species may not survive after the update of position of individuals, these seeds should be copied into the new population and thus enable them to survive.

2.3 Hill-Valley Detection Mechanism

Most niche technique introduce additional parameters, the commonly used parameter is niche radius, which is problem-dependent and difficult to select without priori knowledge. The hill-valley fitness topology function introduced by [12] is better than distance-based technique for it is more adaptable when the fitness landscape of the multimodal problem is unknown. Ursem’s hill-valley detect mechanism is the original approach that divides the population into sub-populations without use of niche radius and distances between individuals. The hill-valley detecting algorithm is simple, easily realized, and it can be easily extended to the high dimensions. The algorithm can be described as Fig. 1.

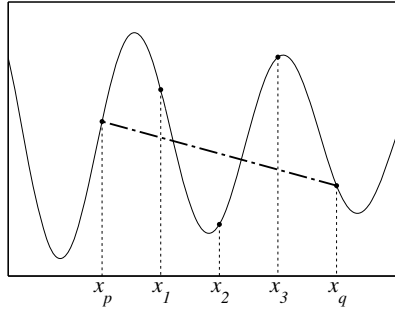


Fig. 1. The scheme for hill-valley function

where x_p and x_q are two given points in the search space. x_1, x_2, x_3 are three sample points (here we use three sample points; it may be vary according to the function considered). If any one of the three points with a fitness evaluation value that is smaller than the minimum of the two ends was found, the function stops and return true, otherwise, it returns false.

3 DSPSO

As mentioned above, like other niche-based techniques for multimodal optimization, the canonical SCGA [7] introduced an additional parameter named as species distance which we denoted by σ_s . The choosing of this parameter is crucial in species determination and species identification. Too many species may be identified in each iteration if we set σ_s to a too small value; meanwhile it will result in the increase of the overhead of this technique. On the other hand, too large value of σ_s will make many solutions indistinguishable in the process of species determination.

Another drawback of the canonical SCGA is the *acceptance threshold* (denoted by rf), which used to identify global optima in the end of iteration. The setting of rf is also need carefully setting, even for the same multimodal optimization problem, different runs may achieve different results with the same rf setting.

The simple using of hill-valley detection mechanism to identify peaks of multimodal optimization result in the increasing of fitness evaluation, and it may be identify the wrong peaks as described in [6]. These problem can be resolved by increasing the sample points in hill-valley function. However, this will further increase the overhead of fitness evaluation. In order to avoid the disadvantage of the canonical SCGA, and to take advantage of the simplicity of species determining and conservation procedure, the DSPSO algorithm combine the species conservation with hill-valley detection mechanism to effectively identify optima. In each iteration, the first round operation using the species conservation to identify species seeds with a roughly setting of species distance. The second round using hill-valley function to refine the species seeds set. The species conservation procedure using the refined species set to check wether a species seed should be conserved or updated in the next generation. The pseudocode of DSPSO are given as algorithm 1.

Algorithm 1. Pseudocode for DSPSO

```

1:  $Xrs \leftarrow \Phi$  //refined species seeds set
2:  $bool\ bSeed=false$ 
3: Use canonical species determining procedure to identify species seeds  $Xs$ 
4: for  $i = 1$  to  $Xs.Length$  do
5:    $bSeed \leftarrow true$ ;
6:   for all  $x \in Xrs$  do
7:     if not  $hill-valley(x, Xs[i])$  then
8:        $bSeed \leftarrow false$ 
9:     end if
10:  end for
11:  if  $bSeed = true$  then
12:     $Xrs \leftarrow Xrs \cup x$ 
13:  end if
14: end for
15: Use  $Xrs$  as species seeds for species conservation procedure

```

4 Results and Analysis

4.1 Test Functions

In order to compare the performance of DSPSO algorithm to identify and maintain optima, we use benchmark functions described in Table 1. All these functions are widely used for comparing the performance of multimodal optimization problems. The number of optima (global or local) of the functions is ranging from 2 to 6, dimensionality is from 1 to 2, the high dimensionality problem is for our future research.

Table 1. Standard benchmark functions adopted in this work

<i>Test function</i>	<i>Expression</i>	<i>Remarks</i>
<i>Deb's function</i>	$f_1(x) = \sin^6(5\pi x)$	5 peaks, where $0 \leq x \leq 1$
<i>Deb's decreasing function</i>	$f_2(x) = 2^{-2 \ln(2)((x-0.1)/0.9)^2} \sin^6(5\pi x)$	5 peaks, where $0 \leq x \leq 1$
<i>Branin RCOS</i>	$f_3(x, y) = (y - \frac{5.1}{4\pi^2} \cdot x^2 + \frac{5}{\pi} \cdot x - 6)^2 + 10 \cdot (1 - \frac{1}{8\pi}) \cdot \cos(x) + 10$	3 peaks, where $-5 \leq x \leq 10$ $0 \leq y \leq 15$
<i>Ursem F3</i>	$f_4(x, y) = \sin(2.2\pi x + 0.5\pi) \cdot \frac{2- y }{2} \cdot \frac{3- x }{2} + \sin(0.5\pi y^2 + 0.5\pi) \cdot \frac{2- y }{2} \cdot \frac{2- x }{2}$	5 peaks, where $-2 \leq x \leq 2$, $-1.5 \leq y \leq 1.5$
<i>Michalewicz</i>	$f_5(x, y) = \sin(x) \cdot \sin^{20}(\frac{x^2}{\pi}) + \sin(y) \cdot \sin^{20}(\frac{2y^2}{\pi})$	2 peaks, where $0 \leq x, y \leq \pi$

4.2 Performance Measures

A number of performance measures for multimodal optimization algorithm have been used in literature. The following performance measures are considered in this study:

- **Success rate**(%) : the percentage of successful runs in which all optima (global and local) are successfully located.
- **Peak ratio** : the sum of the optima identified by the methods divided by the sum of the actual optima in the search space. An optimum is considered to be detected if it is within a Euclidean distance of 0.01 for f_1 , f_2 and 0.05 for f_3 , f_4 , f_5 from the real optimum.
- **NFEs** : The number of fitness function evaluation is recorded when the prespecified threshold is attained. The average and standard deviation of the NFEs are adopted in this study.

4.3 Results and Analysis

All experiments consisted of 30 runs. The PSO parameters setting in all benchmark functions are as following. A population size was set to 30, $c_1 = c_2 = 2.05$. In SCGA, A roulette-wheel selection was applied. A single-point crossover operator with a probability $p_c = 0.6$ was set. The mutation probability p_m was set to 0.05. The species distance σ_s used for SCGA were found by trial and error in a few preparation runs for each functions. The best results achieved with the optimum setting are listed in Table 2.

Table 3 presents descriptive statistics related to the five benchmarks, i.e. the values for the mean of the average success rate, the mean and standard deviation of the fitness evaluation and the mean peak ratio. The results presented in the table are the functions to be tested on each algorithm, the success rate, the mean number of fitness evaluation required and the standard deviation it, and the mean peak ratio.

Table 2. Parameter setting for the test suite

<i>Funcs.</i>	f_1	f_2	f_3	f_4	f_5
σ_s	0.1	0.1	1.0	0.85	2.0

We compare the results obtained using DSPSO with species conservation GA (SCGA) [7]. From the Table 3 we can see both DSPSO and SCGA are able to identify and to maintain multiple optima on all tested benchmark functions (except for f_4 on SCGA). The possibly reason of the failure of SCGA on f_4 due to the specific landscape of f_4 , which makes it difficult for SCGA to distinguish an individual within the same peak from its neighbor that belong to the other peaks. The number of function evaluations of DSPSO is slightly higher than that of SCGA because of the introduce of hill-valley function. It must be taken into account that the distance calculation (not described in results) also increase the overhead of the algorithm, which used in species determination and conservation in both algorithm. With the refining of the species seeds set, the distance calculation reduces rapidly in DSPSO. The preliminary experiments indicate that DSPSO is a competitive candidate for multimodal optimization problems.

Table 3. Experimental results(Averaged over 30 runs)

<i>Funcs.</i>		<i>Success Rate</i> (%)	<i>NFEs</i>	<i>Std. Dev.</i>	<i>Peak Ratio</i>
f_1	DSPSO	100	1342	443	1.0
	SCGA	100	1250	356	1.0
f_2	DSPSO	100	1032	343	1.0
	SCGA	100	786	246	1.0
f_3	DSPSO	100	16221	448	1.0
	SCGA	100	16245	545	1.0
f_4	DSPSO	100	6045	435	1.0
	SCGA	58	5420	520	0.64
f_5	DSPSO	100	5648	452	1.0
	SCGA	100	5503	402	1.0

5 Conclusions

In this paper we presented the DSPSO algorithm based on the idea of hybrid of species conservation and hill-valley detect mechanism. The performance of the DSPSO is compared against the canonical species conservation GA. All experiments have demonstrated that the DSPSO is a competitive candidate for multimodal evolutionary algorithm. Future work will focus on more experimentation in handling complex real-world multimodal optimization problems and problems with higher dimensionality.

References

1. Beasley, D., Bull, D.R., Martin, R.R.: A Sequential Niche Technique for Multimodal Function Optimization. *Evolutionary Computation* 1(2), 101–125 (1993)
2. Cho, H., Kim, D., Olivera, F., Guikema, S.D.: Enhanced speciation in particle swarm optimization for multi-modal problems. *European Journal of Operational Research* 213(1), 15–23 (2011)
3. Das, S., Maity, S., Qu, B.Y., Suganthan, P.: Real-parameter evolutionary multimodal optimization—A survey of the state-of-the-art. *Swarm and Evolutionary Computation* 1(2), 71–88 (2011)
4. Della Cioppa, A., De Stefano, C., Marcelli, A.: Where Are the Niches? Dynamic Fitness Sharing. *IEEE Transactions on Evolutionary Computation* 11(4), 453–465 (2007)
5. Kennedy, J., Eberhart, R.: Particle swarm optimization. In: *Proceedings of ICNN 1995 - International Conference on Neural Networks*, vol. 4, pp. 1942–1948 (1995)
6. Kharma, N.: On Clustering in Evolutionary Computation. In: *2006 IEEE International Conference on Evolutionary Computation*, pp. 1752–1759 (2006)
7. Li, J.P., Balazs, M.E., Parks, G.T., Clarkson, P.J.: A species conserving genetic algorithm for multimodal function optimization. *Evolutionary Computation* 10(3), 207–234 (2002)
8. Li, J.P., Li, X.D., Wood, A.: Species based evolutionary algorithms for multimodal optimization: A brief review. In: *IEEE Congress on Evolutionary Computation*, pp. 1–8 (2010)
9. Li, M., Lin, D., Kou, J.: A hybrid niching PSO enhanced with recombination-replacement crowding strategy for multimodal function optimization. *Applied Soft Computing* 12(3), 975–987 (2012)
10. Shen, D., Li, Y., Wei, B., Xia, X.: Adaptive Forking Multipopulation Differential Evolution Algorithm for Multimodal Optimization. *Journal of Convergence Information Technology* 7(5), 57–65 (2012)
11. Thomsen, R.: Multimodal optimization using crowding-based differential evolution. In: *Proceedings of the 2004 Congress on Evolutionary Computation*, vol. 2, pp. 1382–1389 (2004)
12. Ursem, R.: Multinational evolutionary algorithms. In: *Proceedings of the 1999 Congress on Evolutionary Computation*, pp. 1633–1640 (1999)
13. Ursem, R.: Multinational GAs: Multimodal optimization techniques in dynamic environments. In: *Proc. of the Genetic and Evolutionary Computation Conference*, pp. 19–26 (2000)
14. Vitela, J.E., Castañón, O.: A sequential niching memetic algorithm for continuous multimodal function optimization. *Applied Mathematics and Computation* (2012), doi:10.1016/j.amc.2011.05.051

Cloud Computing: Analysis of Various Services

Nawsher Khan¹, A. Noraziah¹, Tutut Herawan¹, and Mustafa Mat Deris²

¹ Faculty of Computer Systems and Software Engineering Universiti Malaysia Pahang Lebu
Raya Tun Razak, 26300 Gambang, Kuantan, Pahang, Malaysia

² Faculty of Computer Science and Information Technology Universiti Tun Hussein Onn
Malaysia Parit Raja, 86400 Batu Pahat, Johor, Malaysia
nawsherkhan@gmail.com, {noraziah,tutut}@ump.edu.my,
mmustafa@uthm.edu.my

Abstract. Cloud computing fulfills the long-held dream of computing as a utility and fundamentally altering the expectations for how and when computing, storage and networking resources should be allocated, managed, consumed and allow user to utilize services globally. Due to the powerful computing and storage, high availability and security, easy accessibility and adaptability, reliable scalability and interoperability, cost and time effective, cloud computing is the top needed for current fast growing business world. A client, organization or a trade that adopting emerging cloud environment can choose a well suitable infrastructure, platform, software and a network resource, for any business, where each one has some exclusive features and advantages. In this paper, we first present a comprehensive classification for describing cloud computing architectures. This classification help in survey of several existing cloud computing services developed by various projects globally such as Amazon, Google, Microsoft, Sun and Force.com. Then by using this survey results, we identify similarities and differences of the architecture approaches of cloud computing.

Keywords: Cloud Computing, Platform, Virtualization.

1 Introduction

Cloud computing fulfills the long-held dream of computing as a utility and it allows leasing of IT capability, thus represents an modulation point in the natural features of computation and IT services deliverance [1,2]. Cloud computing give services in the form of infrastructure, platform, or software as services in a pay-as-you-go model. With a trend toward cloud based model, the power is shifted to consumers. This paradigm marks an elementary yet massive shifting from the traditional “Desktop-As-A-Platform” to “Internet-As-A-Platform” model. To achieve the infinite scalability, guaranteed performance, easy accessibility and nearly “Always-On” availability demands, these computing platforms typically are deployed in clusters of massive number of servers hosted in dedicated data centers [2]. In the cloud, virtualization occurs at several levels. It can range from ‘what does what’ (server & application virtualization) to ‘what goes where’ (data storage virtualization) to ‘who is where’ (mobility and virtual networking). The beauty of virtualized solutions is that user can run multiple operating systems simultaneously on a single host.

Cloud Computing has emerged recently as a label for a particular kind of datacenter, or most commonly, a group of datacenters. Computing capability has become the bottleneck of systems using traditional grid computing, which demands higher hardware requirements. Cloud computing is a kind of computing platform distributed in large-scale data center, which meets the requirements of scientific research and e-commerce by dynamically providing several types of server resources [3]. Cloud computing platform utilize the virtualization technology to transparently and dynamically supply virtual computing and storage resources for the satisfaction of user's different requirements according to the relative scheduling strategies. In this paper, we first present a comprehensive classification for describing cloud computing architecture. This classification helps in survey of several existing cloud computing services developed by various projects globally such as Amazon, Google, Microsoft, Sun and Force.com. Then by using this survey results, we identify similarities and differences of the architecture approaches of cloud computing.

The rest of the paper is organized as follows. Section 2 describes the rudimentary on cloud and its systems components. Section 3 describes analysis of various services of cloud computing. Finally, the conclusion of this work is presented the Section 4.

2 Rudimentary

2.1 Cloud Computing

Cloud computing brings the difference from traditional IT approaches is the focus on service delivery as well as the consumer utilization model. In the background, service provider's uses system architecture, particular technologies, industry best practices and design to provide and support the delivery of service-oriented and elastically scalable environment to provide better services to multiple customers [4]. Platform-as-a-Service solutions provide applications development platforms and environment for seamlessly incorporate Cloud computing into existing services, application and infrastructure with a market-oriented approach. Cloud computing is emergent based on year's achievement on Grid computing, Virtualization, Utility computing, Web computing and related technologies. Cloud computing provides both platforms and applications on-demand through Internet or intranet [5,6,7,8]. Some examples of emerging Cloud computing platforms are Amazon EC2 [9], IBM Blue cloud [10, 11], Google AppEngine [12] and Microsoft Azure [13].

The Cloud allows sharing, aggregation and allocation of software, storage and computational network resources on-demand. Some of the key benefits of Cloud computing include hiding and abstraction of complexity, virtualized resources and efficient use of distributed resources [6]. Generally Cloud can be classified in three categories as public, private, or hybrid depending on the model of deployment [1].

A Public Cloud is a Cloud made available in a pay-as-you-go manner to the general public. In a typical public Cloud scenario, a third-party vendor delivers services such as computation, storage, networks, virtualization and applications to various customers [14, 15]. Businesses are adopting public Cloud services for sieving capital expenditure and operational cost by leveraging Cloud's elastic scalability and market oriented costing features. Nevertheless, public Cloud computing also raises

concerns about data security, data transfer, management, performance, and level of control. A Private Cloud is a data center of an organization, not made available to the general public. As shown in Figure 1. In a private Cloud environment, internal IT resources are used to serve their internal users and customers. A Hybrid Cloud is a seamless use of public Cloud along with private Cloud when needed. Focus on service delivery and the consumer utilization model makes Cloud computing different from traditional IT approaches [1].

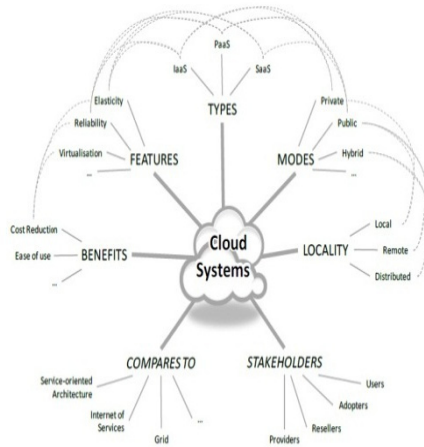


Fig. 1. Non-exhaustive view of cloud system

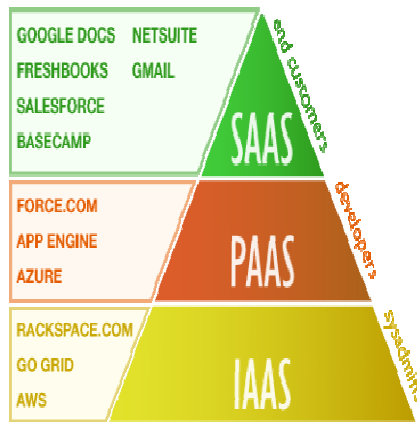


Fig. 2. Structure of Cloud Computing

Platform as a Service (PaaS) of cloud computing is one of the key services in Cloud computing. As illustrated in Figure 2, PaaS is the delivery of a computing platform and environment as well as solution stack as a service without software downloads or installation for developers, IT managers or for end-users, which is also known as “Cloudware”. Through virtualization and other resource sharing

mechanisms, Cloud computing can dramatically reduce costs for the user, which is the need and high demand of today's user. Virtualization techniques [15] make it possible to open a few logical platforms in a single physical machine (Windows, Linux etc); so that resources can be shared better and more users can get benefit and served at a time. Most of Cloud computing platforms are based on virtualized environments. In a virtualized Cloud computing lab, there are four main parts: software and hardware platforms (PaaS resources) provided by virtualized and real servers; resource management node; database servers and users who access these resources through Internet or Intranet [16]. Cloud computing is an evolving tools. In general User can keep their data, which can be stored somewhere, in some part of the world. No need to know about data that where data is going and where their data is residing. For further detail, in the cloud computing world, there is variability in terms of where the physical data resides, where processing takes place, and from where the data is accessed. Like this if a user wants to run some Application (i.e. web services), no need to have special software or hardware on their own machine to run their application. Cloud computing is having the feature of cluster computing, grid computing, service computing and utility computing.

A User has to run one time application means, user can use cloud platform as the best choice. In future users can have dummy terminal along with keyboard and mouse to feed the data. No need to invest on resources (software and hardware), instead they can rent or lease the resources. Moreover they can choose the best available computing platform, since many companies are going to offer cloud based applications in near future, like mobile services. Nowadays, many data centers itself are available in which users can store their data up to 5GB and 10GB, some cloud services provide 30GB space, freely without any cost, through their online operating system. If more space is needed, with less premium amount they can store their own data. Like that, computing and data transferring with fewer premiums also available [11]. Anyhow, Cloud computing is still passing its infancy stage, there are many challenging issues waiting for tackling [5,6,7,8].

2.2 Cloud System Components

In this survey of open-source cloud computing system of Eucalyptus, Open Nebula and Nimbus, we make a quick overview of the entire cloud computing. Generally open-source cloud computing system has six components.

Hardware and Operating System: Hardware and software are the various back bone of any physical machine in cloud system. While proper set up is necessary for any software system. Firstly, for running pure virtualization, Secondly open source frameworks are more flexible to do work with various systems.

Hypervisor: Hypervisor which is also known as Virtual Machine Monitor (VMM). Generally, popular VMMs consists Xen, KVM and Virtual Box, which are open-source, and VMware is commercial.

Framework: Cloud framework itself is an important component of the cloud system. Frameworks, where we can put Eucalyptus, Open Nebula or Nimbus.

Network: Network is an important component of cloud computing and, which includes DNS, DHCP and subnet organization of physical machine. Virtual bridge, which is also a part of network, provides unique virtual MAC address to each virtual machine (VM).

Disk Image: A virtual hardrive is the basic need to be functional a virtual machine. When we need a single VM on a single physical machine, VM installs an operating system and other software after creating a blank disk image.

Front-End: For user request there must be an interface, through which a user can interact with virtual machine (VM), specify parameter in order to login to the created VMs.

3 Analysis of Various Services

3.1 Cloud Platform-as-a- Service (PaaS)

In our Global village, there are various cloud computing platforms; each one has its own characteristics and advantages discussed by our previous work. For better understanding, we analyze these platforms and give comparison with different implementation aspects as described in Table 1.

Table 1. Comparison of Some Cloud Computing Platforms

Property	Different Platforms				
	Amazon Elastic Compute Cloud (EC2)	Microsoft Azure	Google App Engine	Sun Network.com (Sun Grid)	GRIDS Lab Aneka
Focus	Infrastructure	Platform	Platform	Infra-structure	Enterprise clouds
Service Type	Compute, Storage (Amazon S3)	Web and non-web application	Web Application	Computing	Computing
User Access interface	Amazon EC2 command-line tools	Microsoft windows azure portal	Web-based administration	scripts, Sun Grid web portal	Work-bench, web-based portal
Value-added service providers	Yes	Yes	No	Yes	No
Virtualization	OS level running on a Xen hypervisor	OS level through fabric controller	Application container	Job management system (Sun Grid Engine)	Resource manager and scheduler
Web APIs	Yes	Yes	Yes	Yes	Yes
Dynamic negotiation of QoS	None	None	None	None	SLA-base resources reservation
Programing frame-work	Amazon Machine Images (AMI)	Microsoft .NET	Python	Solaris OS. Java, C, C++, FORTRAN	APIs supporting models in c# .Net

3.2 Comparison of Cloud Platforms with Implementation Aspects

We have different kinds of cloud platforms; each one has its own characteristics and advantages. For better understanding, we analyze and give with detail comparison from different implementation aspects as discussed in Table 2.

Table 2. Comparison of cloud platforms with implementation aspects

	Eucalyptus	Nimbus	OpenNebula
Cloud Character	Public	Public	Private
Scalability	Scalable	Scalable	Dynamical, Scalable
Cloud Form	IaaS	IaaS	IaaS
Compatibility	Support EC2, S3	Support EC2	Open, Multi-Platform
Deployment	Dynamical Deployment	Dynamical Deployment	Dynamical Deployment
Deployment Manner	Commandline	Commandline	Commandline
Transplantability	Common	Common	Common
VM Support	VMWare, Xen, KVM	Xen	Xen, VMWare
Web Interface	Web Service	EC2 WSDL, WSRF	Libvirt, EC2, OCCI API
Structure	Module	Lightweight Components	Module
Reliability	-	-	Rollback host and VM
OS Support	Linux	Linux	Linux
Development Language	Java	Java, Python	Java

Table 3. Comparison of open-source cloud platforms

Feature	OpenNebula
Computing Architecture	-Cluster into an IaaS cloud -Focused on the efficient, dynamic and scalable management of VMs within datacenters (private cloud) involving a large amount of virtual and physical servers -Based on Haizea scheduling
Virtualization Management Service	-Xen KVM and on-demand access to Amazon EC2 IaaS
Load Balancing	-Nginx Server configured as load balancer, used round-robin
Interoperability	-Interoperable between intra cloud services -The daemon can be restarted and all the running VMs recovered
Fault Tolerance	-Persistent database backend to store host and VM information
Security	-Firewall, Virtual Private Network Tunnel
Programming Framework	-Java, Ruby
Storage	-Database, persistent storage for ONE data structures -SQLite3 backend is the core component of the OpenNebula internal data structures

Table 3. (continued)

Eucalyptus	Nimbus
-Ability to configure multiple clusters, each with private internal network addresses, into a single cloud.	-Science cloud
-Private Cloud.	-Client-Side cloud-computing interface to Globus-enabled TeraPort cluster
	-Nimbus Context Broker that combines several deployed virtual machines into “turnkey” virtual clusters
	- Heterogeneous clusters of auto-configuring VMs with one command
-Xen hypervisor	-Xen Virtualization
IaaS	IaaS
-Simple load-balancing cloud controller	-Launches self-configuring virtual cluster i.e. the context broker
-Multiple cloud computing interfaces using the same “back-end” infrastructure	-Standards : “rough consensus and working code”
-Separate cluster within the Eucalyptus cloud reduce the chance of correlated failure	-Checking worker nodes periodically and recovery
-WS-security for authentication, cloud controller generates the public/private key	-PKI credential required
-Hibernate, Axis2 and Axis2c, Java	-Works with Grid proxies VOMS, Shibboleth (via GridShip), custom PDPs
-Walrus (the front end for the storage subsystem)	Python, Java
	-Grid FTP and SCP

3.3 Open Source Cloud

The role of open source cloud computing is to build some mechanism around digital identity management, and outlines some technological building blocks are needed for controllable trust and identity verification. Open Nebula and Nimbus are technically sound and popular. Current Cloud is focusing on the issue of interoperability which is essential for enterprise cloud system. Most of the open source clouds are provided IaaS as shown in Table 3.

4 Conclusion

Cloud Computing is the fifth utility after water, electricity, gas and telephony, and it is the promising paradigm for delivering IT services as computing utilities. This paper presents a comprehensive comparison of different aspects of cloud’s platforms. In analysis of these various open-source cloud computing frameworks, we found that there are salient philosophical differences between them regarding the overall scheme of their design. After this analysis user can better understand the characteristic and will be able to do better selection of cloud platform, implementation and deployment requirement. In current cloud still we have challenges i.e. continuously availability, data security and privacy. In current cloud environment, user can’t fine the status of

their data may be someone is using these data for his/her own purposes. Our future work lies in the areas of data replication and data scheduling in cloud computing as well as on the combination of these both, replication and scheduling techniques.

Acknowledgement. This paper is supported by Postgraduate Research Grant Scheme (PRGS) no. vote GRS 090314 from Universiti Malaysia Pahang.

References

1. Buyya, R., Sukumar, K.: Platforms for Building and Deploying Application for Cloud Computing. *CSI Communication* 35(1), 6–11 (2011)
2. Parkhill, D.F.: *The Challenge of the Computer Utility*, 1st edn. Addison-Wesley (1966)
3. Tian, W.: A Framework for Implementing and Managing Platform as a Service in a Virtual Cloud Computing Lab. In: *Proceeding of Second International Workshop on Education Technology and Computer Science (ETCS 2010)*, pp. 273–276 (2010)
4. Filho, O.F.F., Ferreira, M.A.G.F.: Semantic Web Services: A restful Approach. In: *Proceeding of IADIS International Conference WWW/Internet*, pp. 169–180 (2009)
5. Bernstein, D., Vidovic, N., Modi, S.: A Cloud PAAS for High Scale, Function, and Velocity Mobile Applications. In: *Proceedings of the Fifth ACM International Conference on Systems and Networks Communications (SNC 2010)*, pp. 117–123 (2010)
6. Armbrust, M., et al.: Above the Clouds: A Berkeley View of Cloud Computing. Technical Report No. UCB/Eecs-2009-28 (2009)
7. Nurmi, D., et al.: The Eucalyptus Open-source Cloud-computing System. In: *Proceedings of 9th IEEE/ACM International Symposium on Cluster Computing and the Grid (CCGRID 2009)*, pp. 124–131 (2008)
8. Cavoukian, A.: Privacy in the Clouds: Privacy and Digital Identity- Implication for the Internet. Information and Privacy Commissioner of Ontario (May 28, 2008)
9. Amazon: Cloud Computing Issues, Research and Implementations. In: *ITI 2008*, pp. 23–31 (2008)
10. IBM blue cloud, <http://www.ibm.com/grid/> (accessed on May 09, 2011)
11. Boss, G., et al.: Cloud Computing. IBM Corporation white paper (October 2007)
12. Google App Engine, <http://www.aws.amazon.com/ec2/> (May 18, 2011)
13. Microsoft-Azure, <http://www.microsoft.com/windowsazure/windowsazure> (May 22, 2011)
14. Rimal, B.P., et al.: A Taxonomy and survey of cloud computing system. In: *Proceeding of Fifth International Joint Conference on INC, IMS and IDC*, pp. 44–51 (2009)
15. Loganayagi, B., Sujatha, S.: Creating virtual platform for cloud computing. In: *Proceeding of IEEE International Conference on Computational Intelligence and Computing Research (ICCIC 2010)*, pp. 1–4 (2010)
16. Zhang, X., Dong, G.: A New Architecture of Online Trading Platform Based on Cloud Computing. In: *Proceeding of Asia-Pacific Conference on Wearable Computing Systems (APWCS 2010)*, pp. 32–35 (2010)

Quantum Ant Colony Algorithm Based on Bloch Coordinates

Xiaofeng Chen, Xingyou Xia*, and Ruiyun Yu

Software College, Northeastern University, Shenyang, China
{neucxf,neuxiaxy}@163.com, yury@mail.neu.edu.cn

Abstract. Given classic Ant Colony Algorithm only resolves the optimization problem of discrete system, this paper proposed a Quantum Ant Colony Algorithm (QACA) based on the Bloch spherical coordinate by combining Quantum Evolutionary Algorithm and Ant Colony Algorithm. This algorithm applies Bloch spherical coordinate of Qubits to represent the current position information of ants; a new quantum revolving door is designed for updating the position to achieve to watch ants' movement. Quantum doors help to realize the variation of ants' positions, increase the diversity. For different optimization problems, various solution space transformational models and fitness functions are planned, so as to optimally solve the target. Furthermore, simulations of function extreme value and TSP problems were conducted, which indicted that the algorithm is feasible and effective.

Keywords: quantum computing, Bloch Coordinates, Quantum Ant Colony Algorithm.

1 Introduction

Ant Colony Algorithm (ACA) [1] was originally proposed by Macro Dorigo, Italian scholar, in 1990s, which was used to construct a typical NP Hard problem — Traveling Salesman Problem (TSP), then, ACA was applied to combinatorial optimization problem similar to TSP problem, such as Knapsack Problem [2], Assignment Problem [3], Job-shop Assignment [4], Sequential Ordering [5], Network Routing [6], Vehicle Routing [7], Power System [8] and Controls Parameter Optimization [9], etc, which obtained satisfactory effect. However, the biology background of classical ACA limits itself only resolving the discrete system optimization problem, therefore, how to use ACA to achieve problem solving for continuous space optimization effectively is a challenging research work.

Combining quantum evolutionary Algorithm [10] with Ant Colony Algorithm, this paper puts forward Bloch [11] Quantum Ant Colony Algorithm (BQACA), and various solution space transformational models and fitness functions are planned for different optimization problems. Algorithm in this paper is verified by function extreme value problem and Traveling Salesman Problem respectively. The result of simulation shows that the algorithm not only expresses high efficiency of quantum computing, but also maintains the preferable optimizing and robustness of colony algorithm.

* Corresponding author.

2 Quantum Ant Colony Algorithm (QACA)

On the Bloch spherical coordinate, one qubit can be expressed as: $|\phi\rangle = \cos\frac{\theta}{2}|0\rangle + e^{i\phi}\sin\frac{\theta}{2}|1\rangle$; any point on the sphere can be confirmed via θ and ϕ : $|\phi\rangle = [\cos\phi\sin\theta, \sin\phi\sin\theta, \cos\theta]^T$. Suppose there are n ants in the ant colony, where each ant carries a group (m units) of qubit, its Bloch spherical coordinate shows the current location of ant, corresponding to approximate solution of optimization problem. In all locations occupied by current ant colony, the one with highest value is defined as the optimal position, while the ant who takes that position is defined as the optimal ant.

2.1 Initialize Ant Colony

P_i is set as the location of the i th ant, considering that the randomness of coding for ant colony and constraint conditions for probability amplitude of the quantum state, the initialization of BQACA is expressed as:

$$\begin{bmatrix} P_{ix}^j \\ P_{iy}^j \\ P_{iz}^j \end{bmatrix} = \begin{bmatrix} \cos\phi_{i1}\sin\theta_{i1} & \cos\phi_{i2}\sin\theta_{i2} & \cdots & \cos\phi_{im}\sin\theta_{im} \\ \sin\phi_{i1}\sin\theta_{i1} & \sin\phi_{i2}\sin\theta_{i2} & \cdots & \sin\phi_{im}\sin\theta_{im} \\ \cos\theta_{i1} & \cos\theta_{i2} & \cdots & \cos\theta_{im} \end{bmatrix} \quad (1)$$

Where $\phi_{ij} = 2\pi rand$, $\theta_{ij} = \pi rand$, $rand$ are random numbers between $(0,1)$; $i \in \{1, 2, \dots, n\}$, $j \in \{1, 2, \dots, m\}$, n for number of ant; m for number of qubit. 3 coordinates of qubit are regarded as 3 paratactic genes, and each ant contains 3 gene chains, which are called X-chain, Y-chain and Z-chain respectively, each gene chain stands for an optimal solution $P_{ix}^j, P_{iy}^j, P_{iz}^j$.

2.2 Transformation of Solution Space

In BQACA, owing to $[-1,1]$ per dimension of ants' travelling space, solution space changing is needed here to figure out whether ant's current position is superiority or inferiority, which means each ant would occupy three positions, and the unit space of the three position would map to solution space of optimization problem, making every probability amplitude of qubit on ant could correspond to an optimizing variable quantity of solution space. The author would take function extreme value problem and give a further explain in this paper.

Solution space transformation approach for function extreme-value problem: propose the domain of definition of variable X^j is its solution space $[a_j, b_j]$, record the j th qubit as $[\cos\phi_{ij}\sin\theta_{ij}, \sin\phi_{ij}\sin\theta_{ij}, \cos\theta_{ij}]^T$ by using linear transformation, then the corresponding solution space variable is:

$$\begin{bmatrix} X_{ix}^j \\ X_{iy}^j \\ X_{iz}^j \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 + \cos\phi_{ij}\sin\theta_{ij} & 1 - \cos\phi_{ij}\sin\theta_{ij} \\ 1 + \sin\phi_{ij}\sin\theta_{ij} & 1 - \sin\phi_{ij}\sin\theta_{ij} \\ 1 + \cos\theta_{ij} & 1 - \cos\theta_{ij} \end{bmatrix} \begin{bmatrix} b_j \\ a_j \end{bmatrix} \quad (2)$$

Solution space transformation approach for TSP problem: this paper has designed two-layer transformational model in the aspect of solution space aiming at the specific characteristic of TSP problem, the model contains two transformations — linear transformation and lead transformation.

Linear transformation: qubit is transformed from unit space to lead space. Propose the definitional domain of lead message variable, r^j , is $[0, 1]$, formula (2) is used to calculate corresponding lead solution space variable $[\tau_{ix}^j, \tau_{iy}^j, \tau_{iz}^j]^T$.

Lead transformation: impact strength of lead message and inspire message to TSP solution could be regulated by adjusting lead factor and inspire factor. Strategy is selected according to lead probability and roulette to carry out optimal decode. Suppose the current node as i , select node j as the next visiting node:

$$p_{ij}^k = \begin{cases} \frac{r_{ij}^\omega(t) \cdot \lambda_{ij}^\nu(t)}{\sum_{s \in allowed_k} r_{is}^\omega(t) \cdot \lambda_{is}^\nu(t)} & j \in allowed \\ 0 & otherwise \end{cases} \quad (3)$$

Where $r_{ij}^\omega(t) \cdot \lambda_{ij}^\nu(t)$ is for message of path, $r_{ij}(t)$ stands for lead message, ω is lead factor; $\lambda_{ij}(t)$ represents inspire message $\lambda_{ij}(t) = 1/d_{ij}$, d_{ij} means the distance from City i to City j , ν is inspire factor; $allowed_k = \{1, 2, \dots, m\} - tabu_k$ means node set of available city may selected by ant k at the time t ; $tabu_k$ is used to keep the routing table which obtained by transforming ant k .

2.3 Define Fitness Function

A variety of fitness function needs to be designed for different optimal problems, the more fitness it is, the better solution for individual.

Fitness function of extreme-value problem: suppose $f(X_i)$ as the i th solution, $fit(X_i)$ is the adaptive value for the i th solution.

When evaluate the minimum value:

$$fit(X_i) = \begin{cases} \frac{1}{1+f(X_i)} & f(X_i) \geq 0 \\ 1 + abs(f(X_i)) & f(X_i) < 0 \end{cases} \quad (4)$$

When evaluate the maximum value:

$$fit(X_i) = \begin{cases} 1 + f(X_i) & f(X_i) \geq 0 \\ 1 + \frac{1}{1+abs(f(X_i))} & f(X_i) < 0 \end{cases} \quad (5)$$

TSP fitness function: fitness of individual $X_i = \{x_1, x_2, \dots, x_m\}$ of TSP is defined as the reciprocal of path length represented by individual.

$$fit(X_i) = \frac{1}{D(X_i)} \quad (6)$$

2.4 Update of Ant Position

In the solution space of optimal problem, suppose $\tau(X_i)$ is the strength of pheromone of k th ant at X_i , initial moment all set as some constant: $\eta(X_i)$ stands for the visibility at X_i . The basic framework of QACA described as follows:

- (1) Select the target position of ant movement

The transition rule and transition probability of ant k from position X_i to position X_s are:

$$X_s = \begin{cases} \arg \max_{X_s \in P} \{ \tau^\alpha(X_s) \cdot \eta^\beta(X_s) \} & q \leq q_0 \\ \tilde{X}_s & q > q_0 \end{cases} \quad (7)$$

$$p(X_s) = \frac{\tau^\alpha(X_s) \cdot \eta^\beta(X_s)}{\sum_{X_s, X_u \in P} \tau^\alpha(X_u) \cdot \eta^\beta(X_u)} \quad (8)$$

Where $q \in [0, 1]$ is even-distributed random number, $q_0 \in [0, 1]$ is probability parameter, P is the set of occupied points for ant in unit space, \tilde{X}_s is the selected target location as per formula (8); α is the update parameter of pheromone, β is the update parameter of visibility.

- (2) Realize the movement of ant towards target position via quantum revolution door
After ant selected the movement, its movement process may be realized by changing the phase of qubit it brought for quantum revolution door. In unit space, suppose the current position for ant at time t is P_i , selected target position is P_s , argument increment of qubit at P_i is

$$\Delta\phi_{ij}^t = \sigma(\phi_{sj} - \phi_{ij}) \times rand_j \quad (9)$$

$$\Delta\phi_{ij}^{t+1} = \begin{cases} \Delta\phi_{ij}^t + 2\pi & \Delta\phi_{ij}^t < -\pi \\ \Delta\phi_{ij}^t & -\pi \leq \Delta\phi_{ij}^t \leq \pi \\ \Delta\phi_{ij}^t - 2\pi & \Delta\phi_{ij}^t > \pi \end{cases} \quad (10)$$

$$\Delta\theta_{ij}^t = \sigma(\theta_{sj} - \theta_{ij}) \times rand_j \quad (11)$$

$$\Delta\theta_{ij}^{t+1} = \begin{cases} \Delta\theta_{ij}^t + \pi & \Delta\theta_{ij}^t < -\pi/2 \\ \Delta\theta_{ij}^t & -\pi/2 \leq \Delta\theta_{ij}^t \leq \pi/2 \\ \Delta\theta_{ij}^t - \pi & \Delta\theta_{ij}^t > \pi/2 \end{cases} \quad (12)$$

Where $rand_j$ is the random number in between $[-1, 1]$; $\sigma \in [1, 2]$ is crawl speed of ant, here $\sigma = 1.75$.

Update of probability amplitude of qubit based on revolution door

$$U = \begin{bmatrix} \cos \Delta\phi_{ij}^{t+1} \cos \Delta\theta_{ij}^{t+1} & -\sin \Delta\phi_{ij}^{t+1} \cos \Delta\theta_{ij}^{t+1} & \sin \Delta\theta_{ij}^{t+1} \cos(\phi_{ij}^t + \Delta\phi_{ij}^{t+1}) \\ \sin \Delta\phi_{ij}^{t+1} \cos \Delta\theta_{ij}^{t+1} & \cos \Delta\phi_{ij}^{t+1} \cos \Delta\theta_{ij}^{t+1} & \sin \Delta\theta_{ij}^{t+1} \sin(\phi_{ij}^t + \Delta\phi_{ij}^{t+1}) \\ -\sin \Delta\phi_{ij}^{t+1} & -\tan(\phi_{ij}^t/2) \sin \Delta\theta_{ij}^{t+1} & \cos \Delta\theta_{ij}^{t+1} \end{bmatrix} \quad (13)$$

$$\begin{bmatrix} \cos \phi_{ij}^{t+1} \sin \theta_{ij}^{t+1} \\ \sin \phi_{ij}^{t+1} \sin \theta_{ij}^{t+1} \\ \cos \theta_{ij}^{t+1} \end{bmatrix} = U \begin{bmatrix} \cos \phi_{ij}^t \sin \theta_{ij}^t \\ \sin \phi_{ij}^t \sin \theta_{ij}^t \\ \cos \theta_{ij}^t \end{bmatrix} = \begin{bmatrix} \cos(\phi_{ij}^t + \Delta\phi_{ij}^{t+1}) \sin(\theta_{ij}^t + \Delta\theta_{ij}^{t+1}) \\ \sin(\phi_{ij}^t + \Delta\phi_{ij}^{t+1}) \sin(\theta_{ij}^t + \Delta\theta_{ij}^{t+1}) \\ \cos(\theta_{ij}^t + \Delta\theta_{ij}^{t+1}) \end{bmatrix} \quad (14)$$

Apparently, U-gate may rotate the phase of qubit by $\Delta\phi_{ij}^{t+1}$ and $\Delta\theta_{ij}^{t+1}$.

(3) Variation treatment

In BQACA, variation process of ant is achieved by adopting the following V gate

$$V = \begin{bmatrix} 0 & \cot \theta & 0 \\ \cot \theta & 0 & 0 \\ 0 & 0 & \tan \theta \end{bmatrix} \tag{15}$$

$$V \begin{bmatrix} \cos \phi \sin \theta \\ \sin \phi \sin \theta \\ \cos \theta \end{bmatrix} = \begin{bmatrix} \cos(\pi/2 - \phi) \sin(\pi/2 - \theta) \\ \sin(\pi/2 - \phi) \sin(\pi/2 - \theta) \\ \cos(\pi/2 - \theta) \end{bmatrix} \tag{16}$$

Command variable probability as P_m , give a random number rand between (0,1) for each ant, if $rand < P_m$, select several qubit in ant randomly, V gate is used to achieve the rotation of quantum phase along Bloch sphere, and the optimal position of itself of memory remains unchanged.

(4) Update rules for strength and visibility of pheromone

When the ant completes a traverse, mapping the current position from unit space to solution space of optimal problem, calculating fitness function, then update the strength and visibility of pheromone of current position.

$$\begin{cases} \tau(X_i) = (1 - \rho)\tau(X_i) + \rho\Delta\tau(X_i) \\ \Delta\tau(X_i) = Qfit(X_i) \end{cases} \tag{17}$$

$$\eta(X_i) = fit(X_i) \tag{18}$$

Where $(1 - \rho) \in [0, 1]$ is the volatile coefficient of pheromone, Q is the enhancement coefficient of pheromone.

2.5 Discription of BQACA

The implementation process of BQACA algorithm for function extreme value and the traveling salesman problem is basically the same, where the difference only exists in the definition of solution space transform and adaptive function. Following is the case of function extreme value problem, which implemented as below:

- Step1: proposed parameters as the number of ants and maximum number of iterations, randomly give the ant initial position according to (1).
- Step2: transform solution space according to (2), and computing each ant's adaptive function in accordance with (5) and (6). in the light of (17) and (18), renew strength and visibility of pheromone.
- Step3: select moving target for each ant in the colony according to (7) and (8), then achieve to help ants moving by using quantum revolving doors in light of (10), (12) and (14).
- Step4: for each ant's variation probability, achieve the variation of ants' position by quantum door in light of (16).
- Step5: Back to Step2 and conduct cycle calculations until the convergence criteria are met or the maximum number of iterations is reached.

3 Simulation Experiment

In order to verify the effectiveness and feasibility of BQACA algorithm, Function Extreme-value Problem and Traveling Salesman Problem are selected for test, then compared with Common Genetic Algorithm (CGA), Common Ant Colony Algorithm (CACA) respectively. The simulated program is achieved by programming MATLAB 2009a, and the test results show as Intel CoreI i5 3.2GHz in CPU, RAM is gained from the running PC as 2.8GB.

3.1 Function Extreme-Value Problem

Let choose two function extreme-value problem to verify the performance of BQACA algorithm.

Function Shaffer's F6:

$$\max f_1(x, y) = 0.5 - \frac{\sin^2 \sqrt{x^2 + y^2} - 0.5}{(1 + 0.001(x^2 + y^2))^2} \quad (19)$$

Where $x, y \in [-5.12, 5.12]$, the function has infinite local maximum points, only one of them, (0,0) is the maximum in global, values 1. The numeric area of independent variable is (-5.12,5.12). Algorithm convergence is regarded when optimal result greater than 0.995.

Function Shubert:

$$\min f_2(x, y) = \left\{ \sum_{i=1}^5 i \cos[(i+1)x + i] \right\} \left\{ \sum_{i=1}^5 i \cos[(i+1)y + i] \right\} + 0.5[(x + 1.42513)^2 + (y + 0.80032)^2] \quad (20)$$

Where $x, y \in [-10, 10]$, the function has 760 local minimum points, only one of them, (-1.42513,-0.80032) is the least in global, global minimum is -186.73090882259. The local minimum easily fallen of function is -186.34. Algorithm convergence is regarded when optimal result less than -186.34.

Algorithm parameter: the limit algebra of each algorithm is 500, and scale of population is 20, terminal conditions all meet the convergence requirement of algorithm or up to the limit algebra. CGA algorithm parameter: each variable is described by 20 binary bits, possibility of crossover $P_c = 0.8$, probability of variation $P_m = 0.2$; CACA and BQACA algorithm: probability parameter $q_0 = 0.5$, volatile coefficient $1 - \rho = 0.05$, update parameter of pheromone $\alpha = 1$, update of visibility $\beta = 5$, enhancement coefficient of pheromone $Q = 10$, variable probability $P_m = 0.05$, crawl speed of ant $\sigma = 1.75$. Take 50 times of experimental data for each case; see comparison of optimizing results in table 1.

Analyzed from table 1 comprehensively, times of convergence in BQACA algorithm hits the most, while the average step gets the least, and with the best optimal result. In BQACA, 3 gene chain decode schemes are used to boost optimizing ability, while variation operator may prevent algorithm from falling into local optimal solution, so as to improve the optimizing performance of BQACA. Conclusion could be drawn through analyzing results: the QACA proposed in this paper in solving the function extreme-value problem, its convergence speed and global optimizing ability are superior to ordinary genetic algorithm and normal ACA.

Table 1. Test Function Calculation Results

Function name	Algorithm	Optimal solution	Worst solution	Mean value	Standard deviation	Convergence times	Mean time	Average step
$f_1(x, y)$	CGA	0.9995	0.9903	0.9911	0.0022	6	0.5507	464.56
	CACA	0.9999	0.9903	0.9939	0.0038	24	0.1871	373.12
	BQACA	1.0000	0.9903	0.9960	0.0029	37	0.2409	265.82
$f_2(x, y)$	CGA	-184.3660	-49.0466	-119.0155	36.1926	0	0.6029	500.00
	CACA	-186.7157	-147.2586	-181.6533	9.5597	19	0.1818	365.92
	BQACA	-186.7180	-184.0693	-186.3127	0.4472	28	0.3137	324.04

3.2 Traveling Salesman Problem

Take symmetric-distance TSP as example, Oliver30 and Att48 were selected from TSPLIB database as cases to verify the performance of the algorithm.

Algorithm parameters: each algorithm prescribes a limit to algebra of 100 and population of 50. in CGA algorithm, integer encoding is adopted; in CACA algorithm, integer encoding adopted as well, other parameters with the function extreme value problem; In BQACA algorithm, transfer factor parameter = 1, stimulating factor parameter = 5, other parameters with the function extreme value problem. Then in every case, 20 experimental data would be taken, and table 2 shows contrast of optimization results. Figure 1, the best solution of Oliver30 quantum ant colony algorithm, the total distance for 424. Figure 2, the best solution of Att48 quantum ant colony algorithm, the total distance is 34596. Statistics analysis of data based on table 2: from time perspective, the average time of BQACA algorithm is shortest, secondly CGA algorithm, CACA algorithm followed; from steps perspective, the convergence rate of BQACA algorithm is more preferable than the other two algorithms; from the perspective of calculation result, the optimal solution of BQACA algorithm can reach a ideal resolution recommended by TSPLIB database when urban scale is small, while whose optimal

Table 2. TSP Calculation Results

Test library	algorithm	Optimal solution	Worst solution	Mean value	Standard deviation	Mean time	Average step
Oliver30	CGA	482.6800	639.5484	570.9215	37.2503	2.0194	100.00
	CACA	423.7406	429.7853	426.7401	1.4871	3.1717	54.10
	BQACA	423.7406	434.8476	430.2515	3.0685	1.6273	43.55
Att48	CGA	47841.1992	61683.1032	54890.0603	4405.2560	3.2348	100.00
	CACA	34562.4143	35979.9099	35374.9823	413.1165	10.0784	65.0000
	BQACA	34596.4450	37458.9908	36518.7123	603.2932	2.1559	45.90

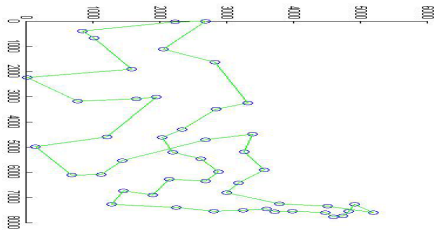


Fig. 1. Time-series of $x(t)$ evolved in system (??) with $\omega = \pi/2$.

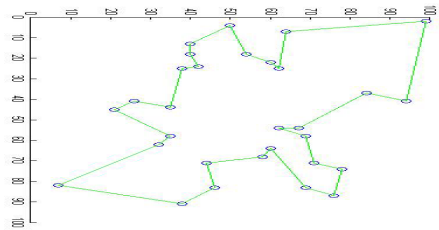


Fig. 2. Time-series of $y_1(t)$ evolved in system (??) with $\omega = \pi/2$.

solution can approach to that of CACA algorithm when the urban scale is large. To sum up, BQACA algorithm in this paper is feasible and effective.

4 Conclusions

Combined with quantum computing and ant colony algorithm, this paper proposes a new Quantum Ant Colony Algorithm which based on Bloch spherical coordinate. The algorithm opened up from the view of quantum computing to regulate the rules of moving ants, and for different optimization problems, various solution space transformational models and fitness functions are planned. The algorithm general idea unique and has wide universality. Research results show that the new algorithm owns particular practical utility which could improve efficiency and accuracy. Compared with traditional intelligent algorithm of CGA and ACA, BQACA has stronger searching ability and higher efficiency, and is applied to complex function optimization and combinatorial optimization problems. At the same time, as a new optimization algorithm, BQACA has a lot to be improved; further study should be followed up in the future.

References

1. Dorigo, M., Maniezzo, V., Colormi, A.: The Ant System: Optimization by a Colony of Cooperating Agents. *IEEE Trans. on SMC* 26(1), 28–41 (1996)
2. Xiong, W.-Q.: Binary ant Colony Algorithm with congestion control strategy for the 0/1 Multiple Knapsack problems. In: *Proceedings of the 8th World Congress on Intelligent Control and Automation (WCICA)*, pp. 3296–3301 (2010)
3. Piao, C., Han, X., Wu, Y.: Improved ant colony algorithm for solving assignment problem. In: *Proceedings of International Conference on Computer Application and System Modeling* (2010)
4. Xing, L., Chen, Y.: A Knowledge-Based Ant Colony Optimization for Flexible Job Shop Scheduling Problems. *Applied Soft Computing* 10(3), 888–896 (2010)
5. Gambardella, L.M., Montemanni, R.: An Enhanced Ant Colony System for the Sequential Ordering Problem. In: *Proceedings of the 41st Annual Conference Italian Operational Research Society* (2010)
6. Hsiao, Y.T.: Computer network load-balancing and routing by ant colony optimization. In: *Proceedings of the 12th IEEE International Conference on Networks*, vol. 1, pp. 313–318 (2004)
7. Gu, Q.H., Jing, S.G.: Study on Vehicle Routing and Scheduling Problems in Underground Mine Based on Adaptively ACA. *Applied Mechanics and Materials* 157, 1293–1296 (2012)
8. Gomez, J.F., Khodr, H.M., De Oliveira, P.M., et al.: Ant colony system algorithm for the planning of primary distribution circuits. *IEEE Trnas. on Power Systems* 19(2), 996–1004 (2004)
9. Yu, Y.Z., et al.: Regulation of PID Controller Parameters Based on Ant Colony Optimization Algorithm in Bending Control System. *Applied Mechanics and Materials* 128-129, 205 (2011)
10. Narayanan, A., Moore, M.: Quantum-inspired genetic algorithms. In: *Proceeding of IEEE International Conference on Evolutionary Computation*, pp. 61–66 (1996)
11. Feng, A.-H., Su, H.-S.: Improved Quantum Genetic Algorithm and Its Application. *Computer Engineering* 37(5), 199–201 (2011)

Energy Efficient VM Placement Heuristic Algorithms Comparison for Cloud with Multidimensional Resources

Dailin Jiang, Peijie Huang^{*}, Piyuan Lin, and Jiacheng Jiang

College of Informatics, South China Agricultural University,
Guangzhou 510642, Guangdong, China
pjhuang@scau.edu.cn

Abstract. Cloud computing provides user utility-oriented IT services, yet accompanied with huge energy consuming, which contributes to the high operational cost as well as CO₂ emission. Making Cloud computing energy efficient can lead to a better tradeoff between profit and environmental impact. In this paper, we formulate the energy efficient VM placement problem in Cloud architecture with multidimensional resources and introduce the objective of this problem. Heuristic algorithms including three traditional local search algorithms and generic algorithm (GA) are presented to provide possible optimized solution. We conduct experiments based on Cloudsim. The result shows that GA sometimes provide the best solution, but with poor stabability. Although the BF provide neither the best nor the worst solution most of time, it have the best stabability.

Keywords: Cloud computing, energy efficient, VM placement, heuristic algorithm.

1 Introduction

Recent year, the rapid growth in demand for computational power driven by modern service applications led to the proliferation of Cloud computing [1], resulting in the establishment of large-scale data centers consuming enormous energy. High-energy consumption not only translates to high-energy cost reducing the profit of Cloud providers, but also high carbon emissions that are not environmentally sustainable [2]. Power has become one of the major limiting factors for a data center [3]. However, the reason for this extremely high-energy consumption is not just lies in the amount of computing resources used and the power inefficiency of hardware, but rather lies in the inefficient usage of these resources. Many data centers often operate at low utilization. Data collected from more than 5000 production servers over a six-month period showed that servers operate only at 10-50% of their full capacity most of the time [4]. But even at a very low load, such as 10% CPU utilization, the power consumed of a server is over 50% of the peak power, because the power needed to run the OS and to maintain hardware peripherals is not negligible [5]. Similarly, if the disk, network, etc. is the performance bottleneck, the idle power wastage in other

^{*} Corresponding author.

resources goes up. There is a need to shift the focus of the Cloud data center from optimizing resource management for pure performance to optimizing them for power and energy efficiency, while meeting performance guarantees.

In Cloud, there are potentially two types of VM placement decisions to be made in VM consolidation: (1) initial placement and (2) migration (and/or resizing) of VMs over time [6]. This paper formulates the problem of initial placement optimization and introduces four heuristic algorithms, FF (First Fit), NF (Next Fit), BF (Best Fit), and generic algorithm to provide possible optimized solution. Then we conduct our experiment on an extended simulator based on Cloudsim [7], which supports simulating major components of a host.

The remainder of this paper is organized as follows. In the next section, we will introduce how we formulate the VM placement-optimizing problem. The details of heuristic algorithms are presented in Section 3. Section 4 introduces how we conduct our experiments and compare the results of the experiments. Finally, Section 5 lists some conclusions and discusses some areas of future research.

2 Problem Formulation

2.1 Research Goal

In modern datacenters, it is essential to take into account the usage of multiple system resources, including CPU, memory, disk storage and network interfaces, etc, in the energy efficient VM consolidation. However, most of the existing literature only focuses on managing power consumption and efficient usage of CPU [2, 8]. The energy model of processors can be derived from the power consumption model in Complementary Metal-Oxide Semiconductor (CMOS) logic circuits given by $P = \alpha f^3 + \beta$, where f is the frequency of the processor, and some researches pointed out that there exists a linear model between utilization of the processor and the energy [9]. However, as far as we are concerned, we have not found any research that indicates there an energy model of other components (RAM, disk, etc.) can be easily implemented in the simulator. That is why we shifted our goal from calculating the energy consumption directly to minimizing the number of hosts used.

2.2 Multidimensional VM Placement Problem Formulation

The datacenter here represents the set of hosts. Each host is characterized by the CPU performance defined in millions instructions per second (MIPS), amount of RAM, network bandwidth and disk storage. Users submit requests for provisioning of VMs characterized by the corresponding resources. One can view this problem as a multi-dimensional bin-packing problem with differently sized bins [3]. The hosts are bins with each resource being one dimension of the bin, and balls represent VMs.

The multidimensional VM placement problem is formulated as follows:

Given n VMs, each of which is characterized by a vector w , which represents the “demand” of its object, and m hosts, each of which is characterized by vector y ,

which represents its capacity, and given an objective function f , our goal is to find a solution to allocate the n VMs into m hosts. Both w and y have the same dimensions representing the resources. To compare the dimensions representing the resources, we redefined the “ \leq ” operator to be a binary operator comparing the dimensions representing the resources of the two vector variables.

In mathematical terms, the above problem can be written as:

$$\begin{aligned} \min/ \max Z(W) &= \sum_{i=1}^n f(w_i) \\ \text{st : } \quad &\sum_{i=1}^n w_i x_{ij} \leq y_j, \quad j \in \{1, 2, \dots, m\} \\ &\sum_{j=1}^m x_{ij} = 1, \quad i \in \{1, 2, \dots, n\} \\ &x_{ij} = 0 \quad \text{or} \quad 1, \quad i \in \{1, 2, \dots, n\}, j \in \{1, 2, \dots, m\} \end{aligned} \quad (1)$$

Here, $x_{ij} = 1$ means the VM i is allocated in host j . $x_{ij} = 0$ means VM i is not allocated in host j .

3 Heuristic Algorithms

3.1 Selection of Viable Solutions

Some researches have used heuristic in VM consolidation optimizing [8, 10]. However, there is a lack of horizontal comparison of these heuristic algorithms applying in multidimensional resource provisioning. In this research, we focus on minimizing energy consumption on the basis of ensuring service level agreement (SLA) and quality of service (QoS). That is to say, we find the Pareto optimal solution by applying the lexicographic ordering approach. Practically, with the precondition of allocating the VMs as many as possible, we find the optimal solution by find the solution which consuming the least energy.

3.2 Traditional Local Search Algorithms

Three are traditional local search algorithms, FF, NF, and BF, which only generate one solution throughout its process. The objective of VM allocation mentioned in Subsection 3.1 will have no use in these three algorithms. The three algorithms are based on different strategies. The FF algorithm always tries to allocate the VM from the first host to the last host. That is to say, the probability of whether a VM will be allocated into a host is descending according to its serial number. Nevertheless, the NF algorithm is based on a relatively same-probability strategy. The BF algorithm is based on the greedy strategy.

Algorithm 1: First Fit (FF)

```

1 Input: hostList, vmList Output: allocation of VMs
2 foreach vm in vmList do
3   foreach host in hostList do
4     if (isAllocable(vm,host)) then
5       allocate (host, vm)
6       break
7 return allocation

```

Algorithm 2: Next Fit (NF)

```

1 Input: hostList, vmList Output: allocation of VMs
2 currentHost←hostList[0]
3 foreach vm in vmList do
4   for host←currentHost in hostList do
5     if (isAllocable(vm,host)) then
6       allocate (vm,host)
7       break
8     else
9       currentHost++
10 return allocation

```

Algorithm 3: Best Fit (BF)

```

1 Input: hostList, vmList Output: allocation of VMs
2 foreach vm in vmList do
3   foreach host in hostList do
4     if (isAllocable(vm,host) and host.spare<
bestHost.spare) then
5       bestHost←host
6     allocation (vm,bestHost)
7 return allocation

```

As we are trying to solve a multidimensional problem, we define the “best” function in BF as follows: without considering the instability of high utilization, we set the optimal point of each host with 100% utilization of each dimension. Then we calculate the Euclidean distance between the optimal point and the current utilization point after allocation of the to-be-allocated VM, and we define the “best” host that is closest to the optimal point.

3.3 Generic Algorithm

Generic algorithm is a search-heuristic algorithm, which will generate viable solutions throughout its process. Moreover, we will implement the approach mentioned in Subsection 3.1 to select the Pareto optimal solution. The algorithm 4.1 describe the main flow of the generic algorithm. Algorithms 4.2 to 4.6 are the pseudo-code of the operators used in algorithm 4.1. The fitness in algorithm 4.6 is the implement of Subsection 3.1.

Algorithm 4.1: The main algorithm of genetic algorithm

```

1 Input: Host_List, VM_List Output: allocation of VMs
2 P(I)←Initialize P(I)
3 Evaluate P(I)
4 while I <= End_Generation do
5     parents←Choose P(I)
6     children←Cross (parents)
7     Mutate (children)
8     Repair (children)
9     P(I)←Initialize P(I++)+children
10    Evaluate P(I)
11 return allocation

```

Here, “Initialize” means initializing the population, which uses create operator to generate individual based on minimum host number. “Choose” means selecting two individuals as parents according to their fitness, the selection algorithm we employed is Roulette pick. “Cross” means generate a child with the DNA provided by parents. “Mutate” simulate the process of mutation. “Repair” make the individual—a solution—to be a viable one. “Evaluate” means calculating the fitness of each individual in the population in the current generation.

Algorithm 4.2: Create operator

```

1 Input: Host_List, VM_List, min_host
2 Output: allocation of VMs
3 Allocate_Host_List←rand min_host hosts in Host_List
4 foreach vm in VM_List do
5     vm.allocatedHost←rand one host in
Allocate_Host_List
6 return allocation

```

Algorithm 4.3: Cross operator

```

1 Input: parent1, parent2 Output: children1, children2
2 children1←parent1
3 children2←parent2
4 host1←rand one host in parent1.Host_List
5 host2←rand one host in parent2.Host_List
6 foreach vm in parent2.VM_List do
7     if (vm.allocatedHost==host2) then
8         children1.vm. allocatedHost←host2
9 Variation (children1)
10 Repair (children1)
11 foreach vm in parent1.VM_List do
12     if (vm.allocatedHost==host1) then
13         children2.vm. allocatedHost←host1
14 Variation (children2)
15 Repair (children2)
16 return children1, children2

```


Algorithm 4.4: Mutate operator

```

1  foreach host in Allocate_Host_List do
2    if(rand()<Host_Variation_Rate) then
3      host2←rand one host in Host_List
4      foreach vm in VM_List do
5        if(vm.allocatedHost==host) then
6          vm.allocatedHost←host2
7  foreach vm in VM_List do
8    if(rand()<VM_Variation_Rate) then
9      host2←rand one host in Host_List
10   vm.allocateHost←host2

```

Algorithm 4.5: Repair operator

```

1  Calculate every host's used resources
2  foreach vm in VM_List do
3    if (vm.allocateHost.used>
4     vm.allocatedHost.resources) then
5     remove (vm)
6  foreach vm in VM_List do
7    if (vm not allocated) then
8     FF (vm,Host_List)

```

Algorithm 4.6: Evaluate operator

```

1  Input: P(I)   Output: Fitness of P(I)
2  Calculate the Euclidean distance of each host in P(I)
and sort by Euclidean distance
3  Calculate the used number of hosts in P(I) and sort
by the used number of hosts
4  return Fitness of P(I)

```

4 Heuristic Algorithms Comparison

4.1 Experiment Setup

As our targeted system is generic Cloud computing environment, it is essential to evaluate it on a datacenter infrastructure in certain large scale. However, it is relatively difficult to replica large-scale experiments. Therefore, we conduct our experiment on CloudSim. In addition, we extended this simulator to better support multidimensional resource provisioning simulation.

We have simulated a datacenter composed with 50 homogeneous hosts. Each host is modeled to have 100000 MIPS, 16GB of RAM, 1TB Storage, and 100GB bandwidth. Then, we have simulated 50 to 370 VMs of the initial placement. Each VMs is requiring 250, 500 ... or 2500 MIPS, 512 MB, 1GB ... or 4GB RAM, 1, 2 ... or 10GB of Storage, and 1, 2 ... or 10GB bandwidth. We start our simulation by creating 50 VMs. Then we add 5 to the previous amount of the VMs for another new simulation until the amount of VMs reaches 370. The main parameters of generic algorithm are set as follow: End_Generation=5000, Population_Size=100, Host_Variation_Rate=0.01, and VM_Variation_Rate=0.02.

4.2 Experiment Result Analysis

By applying the four heuristic algorithms and selection strategy mentioned in Section 3, we have the result showed in Fig. 1. Fig. 1 (a) shows the algorithm stability of 4 heuristic algorithms. Fig. 1 (b) is the result averaged over 10 independently runs of the experiments, which is used to show the performance comparison, i.e. the number of hosts used.

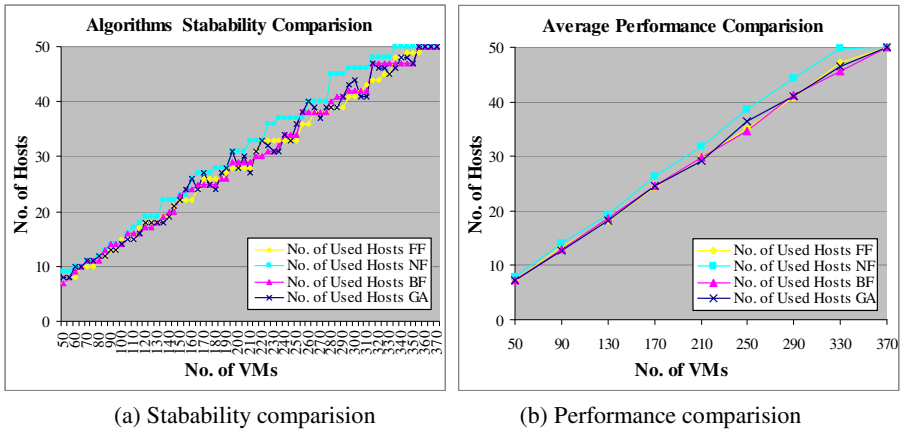


Fig. 1. Heuristic algorithms comparison

As we can see from Fig. 1 (a), GA sometimes provide the best solution among the 4 algorithms, but its stability is worst, which may be affected by the population size, amount of generations, etc. The corresponding R^2 of FF, NF, BF and GA are 0.9940, 0.9890, 0.9959, 0.9885, respectively. Although the BF provide neither the best nor the worst solution most of time, it have the best stability.

Fig. 1 (b) shows that BF has the best overall performance. Although GA gets the best performance when the number of VM is relatively small, it cannot maintain good performance with large number of VM, which may cause by the limited amount of generations used in the experiment. The NF has the worst performance, which is understandable for its relatively same-probably strategy, which causing discrete VM allocation, and thus more hosts in used.

5 Conclusions and Future Work

In this paper, we have carefully formulated one sub-problem, initial placement of VM consolidation optimizing, in Cloud computing model with multidimensional resources. Also we have presented 4 heuristic algorithms to provide solution of initial placement of VM consolidation optimizing, and compared the performances of these algorithms.

In our future work, we will try to do some experiments on the real Cloud computing environments based on XenServer to research the impacts of live migration especially on

energy consumption. Being able to evaluate the impacts of live migration, we are able to extend our works to support continuous optimization.

Moreover, service provider may consider one or more objectives (e.g. minimizing the rate of violating SLAs, maximizing the performance, minimizing the cost and energy, etc.) in VM consolidation. More specifically, the objective function f in equ. (1) is actually a set of objective function $\{\mu_1(x), \mu_2(x), \dots, \mu_n(x)\}$. One solution represents a Pareto point. The “min/max $Z(W)$ ” is actually trying to find Pareto optimal solutions. By considering multiple objectives, it is possible to optimize VM consolidation according to different orientations such as considering tradeoff between CO₂ emission and profits.

Acknowledgments. This work is supported by the Industry-Education-Research Cooperation Project of Guangdong Province and Ministry of Education under Grant No. 2011A090200072, the Foundation for Distinguished Young Talents in Higher Education of Guangdong, China under Grant No. LYM09034, the Soft Science Research Project of Guangdong Province under Grant No. 2011B070400009, and the College Students Innovation Experiment Project of Guangdong Province under Grant No. 1056411060.

References

1. Buyya, R., Yeo, C.S., Venugopal, S., et al.: Cloud computing and emerging IT platforms: vision, hype, and reality for delivering computing as the 5th utility. *Future Gener. Comp. Sy.* 25, 599–616 (2009)
2. Garg, S.K., Yeo, C.S., Anandasivam, A., Buyya, R.: Environment-conscious scheduling of HPC applications on distributed Cloud-oriented data centers. *J. Parallel Distr. Com.* 71, 732–749 (2011)
3. Liao, X.F., Jin, H., Liu, H.K.: Towards a green cluster through dynamic remapping of virtual machines. *Future Gener. Comp. Sy.* 28, 469–477 (2012)
4. Barroso, L.A., Holzle, U.: The case for energy-proportional computing. *Computer* 40, 33–37 (2007)
5. Chen, G., He, W., Liu, J., et al.: Energy-aware server provisioning and load dispatching for connection-intensive internet services. In: 5th USENIX Symposium on Networked Systems Design and Implementation (USENIX NSDI 2008), pp. 337–350 (2008)
6. Mills, K., Filliben, J., Dabrowski, C.: Comparing VM-placement algorithms for on-demand Clouds. In: 3rd IEEE International Conference on Cloud Computing Technology and Science (CloudCom 2011), pp. 91–98 (2011)
7. Calheiros, R.N., Ranjan, R., Beloglazov, A., et al.: CloudSim: a toolkit for modeling and simulation of Cloud computing environments and evaluation of resource provisioning algorithms. *Software Pract. Exper.* 41, 23–50 (2010)
8. Beloglazov, A., Abawajy, J., Buyya, R.: Energy-aware resource allocation heuristics for efficient management of data centers for Cloud computing. *Future Gener. Comp. Sy.* 28, 755–768 (2012)
9. Kusic, D., Kephart, J.O., Hanson, J.E., et al.: Power and performance management of virtualized computing environments via lookahead control. *Cluster Comput.* 12, 1–15 (2009)
10. Kessaci, Y., Melab, N., Talbi, E.: A pareto-based GA for scheduling HPC applications on distributed Cloud infrastructures. In: 2011 International Conference on High Performance Computing and Simulation (HPCS 2011), pp. 456–462 (2011)

A Python Based 4D Visualization Environment

Lin Jing¹, Xipei Huang¹, Yiwen Zhong¹, Yin Wu², and Hui Zhang^{3,*}

¹ College of Computer and Information Science,
Fujian Agriculture and Forestry University, Fuzhou, China

² School of Informatics and Computing,
Indiana University, Bloomington, USA

³ Pervasive Technology Institute,
Indiana University, Indianapolis, USA
hui Zhang@iu.edu

Abstract. One of the challenges in 4D math visualization is to develop an interactive and integrated computational environment to quick-prototype, simulate, and experiment with the abstract mathematical concepts. We have investigated several areas in 4D visualization including 4D surface rendering, various user interface elements to manipulate mathematical objects in the higher-dimensional space, and physically based modeling of cloth-like 4D objects to understand the math phenomena in the fourth dimension. In this paper, we present one such Python based 4D visualization environment that achieves a high level of integration between 4D math, physics computation and interactive visualization.

Keywords: Math visualization, 4D visualization, Python.

1 Introduction

Mathematical visualization is the art of creating an interactive experience with abstract mathematical objects and concepts [1]. Typical geometric problems of interest to mathematical visualization applications involve both static structures, such as real or complex manifolds, and changing structures requiring computer animation, such as sphere eversion. In practice, much emphasis has been on manifolds of dimension two or three embedded in three or four-dimensional space. Due to the practical limitations of holistic human spatial perception, it is nearly impossible to construct useful physical models for visualization purposes. For example, Figure 1(a) shows a mathematically accurate 3D plastic model of a 2D manifold embedded in four dimensions [4], which may contain no abstract mathematical information at all to someone unfamiliar with its mathematics.

The limitations of using physical models to learn abstract math concepts can be overcome by using computer based interactive systems. For example, the recent advent of high-performance interactive computer graphics systems sparked a renewed interest in visual mathematics in various interactive forms. There are many tasks that can be performed better with an interactive visualization system. For example, the problem of turning a sphere inside out without tearing or creasing, known in mathematics as a regular homotopy that *everts the sphere*, is a classic puzzle that has been solved in many

* Corresponding author.

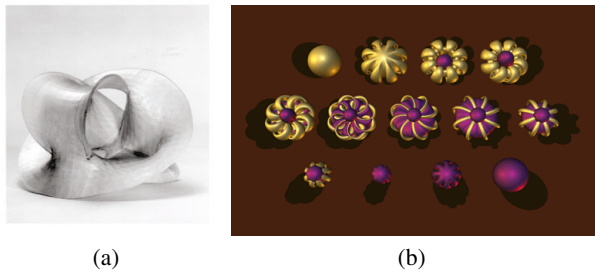


Fig. 1. (a): Examples of physical models representing a highly self-intersecting surface, constructed by 3D projection from the four-dimensional mathematical description (Model courtesy of Stewart Dickson.) (b): Thurston's method for everting the sphere.

ways since Smale first proved it must be possible. Even though it is mathematically possible to evert the sphere, there is not such a physical model that allows the eversion to be achieved in a natural and physically intuitive manner and for the broad public to appreciate (simply because there is no physical materials allowing self-intersections in our everyday life). The visual solution of a particular sphere eversion was due to Thurston (see e.g., Figure 1(b)), implemented in computer graphics at the Geometry Center, and discussed in detail in the file *Outside In* [3]. Other examples include many interesting mathematical phenomena that have not only standard 3D geometric concepts but their higher-dimensional counterparts [6]. Although mathematical visualization helps one to understand and explore mathematical phenomena in four-dimensional space intuitively, building such software for math education has proved very challenging due to both the need for considerable domain knowledge and insight of the math data plus the need for considerable technical expertise for the undertaking.

2 Background of 4D Visualization

During the last several years, the research efforts on 4D visualization at Indiana University have examined a family of methods of physically constructing and interacting with mathematical objects embedded in four-dimensional space. We presented methods for the computer construction, multimodal exploration and interactive manipulation of a wide variety of mathematical 4D objects. The basic problem is that, just as 2D shadows of 3D curves lose structure where lines cross, 3D graphics projections of smooth 4D topological surfaces are *interrupted* where one surface intersects another. Furthermore, if one attempts to trace *real* knotted ropes or a plastic model of self-intersecting surfaces with fingertip, one inevitably collides with parts of the physical artifact. In our research, we exploited the free motion of a computer-based haptic probe to support a continuous motion that follows the *local continuity* of the object being explored. The proposed haptics techniques to explore mathematic objects in 4D has also led to a new interaction technique without a real haptic interface ([7], [6]).

By combining graphics and various interaction techniques, we have found a way to enhance our experience of interacting with 4D object by producing a reduced-dimension 3D tool for manipulating objects embedded in 4D. By physically modeling the correct

properties of 4D surfaces, their bending forces, and their collisions in the 3D haptic controller interface, we can support full-featured computer-based exploration of 4D mathematical objects in a manner that is otherwise far beyond the experience accessible to human beings.

3 A Python Based 4D Visualization Environment

Different from many other general graphics applications, 4D math visualization systems as such often require customized hybrid approaches and there are a number of considerations when choosing a platform or language for the undertaking. For example, the undertaking usually spans several fields including computer based visualization, graphics, numerical computation, and physically based modeling. In our experiments, we have had very good success using Python. Python is a very high-level, interpreted language that is perfect for mathematician's quick experimentation and rapid prototyping. A number of freely available contributed libraries for Python make it a wonderful platform language for 4D visualization [5]:

- **Math and Physics aware:** Numerical libraries (i.e., *ScientificPython*, *NumPy* and *SciPy*) have added array and matrix support to the Python language. High-performance numerical routines, matrix and multidimensional array operations, optimizations, and linear algebra are now seamlessly achievable with Python.
- **Graphics-supporting:** OpenGL package and VPython module allow access to the underlying graphics hardware for fast 3D rendering, and easy creations of navigable 3D scene graph and animation (even for those with limited programming experience).
- **Interactive:** Python is designed for quick experimentation and rapid prototyping. This is one of the key facets we are looking for when building math visualization environment.
- **Open environment:** Python is open source, easy-to-integrate, and therefore able to interact with many other tools.

As shown in Figure 2, our 4D visualization environment is based on a wide variety of contributed libraries developed on Python. In the next section, we will describe the families of modules we developed to manipulate and render the projected images of mathematical objects embedded in four dimensions.

4 Implementation Methods and Modules

4.1 Geom4D Module

Vector Operations and Points in 4D. For the most part, vector operations in four space are simple extensions of their three-space counterparts. For example, computing the addition of two four-vectors is a matter of forming a resultant vector whose components are the pairwise sums of the coordinates of the two operand vectors. In the same fashion, subtraction, scaling, and dot-products are all simple extensions of their more common

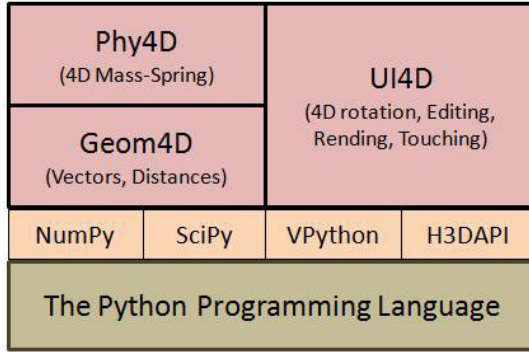


Fig. 2. Architecture and modules overview of our 4D visualization system

three-vector counterparts. In addition, operations between four-space points and vectors are also simple extensions of the more common three-space points and vectors. For example, computing the four-vector difference of four-space points is a simple matter of subtracting pairwise coordinates of the two points to yield the four coordinates of the resulting four-vector. For completeness, the equations of the more common four-space vector operations follow.

$$\begin{aligned}
 U + V &= [U_0 + V_0, U_1 + V_1, U_2 + V_2, U_3 + V_3] \\
 U - V &= [U_0 - V_0, U_1 - V_1, U_2 - V_2, U_3 - V_3] \\
 kV &= [kU_0, kU_1, kU_2, kU_3] \\
 U \cdot V &= U_0V_0 + U_1V_1 + U_2V_2 + U_3V_3
 \end{aligned}$$

Distances in 4D. To understand the non-intuitive mechanisms of 4D collision, let us start with a pair of two-dimensional planes through the origin in four-dimensional space (see Figure 3). The two squares intersect in a single 4D point at the origin. In the three-dimensional projection, although the planes appear to intersect along an entire line due to the annihilation of the w dimension, when the surfaces are 4D depth color-coded, we can see that there is just one pair of points with the same fourth coordinates as well as the same coordinates in the 3D projection. Figure 3(b) illustrates the basic case for 4D collision detection; a 4D depth collision test must be performed along the intersecting lines of the projected images of 4D surfaces. 4D collision occurs if, and only if, one or more pairs of points are located with same fourth coordinate along the intersecting line. Now let 4D surfaces S_A and S_B intersect in the 3D projected image along the line segment L_{se} , from point $P_s = (x_s, y_s, z_s)$ to point $P_e = (x_e, y_e, z_e)$. Suppose the pair of 4D points sharing P_s as shadow points are P_0 on S_A , and Q_0 on S_B ; likewise, we have P_1 on S_A , and Q_1 on S_B sharing P_e as shadow point. Obviously, $P_0, P_1, Q_0,$ and Q_1 can be represented as:

$$\begin{aligned}
 P_0 &= (x_s, y_s, z_s, w_{sA}) \\
 Q_0 &= (x_s, y_s, z_s, w_{sB}) \\
 P_1 &= (x_e, y_e, z_e, w_{eA}) \\
 Q_1 &= (x_e, y_e, z_e, w_{eB}) .
 \end{aligned} \tag{1}$$

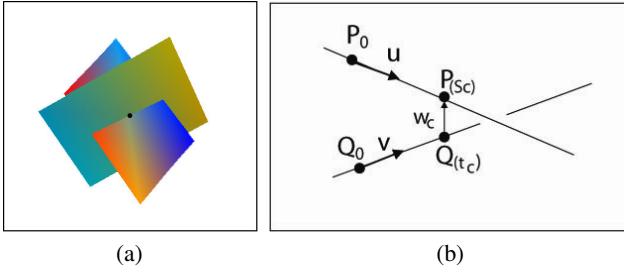


Fig. 3. (a) 4D collision: surfaces are color-coded to indicate their depth in four-space relative to the projection point, so we observe that there is just one pair of points with the same fourth coordinate as well as the same first three coordinates. (b) Closest points between 4D lines.

The next two sections outline two major distance-calculation algorithms: the distance between two 4D line segments (in order to avoid edge-edge collision), and the distance between a point and a triangle embedded in 4D (in order to avoid point-triangle collision).

Closest Points between 4D Line Segments. We first consider two infinite lines L_1 : $\mathbf{P}(s) = \mathbf{P}_0 + s(\mathbf{P}_1 - \mathbf{P}_0) = \mathbf{P}_0 + s\mathbf{u}$ and L_2 : $\mathbf{Q}(t) = \mathbf{Q}_0 + t(\mathbf{Q}_1 - \mathbf{Q}_0) = \mathbf{Q}_0 + t\mathbf{v}$. Let $\mathbf{w}(s, t) = \mathbf{P}(s) - \mathbf{Q}(t)$ be a vector between points on the two lines. We want to find the $\mathbf{w}(s, t)$ that has a minimum length for all s and t . In any N -dimensional space, the two lines L_1 and L_2 are closest at unique points $\mathbf{P}(s_c)$ and $\mathbf{Q}(t_c)$ for which $\mathbf{w}(s_c, t_c)$ attains its minimum length. Also, if L_1 and L_2 are not parallel, then the line segment $\mathbf{P}(s_c) \leftrightarrow \mathbf{Q}(t_c)$ joining the closest points is uniquely perpendicular to both lines at the same time. No other segment between L_1 and L_2 has this property (see Figure 3(b)). That is, the vector $\mathbf{w}_c = \mathbf{w}(s_c, t_c)$ is uniquely perpendicular to the line direction vectors \mathbf{u} and \mathbf{v} , and thus it satisfies the equations:

$$\begin{aligned} \mathbf{u} \cdot \mathbf{w}_c &= 0 \\ \mathbf{v} \cdot \mathbf{w}_c &= 0 . \end{aligned} \tag{2}$$

We can solve these two equations by substituting $\mathbf{w}_c = \mathbf{P}(s_c) - \mathbf{Q}(t_c) = \mathbf{w}_0 + s_c\mathbf{u} - t_c\mathbf{v}$, where $\mathbf{w}_0 = \mathbf{P}_0 - \mathbf{Q}_0$, into each part of Eq. (2) to get two simultaneous linear equations. Then, letting $a = \mathbf{u} \cdot \mathbf{u}$, $b = \mathbf{u} \cdot \mathbf{v}$, $c = \mathbf{v} \cdot \mathbf{v}$, $d = \mathbf{u} \cdot \mathbf{w}_0$, and $e = \mathbf{v} \cdot \mathbf{w}_0$, we solve for s_c and t_c as:

$$s_c = \frac{be - cd}{ac - b^2}, \quad t_c = \frac{ae - bd}{ac - b^2} . \tag{3}$$

Having solved for s_c and t_c , we have the points $\mathbf{P}(s_c)$ and $\mathbf{Q}(t_c)$ where the two lines L_1 and L_2 are closest. Then the distance between them is given by:

$$d(L_1, L_2) = \left| (\mathbf{P}_0 - \mathbf{Q}_0) + \frac{(be - cd)\mathbf{u} - (ae - bd)\mathbf{v}}{ac - b^2} \right| . \tag{4}$$

Now we represent a segment S_1 (between endpoints \mathbf{P}_0 and \mathbf{P}_1) as the points on L_1 : $\mathbf{P}(s) = \mathbf{P}_0 + s(\mathbf{P}_1 - \mathbf{P}_0) = \mathbf{P}_0 + s\mathbf{u}$ with $0 \leq s \leq 1$. Similarly, the segment S_2 on L_2 from

\mathbf{Q}_0 to \mathbf{Q}_1 is given by the points $\mathbf{Q}(t)$ with $0 \leq t \leq 1$. The distance between segments \mathbf{S}_1 and \mathbf{S}_2 may not be the same as the distance between their extended lines \mathbf{L}_1 and \mathbf{L}_2 . The first step in computing a distance involving segments is to get the closest points for the lines they lie on. So, we first compute s_c and t_c for \mathbf{L}_1 and \mathbf{L}_2 , and if these are in the range of the involved segment, then they are also the closest points for them. But if they lie outside the range, then they are not and we have to determine new points that minimize $\mathbf{W}(s,t) = \mathbf{P}(s) - \mathbf{Q}(t)$ over the ranges of interest.

To do this, we first note that minimizing the length of \mathbf{w} is the same as minimizing $|\mathbf{w}|^2 = \mathbf{w} \cdot \mathbf{w} = (\mathbf{w}_0 + s\mathbf{u} - t\mathbf{v}) \cdot (\mathbf{w}_0 + s\mathbf{u} - t\mathbf{v})$ which is a quadratic function of s and t . In fact, this expression defines a paraboloid over the (s,t) -plane with a minimum at $C = (s_c, t_c)$, and which is strictly increasing along rays in the (s,t) -plane that start from C and go in any direction. However, when segments are involved, we need the minimum over a subregion \mathbf{G} of the (s,t) -plane, and the global minimum at C may lie outside of \mathbf{G} . An approach is given by [2], suggesting that the minimum always occurs on the boundary of \mathbf{G} , and in particular, on the part of \mathbf{G} 's boundary that is visible to C . Thus by testing all candidate boundaries, we can compute the closest points between 4D line segments.

Closest Points between 4D Point and Triangle. The problem now is to compute the minimum distance between a point \mathbf{P} and a triangle $\mathbf{T}(s,t) = \mathbf{B} + s\mathbf{E}_0 + t\mathbf{E}_1$ for $(s,t) \in D = (s,t) : s \in [0, 1], t \in [0, 1], s+t \leq 1$. The minimum distance is computed by locating the value $(\bar{s}, \bar{t}) \in D$ corresponding to the point on the triangle closest to \mathbf{P} .

The squared-distance function for any point on the triangle to \mathbf{P} is $Q(s,t) = |\mathbf{T}(s,t) - \mathbf{P}|^2$ for $(s,t) \in D$. The function is quadratic in s and t ,

$$Q(s,t) = as^2 + 2bst + ct^2 + 2ds + 2et + f, \quad (5)$$

where $a = \mathbf{E}_0 \cdot \mathbf{E}_0$, $b = \mathbf{E}_0 \cdot \mathbf{E}_1$, $c = \mathbf{E}_1 \cdot \mathbf{E}_1$, $d = \mathbf{E}_0 \cdot (\mathbf{B} - \mathbf{P})$, $e = \mathbf{E}_1 \cdot (\mathbf{B} - \mathbf{P})$, and $f = (\mathbf{B} - \mathbf{P}) \cdot (\mathbf{B} - \mathbf{P})$. Quadratics are classified by the sign of $ac - b^2$, which here becomes

$$ac - b^2 = (\mathbf{E}_0 \cdot \mathbf{E}_0)(\mathbf{E}_1 \cdot \mathbf{E}_1) - (\mathbf{E}_0 \cdot \mathbf{E}_1)^2 = |\mathbf{E}_0 \times \mathbf{E}_1|^2 > 0. \quad (6)$$

The positivity is based on the assumption that the two edges \mathbf{E}_0 and \mathbf{E}_1 of the triangle are linearly independent, so their cross product is a nonzero vector.

In calculus terms, the goal is to minimize $Q(s,t)$ over D . Since Q is a continuously differentiable function, the minimum occurs either at an interior point of D where the gradient $\nabla Q = 2(as + bt + d, bs + ct + e) = (0,0)$, or at a point on the boundary of D .

4.2 Phy4D Module

In *Phy4D* module, we provide a library to model 4D surfaces with a 4D mass-spring system supporting physical interaction. In practice, we focus on 2-manifold deformable objects embedded in 4D. The *Phy4D* module extends the 3D mass-spring system to four dimensions, as can be used to model the 4D mathematical and physical behavior of 4D topological surfaces. When applying the 4D mass-spring system to 2-manifold deformable objects, we assume that each mass point i is linked to all the others with

(linear) springs of rest length $l_{i,j}^0$ and stiffness $k_{i,j}$. This stiffness is set to zero if the actual model does not contain a spring between mass i and j . We note that our mass-spring system is configured in 4D, so the internal forces exerted by 4D springs are 4D vectors. Once a 4D spring is enabled for dynamics, the fundamental dynamical equation is explicitly integrated across time by the Euler method:

$$\begin{cases} \alpha_{i,j}(t + \Delta t) = \frac{1}{m_{i,j}} F_{i,j}(t) \\ V_{i,j}(t + \Delta t) = V_{i,j}(t) + \Delta t \alpha_{i,j}(t + \Delta t) \\ P_{i,j}(t + \Delta t) = P_{i,j}(t) + \Delta t V_{i,j}(t + \Delta t) \end{cases} \quad (7)$$

5 User Scenarios

Our 4D visualization environment can be used to develop a correct interactive experience with the intuitive nature of unfamiliar 4D geometry. One value of our 4D visualization environment is that by zeroing sub-dimensions in the 4D mass-spring system, we can also simulate mathematical and physical phenomena in 2D and 3D, using deformable curves and manifolds in 2D and 3D. This makes it possible for mathematicians to prototype and experiment with dimensional progress analogies.

One such interesting example is shown in Figure 4(a)-(d). The point and the ring have non-zero geometric information only in 2D, i.e., (x,y) initially, and the point will be trapped inside the ring when exposed to 2D external forces only (both the ring and the point are sharing the same (i.e., zero) 3D depth and 4D depth initially). Interestingly,

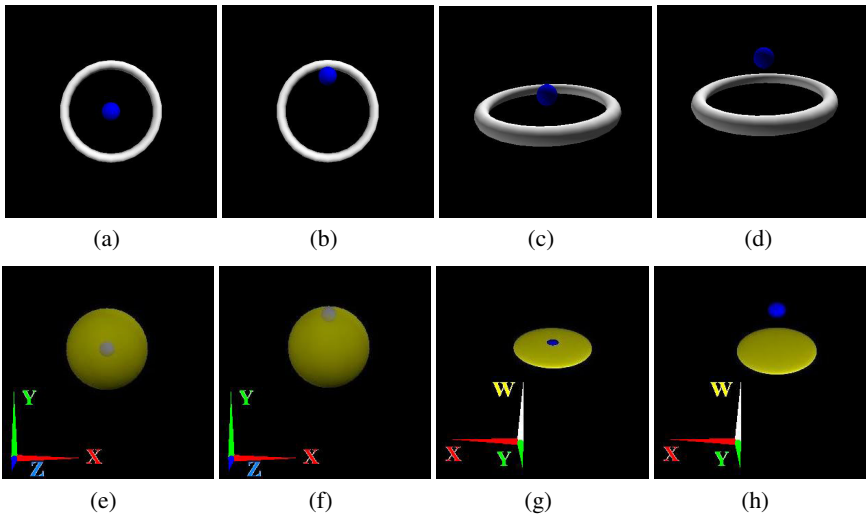


Fig. 4. (a)-(b) A point is trapped inside a ring when embedded in 2D. (c) With a general 3D rotation applied, a 2D force can be applied to the point then to pull the point outside of the ring, which are now embedded in 3D. (e)-(f) A point trapped inside a 2-sphere when embedded in 3D. (g)-(h) With a general 4D rotation applied, the point can escape from the 2-sphere embedded in 4D.

this can be trivially extended to an interesting 4D analogy in Figure 4(e)-(h). A point inside a 2-sphere can't seem to escape when embedded in 3D initially, we all know this is true since in our everyday life we use containers to hold various stuffs (a sealed container is topologically equivalent to a 2-sphere). Just like a point can escape from a ring in 3D, a point can escape from a 2-sphere in 4D. If we apply a general 4D rotation, each vertex will acquire a non-vanishing 4D "eye coordinate" or depth w ; if we again apply 2D forces to the point (e.g., along w direction), we see that we are able to "pull" the point from the 2-sphere in 4D just like we pull the point from a ring in 3D (see Figure 4(g)). The apparent intersections between the point and the 2-sphere are just an illusion in 4D, and a 2-sphere simply can't trap a point any more in 4D (see Figure 4(h)).

6 Conclusion

We have demonstrated a Python-based 4D visualization environment for prototyping and experimenting with mathematical and physical phenomena in 2D, 3D and 4D. A family of computational modules have been developed for generating interactive experience with mathematical objects in higher-dimensional space. By building such a framework based on the Python programming language and the contributed libraries, we provide a simple-to-use and interactive programming environment for mathematicians, research software developers and students to enjoy the excitement of exploring higher-dimensional phenomena.

References

1. Azman Abu, N., Daud Hassan, M., Sahib, S.: Mathematical animations: The art of teaching. In: 31st Annual Frontiers in Education Conference, vol. 3, pp. S1C-10-15 (2001)
2. Eberly, D.: 3D Game Engine Design. Morgan Kaufmann Publisher (2001)
3. Francis, G., Sullivan, J.M.: Visualizing a sphere eversion. *IEEE Transactions on Visualization and Computer Graphics* 10(5), 509-515 (2004)
4. Hanson, A.: A construction for computer visualization of certain complex curves. *Notices of the Amer. Math. Soc.* 41(9), 1156-1163 (1994)
5. Zelle, J.M., Figura, C.: Simple, low-cost stereographics: VR for everyone. *SIGCSE Bull.* 36, 348-352 (2004)
6. Zhang, H., Hanson, A.: Shadow-driven 4D haptic visualization. *IEEE Transactions on Visualization and Computer Graphics* 13(6), 1688-1695 (2007)
7. Zhang, H., Weng, J., Hanson, A.J.: A pseudo-haptic knot diagram interface. In: *VDA 2011: Proceedings of the Conference on Visualization and Data Analysis 2011, San Francisco, CA (January 2011)*

Study of Trustworthiness Measurement and Kernel Modules Accessing Address Space of Any Process

Ce Zhang^{1,2}, Gang Cui¹, Bin Jin², and Liang Wang²

¹ School of Computer, Harbin Institute of Technology, Harbin 150001, China

² School of Computer, Harbin Institute of Technology at Weihai, Weihai 264209, China
zhangce@hitwh.edu.cn, cg@ftcl.hit.edu.cn, jbsumit@126.com,
wangxijue82@gmail.com

Abstract. Trustworthiness measurement is the base and important supporting technology of Trusted Computing. The main objective of trustworthiness measurement is that, how to estimate the trustworthiness of different objects by appropriate policies. In measurement, accessing the address space of measured objects and obtaining the various datum and evidences are considered firstly. Aiming to this problem, this paper presents the primary measurement system architecture, and puts forward three methods of MA(Measurement Agent)in user space invoking MMK(Measurement Module in Kernel)in kernel space. In addition, the principal and realization of accessing a process address space is proposed, including address remapping, switching the CR3 manually and by kernel thread. Finally, three methods are compared qualitatively, and performance consumption is listed by experiment.

Keywords: measurement, kernel module, process, address space, CR3.

1 Introduction

Software trustworthiness is the most important not-function attributes in the software quality characteristics[1]. Reliable and secure measurement is the quantitative characterization for the trustworthiness of software object[2-3], it need to procure kinds of credible evidence in this process[4-5]. Static measurement pay attention to code data integrity of software, but it does not mean that the integrity of the running software[6-9]. Dynamic integrity pay close attention to the integrity change of the software in running[10], corresponding to the dynamic trustworthiness essence and the security of the measurement. In order to obtain credible evidence, it need to access the running software, this paper do the research for the measurement modules of kernel of dynamic trustworthiness measurement and how to invoke measurement module, and how to access the address space of process measured.

2 Primary Measurement System Architecture

First, In order to get the basic framework for measure the process it need to increase the measurement module in the kernel.

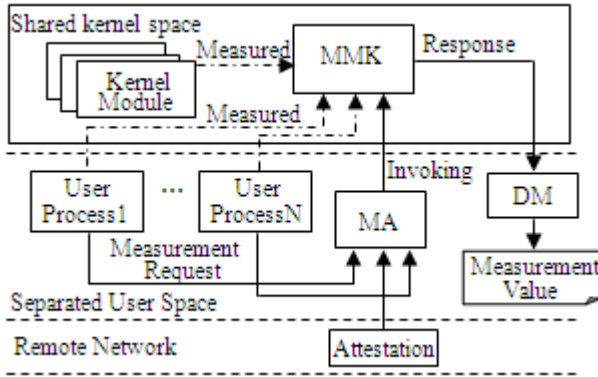


Fig. 1. Primary measurement system architecture

The basic measurement architecture is represented as fig. 1, composing three parts, Measurement Agent (MA), located in Measurement Module in Kernel (MMK) and Decision Module (DM). MA running in the user space, accepting measurement request (MR), send measurement request to MMK, trigger the MMK to measurement. MR can from within the system. it also can be Attestation/Challenger. MA need to make necessary reduction on MR, submitting to MMK in a specific format of the system call way. MMK as system call of MA and executed, running in the kernel space, MMK use a suitable measurement strategy for measuring integrity of user process and other credible attribute or kernel module, the response of measure back to the Decision Module(DM), then it giving the finally measure result ,according to the anticipate requirements.

When the coda of measurement module in kernel execute as a part of kernel, it also has context of current process, the MA in the user space has three methods about how to trigger MMK in the kernel space , three methods are as follows:

Method 1: add a new system call function for the kernel, of course, it need to recompile the kernel, it doesn't work amenity.

Method 2: adopt the devise driven architecture of bogus character, namely, MMK be a devise driven module of character, application of the user mode which is MA call it through the system call ioctl. This method is like the traditional driver development method.

Method 3: according to the shortage in method 2 is not suitable for the transmission of large quantities of data, it still can use "NETLINK" mechanism ,sending measurement request to MMK in the kernel module, MA fall in the kernel mode and execute coda of MMK module which in the kernel, trying to access the user space of target process.

The above three methods in essence are based on the method of system call fall into kernel automatic, but it has big difference in its realization and in the effect to the system performance. Table 1 gives the difference comparison of three methods.

This paper focused on the research of MMK triggered by MA, and of the methods how MMK access the address space of the process in the user mode, in order to be able to provide the necessary support for the measure.

Table 1. Comparisons of MA invoking MMK methods

Method	Method of trapping into kernel	Influence to system performance	Speed	Data throughput
Method 1	System call	bigger	faster	small
Method 2	System call based on character device driver module	bigger	faster	small
Method 3	System call based on "NETLINK" mechanism	little	faster	Bigger

Consider MMK will obtain measure information of various process is larger and it is need to be returned to the MA. So use method 3 here. By the way points out that Linux is a modern operating system in using protection mode, it is unable to directly access the physical memory address. Therefore, it is absolutely unable to converse the virtual address of target process to physical address, then read the target process image content by using these physical address. The accessing address used by CPU is virtual address in the protection mode, the conversion of virtual address to physical address is made by MMU hardware, and this process is transparent for programmer.

This shows, in actually, it is a process (Measurement Agent) access user space of another process (measure target process), to be frank, it is cross-border access. Cross-border access, we all know that must be a fatal exception error page^[11], Is there has some way to achieve the requirements in safely? The answer is sure, for the program itself in the kernel space and it has the privilege level as same as kernel, it can free to access all the kernel space, it also should be able to access any physical memory. We put forward three methods to some degree all of them could meet the requirements.

3 Theory and Realization of Accessing Address Space of Any Process

3.1 Method 1: The Method of Address Remapping

Operating system in protection mode use virtual address when it access memory, when create a process , Linux system mapping virtual address of process to physical memory through the structure of the page directory table, page table, page and offset. To access the address space of other process, after all, in essence, is to access its corresponding physical memory space. So we can suppose that if we can establish a mapping from the virtual address of current process" of MA to target physical address, then should be able to achieve the "cross-border" access. As a result, the target process image corresponding to physical address falls on range in virtual address mapping of the current process of MA's , which is equivalent to access its own address space. Alternatively, the physical address of the target process is mapped to the kernel space, while kernel space is shared among all processes, so that we can assess this section of physical memory. We use the latter means, namely, is about to access the physical page by mapped them to the public kernel space.

However, to achieve this process, there are two major difficulties:

(1) Some pages of the target process may not be in memory. In order to complete re-mapping of physical address to virtual address, all the pages of target process must exist in the memory and can't be swapped out;

(2) How to establish a mapping from physical pages to the kernel address space.

Fortunately, the above two difficulties in the Linux kernel source code can be found in the corresponding solution. The page will be transferred to the memory, can be accomplished through a function `get_user_pages()`; the page is mapped to the kernel space can be `kmap()` function to be complete. These two functions have been exported by `EXPORT_SYMBOL` macro, it can be called directly in the kernel module. It is worth mentioning is that `get_user_pages()` explicitly calls `handle_mm_fault()` function, transferring the swapped out pages into memory, in normal circumstances it is done automatically in the page fault exception handling process. If you are interested to their processes details, can be found the source code in the `<mm/memory.c>`, readers can read on their own.

In addition, a buffer will be needed to store the copied contents, because `kmap()` can only be mapped very small amount of high memory pages to the kernel space at one time, and then as soon as possible to lift the mapping.

Complete procedure function in this method has achievement in `access_remote_vmspace` depicted in Fig.2. , the prototype is as follows:

`static int access_remote_vmspace(struct task_struct *tsk, struct mm_struct *mm, unsigned long addr, void *buf, int len, int write)`. Meaning of each parameter:

@tsk: `task_struct` of target process, only used to increase the page exception count of process; @mm: `mm_struct` of target process; @addr: the begging accessing address(relative to virtual address of target process); @buf: destination or source buffer; @len: data length, in bytes semaphore; @write: read / write flag

Fig.2 gives the function flowchart of "access_remote_vmspace":

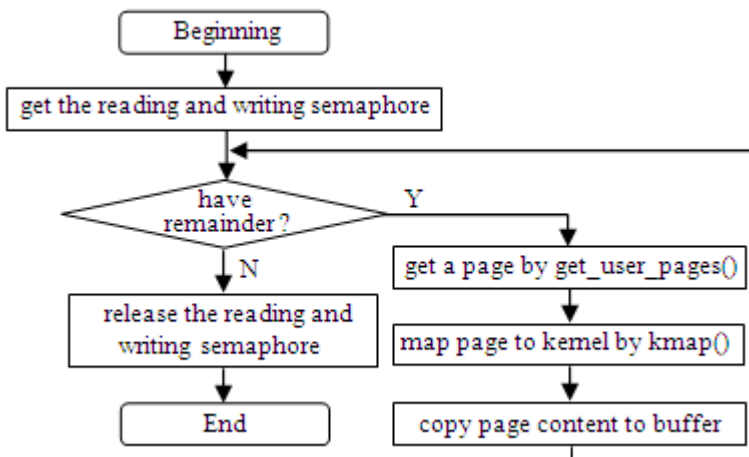


Fig. 2. Function flowchart of "access_remote_vmspace"

In order to verify the proposed method, we make the following experiment of reading code segment length of bash, cat, python and gnome-terminal (Ubuntu 11.10 system). Table 2 gives the results.

Table 2. Reading code segment length

Process name	Real code segment length(Byte)	Read code segment length(Byte)
bash	894,148	894,148
cat	35,372	35,372
python	2,238,912	2,238,912
gnome-terminal	276,168	276,168

As showed in fig. 2, the content of target process address space can be read integrity. Furthermore, considering the influence to system performance, the method can be realized in system kernel thread

The implementation of method 1 is as follows:

Calculate the length of code segment, data segment, and stack segment of object process;

Get buffer by calling kmalloc() function;

Copy code segment, data segment and stack segment from object process to buffer requested above;

Release the buffer.

3.2 Method 2: Switching Controlling Register CR3 Manually

To access the memory in the running process, MMU is responsible for conversion of virtual address to the physical address. MMU address translation's first step is to go to the CR3 register to get directory basic address of current process (physical address), followed by taking a page directory, page table entries, the last page and offset address within the page, in can be seen that the source of address translation is the CR3 register, whenever the kernel is ready to dispatch a new process, at the same time it will set the CR3 register. The procedure of switch this register complete by the switch_mm() function in the kernel, key statement are as follows:

```
asm volatile ("movl %0, %%cr3" : "r" __pa(next->pgd));
```

This is a assembly instruction embedded in C language, and the "next" is mm_struct of running process which is upcoming switch to enter the CPU. Obviously, the effect of this instruction is removed pgd field (basic address of the page directory table), and convert pgd content to physical address (the role of macro _pa() is the translation of the kernel virtual address into physical address), and finally transfer this physical address into register CR3. In this way, the address mapping of process switched from the source.

In fact, it is sure for the programmer to switch the CR3 register, for it is in kernel space, After CR3 register switched to the target process, you can access its address space data by using the virtual address of target process, because at this time the CPU memory access the data has been completely in accordance with the address mapping of target process of measure when it translate the address. Make a popular analogy is

to get someone's home key (CR3 register in the program), and then open the door of someone's home to get something (remove the data of the address space), natural unimpeded. Of course, it must be return to the key, after reading the data, it need to switch CR3 register back to physical address corresponding its process's page directory basic address.

However, there is a defect in this methods, namely, if the address space of target process has page fault exception, how to do? The program will collapse, and can not be read out the entire contents of the target process image.

Programmed in accordance with the method above, after being repeated tests to read the code segment of system login process, it can be read completely with the expected. In the subsequent tests, when trying to access a larger process code segment (such as the process gnome-terminal and so on.), it will occur page exception, and the measurement agent (MA) process is killed.

The implementation of method 2 is as follows:

Switch CR3 to page directory physical base address of object process;

Assess the address space of object process using linear address, and read code segment, data segment and stack segment of object process image;

Switch CR3 to page directory physical base address of initial process.

3.3 Method 3: Switching Controlling Register CR3 in Kernel Thread

Method 2 is an elegant and simple way, but there are indeed deeply flawed, after all, we can not expect all the memory page of the measure target process in the memory, in particular, when the process consumes large memory space, accessing to their address space, it is inevitable to occur page fault exception, Is there a better way to make up for deficiencies in method 2?

Linux memory management mechanism provides a pair of function similar method 2: `use_mm()` and `unused_mm()`.

These two functions are defined in the file `/mm/mmu_context.c`, according to the code comments, these two functions should only be used in the kernel thread.

The `task_struct` ,a process descriptor which contains two segment related to the process address space "mm" and "active_mm", to the normal user process, "mm" point to the user space of virtual address space, for the kernel thread , "mm" is NULL. "active_mm" is mainly used to optimization. the kernel thread is not associated with any particular user-level process, the kernel does not require exchange of the user-level part of virtual address, retaining the old settings is well, May some user-level process be antedate executed than the kernel threads, so the content of user-level is essentially random, kernel modules must be not modify its contents, so "mm" is set to NULL, at the same time, if it switched out is user process, kernel place the "mm" of original process into the "active_mm" of the new kernel process, although, the kernel thread, in principle, should not access user space of virtual address, but in fact it can be done, to access the user space of one process, just the "mm" field in the "task_struct" of kernel thread is set to "mm_struct" address of target process.

Note that in "use_mm", first, seting the "mm" of `task_struct` to "mm_struct" of target process, then calls "switch_mm()" to switch CR3 register of the global page

directory. As a result, the kernel can free to access user space data of target process (but not write), the page fault exception happened in the access procedure can be handled automatically, after the completion of reading data, then call “unused_mm()”, the “mm” field in the “task_struct” of kernel thread is come back to NULL. The implementation of method 3 is as follows:

- Create kernel thread;
- Call use_mm() function in the kernel thread;
- Read the data of virtual address space of object process in the kernel thread;
- Call unused_mm() function, reset the “mm” of thread as NULL.

3.4 Comparison of Three Methods

Specifically pointed out that, the CPU uses rotary mechanism of time slice when it deals with the switch of process. When the measurement agent (MA) access resources of one process, the registers in the CPU are stored in the current information of MA, and can be not access the register information of target process.

Table 3. Comparison of three methods of accessing specified process address space

Method	difficulty	complexity	security	Dependency in platform	Speed
Method1	high	complex	high	no	faster
Method2	high	simple	low	yes	faster
Method3	general	simple	high	no	faster

4 Realization Comparison of Three Methods of Accessing Process Space

In order to have a deep insight into three methods, we make the following performance measurement experiments of bash, cat, python and gnome-terminal respectively, including CPU occupation rate and time overhead. The experiment values are average values of 20 times. The experiments environment is that, Intel(R) Core(TM)2 Duo CPU 2.00GHZ, 1GB DRAM, and Ubuntu 11.10. Table 4 gives the comparison

Table 4. Comparison of three methods of CPU occupation ration and time overhead

Method	App	CPU occupation ratio (%)		CPU occupation ratio value added (%)	Time overhead (μs)
		closing Measurement mechanism	Opening Measurement mechanism		
Method1	bash	0.32%	0.70%	0.38%	4,999,556
	cat	0.18%	1.00%	0.82%	4,999,567
	python	0.15%	1.71%	1.56%	5,002,582
	gnome-terminal	0.65%	1.20%	0.55%	4,997,600
Method3	bash	0.63%	1.40%	0.77%	4,999,522
	cat	0.20%	1.10%	0.90%	4,996,710
	python	0.20%	1.30%	1.10%	5,000,228
	gnome-terminal	0.60%	1.31%	0.71%	4,997,295

Note: Due to the defect in needing all pages of object process to be in memory, method 2 is not included here. In theory, the performance of method 2 is the same to method 3.

5 Conclusions

In this paper, In order to meet the need to trustworthiness measurement, the basic measurement system architecture is proposed, and three methods of LKM (Loadable Kernel Module) accessing specified process address space are presented. Among these, method 1(address remapping) technology difficulty is higher with high security, method 2 adopt the way switching process page directory address register CR3 with some defect in page fault abnormal and dependency in X86 platform, and method 3 is a recommendation choice with no dependency in hardware platform and high efficiency.

References

1. Avizienis, A., Laprie, J.-C., Randell, B., Landwehr, C.: Basic concepts and taxonomy of dependable and secure computing. *IEEE Transactions on Dependable and Secure Computing* 1(1), 11–33 (2004)
2. Shen, C.X., Zhang, H.G., Feng, D.G., et al.: Research and development of trusted computing. *China Science* 40(2), 139–166 (2010) (in Chinese)
3. Liu, Z.W., Feng, D.G.: TPM-Based Dynamic Integrity Measurement Architecture. *Journal of Electronics & Information Technology* 32(4), 875–879 (2010) (in Chinese)
4. Li, X.Y., Gui, X.L., Mao, Q., et al.: Adaptive Dynamic Measurement and Prediction Model Based on Behavior Monitoring. *Chinese Journal of Computer* 32(4), 664–674 (2009) (in Chinese)
5. Cai, S.B., Zou, Y.Z., Shao, L.S., et al.: Framework Supporting Software Assets Evaluation on Trustworthiness. *Journal of Software* 21(2), 359–372 (2010) (in Chinese)
6. Sailer, R., Zhang, X., Jaeger, T., et al.: Design and implementation of a TCG-based integrity measurement architecture. In: 3th Conference on USENIX Security Symposium, pp. 223–238. USENIX Association, Berkeley (2004)
7. Trusted Computing Group. TCG Specification Architecture Overview [DB/OL] (March 01, 2005), <https://www.trustedcomputinggroup.org/>
8. Lin, H., Lee, G.: Micro-Architecture Support for Integrity Measurement on Dynamic Instruction Trace. *Journal of Information Security* 1, 1–10 (2010)
9. Maruyama, H., Nakamura, T., Munetoh, S., et al.: Linux with TCGA Integrity Measurement. IBM Japan, Ltd. (January 28, 2003)
10. Jaeger, T., Sailer, T., Shankar, U.: PRIMA: Policy-reduced integrity measurement architecture. In: The 11th ACM Symposium on Access Control Models and Technologies, pp. 19–28. ACM, New York (2006)
11. Azab, A.M., Ning, P., Sezer, E.C., et al.: HIMA: A Hypervisor-Based Integrity Measurement Agent. In: The 2009 Annual Computer Security Application Conference, pp. 461–470. IEEE, Honolulu (2010)

Human Resource Management System Based on Factory Method Design Pattern

Xing Xu¹, Hao Hu², Na Hu¹, Lin Xiao¹, and Weiqin Ying³

¹ College of Information & Engineering, Jingdezhen Ceramic Institute, Jiangxi, China

² Library of Huaihua University, Huaihua University, Huaihua 418000, China

³ School of Software, South China University of Technology, Guangzhou 510006, China
whuxx84@yahoo.com.cn

Abstract. Analyzing the human resource management system (HRMS), many subtle differences in the user interface design (UI) result in a large number of repetitive works for developers. In order to avoid such a situation, design pattern, as an important software reusable technology, is broadly applied to many kinds of information management platform for effectively economizing development costs. This paper introduces a kind of factory method design pattern as an essential aid of system design. The factory method is imported in HRMS for analyzing and designing the UI. It shows better reusability, scalability, maintainability and provides a strong support for meeting the growing needs of different business in HRMS, so that people can more simply and conveniently reuse successful design and architecture and provides further details of its application in UI of HRMS.

Keywords: Design Pattern, Factory Method, Human Resource Management System, User Interface Design.

1 Introduction

In accordance with increasing maturity of market economy system and tightening market competition at present, human resource management system (HRMS) is playing a more and more important role in market competition [1]. However, the more affairs HRM department handles, the more tough problems in management it will encounter. Such problems render less time devoted to deal with a company's strategic development, which results in much more time for trivial daily affairs. Due to the nature of different trades, HRMS differs from each other in terms of their major concerns [2]. For example, technique-oriented companies usually place much more emphasis on staff's technological level, while manufacturing-oriented ones attach more importance to workers' skills and experiences [3].

Consequently, it is rather difficult for HRMS to solve all these problems with traditional technology in design. Thus, it will cost a large proportion of time to design HRMS, in particular interface design. And when the concept of role is concerned in HRMS, the difference of role will lead to the difference of interface and accordingly make systems interface design more complicated.

2 Difficulties in HRMS Interface Design

According to authoritative investigation, in the process of HRMS development, time consumption for different procedures is as follows.

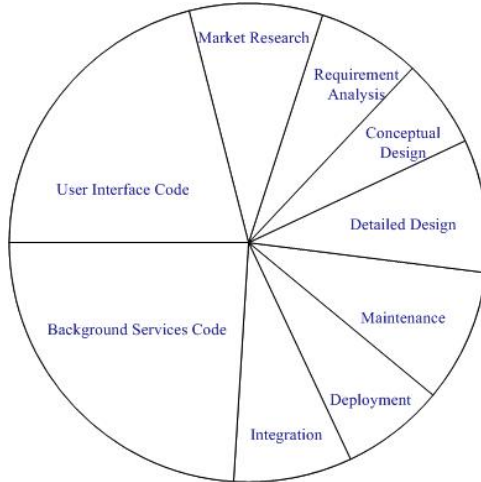


Fig. 1. The time occupation proportion of HRMS development

The fig. 1 shows that interface design and encoding account for the most time. Particularly as for a multi-role and multi-access system, when very complex business is involved, the amount of time spent on these two parts will overtake the amount on back-stage service, which can be caused by the following: (1) the system processing personnel have high expectation towards the operability and user-friendliness; (2) it is much necessary to take individual factors into account in interface design and development in order to simulate human's real mind, which will inevitably increase the overall difficulty of system design.

Traditional HRMIS framework is basically like what Fig. 2 shows [4]. It is known through analysis that system users possess different roles, and their interface is different according to their access. For instance, generally speaking, users in HRMS share a lot of common features, including baseMessages such as name, gender, age, department, telephone number, email; externMessages such as ID number, residence address, occupation, company registration number; loginMessages like login name and code, workExperienceMessages and studyExperienceMessages. The display interface varies according to different roles as follows:

Administrator: baseMessages, externMessages, loginMessages;

Employee: workExperienceMessages, studyExperienceMessages, baseMessages, externMessages, loginMessages;

Company leader: baseMessages, externMessages;

Project manager: baseMessages, externMessages, workExperienceMessages, studyExperienceMessages.

Different roles have different access, so does the content of system display interface. Therefore, difficulty in interface design will arise. Programmers have to make design for every interface and even though such potential interfaces have a lot in common, they have no other choice but to repeat their work. Such system, though with low coupling degree, is weak reusability and scalability.

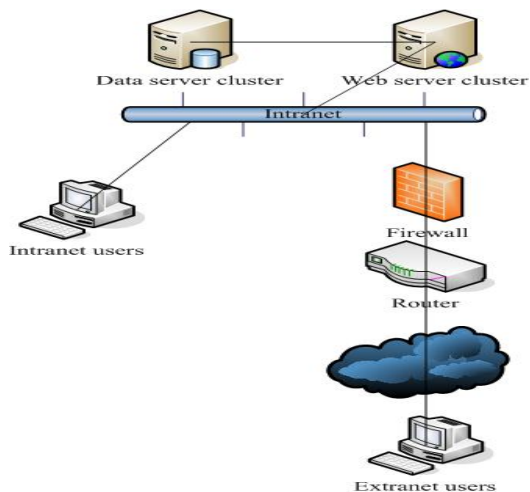


Fig. 2. Network architecture of HRMS

3 Factory Method Design Pattern

In 1987 Kent Beck and Ward Cunningham introduced the thought of design pattern in architecture to the field of program design [5]. So far the design pattern has been acknowledged as a powerful framework tool in the field of object-oriented software design [6], [7], [8], [9]. Design patterns and traditional development methods are very different in terms of emphasis on problem solution. For traditional development method of general system of, the key issue is how to use computer system as a platform to build up a stable and high-efficiency execution system, while design patterns need to address the problem of how to meet users' changing needs and how to solve software expansion in the software development process. In short, the core technology is how to achieve software reuse in a better way [10]. Software reuse technology is an important technology for the development of software, which makes programmers avoid repeated work in later development process, and just use the existing work, based on the existing work, to change and extend in order to cater for new demand.

Design patterns are mainly in the form of template for developers' use. All kinds of template solutions are put forward specific to different problems. For example, facade pattern is mainly used for solving the problems at the database access layer, state pattern and command pattern are essentially used to solve the process problems at the business layer. Factory design pattern mentioned in this article is basically used in the interface design at user presentation layer.

Usually there are four elements in every design pattern [4, 11]:

pattern name: the name of a design pattern;

problem: to clarify what kind of conditions to use a certain pattern;

solution: to specify the method to realize design pattern, the class, object and relations between the two, as a template, and not just for a specific system;

consequence: the effect of using a design pattern, including its merits and shortage.

In the real world there are 23 standardized design patterns. Each pattern can attack a problem, or be combined with other patterns to address other problems. New design patterns can also be created. Since design patterns are different in particle size and abstract level, standardized design patterns can be categorized into three groups in light of different application principles:

Creational Patterns: to create object like Factory Method, Builder, etc;

Structural Patterns: combine class and object to form correspondent structures like Adapter, Proxy;

Behavioral Patterns: describe how class and object interact, how a task is dealt with by a different object, such as Command, Visitor.

Factory method of design patterns, also called factory pattern, or virtual constructor pattern or polymorphism factory pattern, belong to the creational design pattern. Factory pattern mainly provides the creational pattern for class, the father class is responsible for the definition of public interface of creational object, and a subclass is responsible for generating concrete objects in order to delay the instantiation operation to subclasses, i.e. the subclasses will decide which instantiation creates which class. A typical factory pattern class structure is demonstrated as follows:

The relationship between class object is presented below:

Product: the role of product, defining the interface of product;

Realproduct: the real-life product, realizing the interface class;

Creator: the role of factory, declaring factory method and return a product;

RealCreator: the real factory, realizing factory method, used by customers and return a product instantiation.

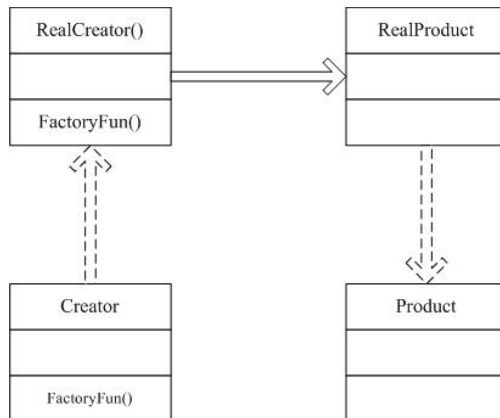


Fig. 3. Class Diagram of Factory Pattern

It is known from the analysis of this model that factory pattern defines the main Product, and defines many virtual interfaces within it, and product is the realization class of this class. Creator is the real Product factory class, inherited from the Product class, also defines many virtual interfaces within it. RealCreator is the realization class for Creator, which realize the interface in Creator. Users will invoke Creator to obtain the instantiation of designated product. So, users don't have to care about the contents of the Product, and only need to pay attention to how to invoke RealCreator to construct real products. In fact, the product instantiation constructed by RealCreator possesses all attributes of Product. Accordingly, many common attributes can be extracted and be managed by Product, which achieves the reuse of software [12].

4 Factory Method Applied to HRMS

For any set of human resource management system, developers often think that user interface design and development is a very complicated task. According to the foregoing discussion, once the initial design is not good enough, it is likely to lead to a large number of labor resources waste in later system updating and in meeting users' needs. Therefore, we introduce factory method design pattern to system design, so as to enhance the system reusability and scalability, thus effectively reducing the waste of resources.

First, we can create base class which contains no information in the system called Product. The correspondent RealProduct is the interface realization class for the base class. In the process of interface design and development, we can create the base class in the system, and then build the correspondent factory roles of attributes, thus correspondent interface elements can be used. Later, when some properties need to change, we don't need to change the properties of the specific operation part. We can just modify the Product and the part which uses the Product will automatically change accordingly, which can effectively reduce repetitive labor in modifying interface.

The information shown by HRMS contains the following attributes: name, sex, age, telephone, address, position. And the interface not only has the function of looking through all basic attributes of ordinary employees, but also expands such a group of interface elements as work time, individual salary standard, employee performance in a certain month. Therefore, when the factory design pattern is used in HRMS interface design, the basic Product of the system will be defined as Resource class, that is employees' basic interface class, the common staff attributes of the system will be defined as basic product class—ResourceEmployee, inherited from the Resource class, including such specific attribute as name, sex, age, phone, address, position. It is in constructing the interface display method in ResourceEmployeeImpl to realize the layout and operation of the interface. Then, we can construct ResouceFinance on the basis of various and extended attributes of the financial personnel, which adds time, payStandard, performance to it. It will not only inherit from the ResourceEmployee class, but its interface layout will also descend from ResourceEmployee class. Thus, we realize the correspondent specific interface layout in ResouceFinance class. Class attribute and layout diagram are shown in Fig. 4.

In actual operation, the design personnel only need to focus on the actual operation interface class of ordinary staff, the ResourceEmployee class. In the development of staff performance module, they only need to care about the display, design and processing of the three parameters of performance. In later development period, if the attribute of common employee alters, it is just necessary to modify the ResourceEmployee class which displays the information of ordinary staff based on early development work, and the performance module will automatically change accordingly. There is no need for developers to make further revision.

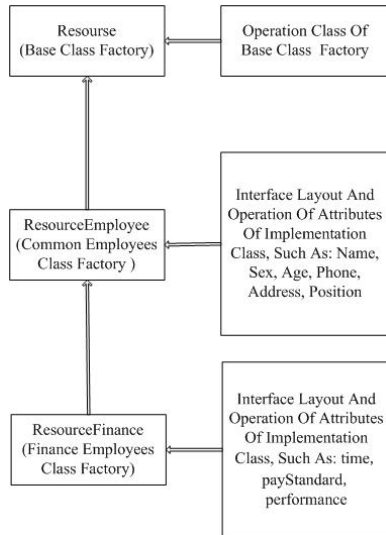


Fig. 4. The personnel interface design of two roles

When a part of the system requires the display of basic attributes of staff, we just need to create an instantiation in ResourceEmployee class in a specific place on the interface, where the system will present the interface of basic personnel attributes. Developers can invoke the operation method defined by the ResourceEmployee class to achieve the input and output in the interface, such as defining resourceLogin, which means to show the basic information of the first logging-in person, and then developers can invoke the resourceLogin.getName() method to obtain user’s name information on the interface, and invoke resourceLogin.GetSex() method to get user’s gender information on interface. When the interface switches to financial personnel management page, the page because inherited from the ResourceEmployee class, can directly use its operating methods to gain the basic information of the employees on the interface, and meanwhile, it also needs to define the invoking function getTime() for the working time of employees in a week and the invoking function setTime() to for the actual working time of employees in a week. After constructing the interface based on factory design pattern, if the system is required to display the actual performance salary information of a staff in a month based on the performance module, developers only need to add realPay attribute for ResourceFinance class factory, and design the operation method of this attribute realPay(). Any interface using this class will display the attributes actual performance salary of employees. There is no need to modify other modules.

Summarizing the above analysis, we can see that different interfaces of HRMS share a lot of common operating elements. At the same time, many seemingly different interfaces actually have great similarities. Therefore, we can extract the same element of these interfaces. And for the different elements, we can use subclass inheritance method and invoke different interfaces through different configuration parameters in order to avoid repetitive work. In later development and maintenance, once the same part of interfaces changes, we can modify the operation elements of father class in order to modify the correspondent operation interface of all subclasses inherited from this class, so as to reduce repetition, and improve system reusability. The internal structure and operation flow chart based on the factory design pattern is shown in Fig. 5 and Fig. 6.

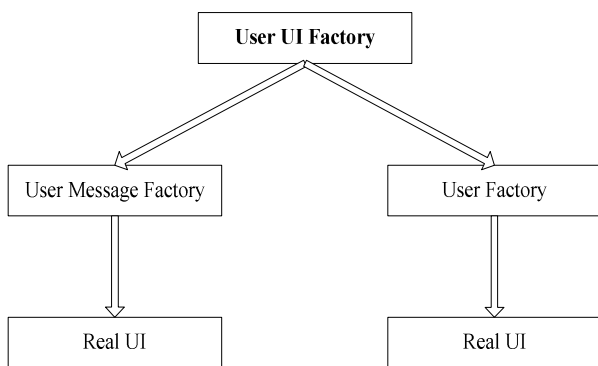


Fig. 5. Internal structure of factory design pattern

It is observed, from the above two fig.s, that developers, in the process of development, don't have to be concerned about the design and development of some public interface. If there is a need for any part of the system to be connected with the public part, developers can invoke correspondent user interface to generate factory and create the needed interface. Designers don't have to worry about public interface change. They only need to care about their own part of a specific attribute.

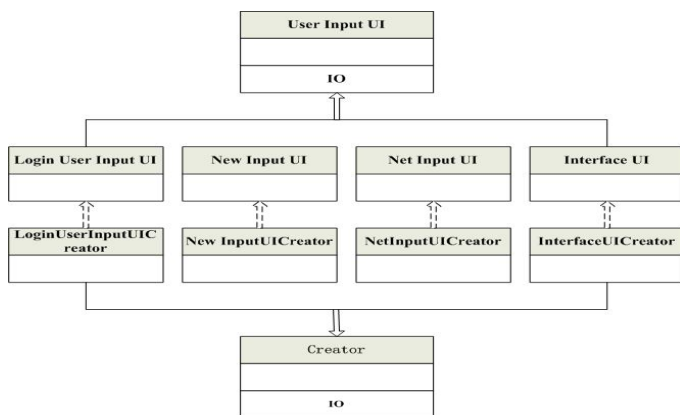


Fig. 6. Operation process of factory pattern

Acknowledgments. This paper was supported by the National Science and Technology Support Plan (2012BAH25F02), the Project of Jingdezhen Science and Technology Bureau (2011-1-47), the National Natural Science Foundation of Jiangxi Province (2009GZS0065), the Youth Science Foundation of Jiangxi Provincial Department of Education (GJJ12514), the Natural Science Foundation of Guangdong Province (S2011040002472), the Specialized Research Fund for the Doctoral Program of Higher Education (20110172120035), the Fundamental Research Funds for the Central Universities (2011ZM0107).

References

1. Cai, Z.: The human resource management system. *Journal of Qiqihar University (Natural Science Edition)* 24(1), 51–54 (2008)
2. He, C., He, K.: A Role-Based Approach to Design Pattern Modeling and Implementation. *Journal of Software* 17(4), 658–669 (2006)
3. Dessler, G.: *Human Resource Management*, 10th edn. Prentice Hall, New Jersey (2004)
4. Wang, K.: Design and Implementation of Human Resource Management Information System Based on Success Model. *Sci.-Tech. Information Development & Economy* 18(14), 152–155 (2008)
5. Chung, C.: *Pro Objective-C Design Patterns for iOS*. Apress, New York (2011)
6. Wu, D., Lu, L., Yang, F., Zhang, C.: Discussing of factory pattern in embroidery CAD. In: *IEEE International Conference on Software Engineering and Service Sciences*, pp. 234–236. IEEE Press, New York (2010)
7. Ni, Y., Miao, M., Li, M.: Design and implementation of the business layer of the Tourism Business Information Collection and Distribution System based on the proxy and the factory patterns. In: *International Conference on Computer Research and Development*, pp. 199–202. IEEE Press, New York (2011)
8. Ko, J., Song, Y.: Test Driven Development of Model Transformation with Reusable Patterns. In: Park, J.J., Chao, H.C., Obaidat, M.S., Kim, J. (eds.) *Computer Science and Convergence*. LNEE, vol. 114, pp. 597–605. Springer, Heidelberg (2012)
9. Mannava, V., Ramesh, T.: A Novel Adaptive Monitoring Compliance Design Pattern for Autonomic Computing Systems. In: Abraham, A., Mauri, J.L., Buford, J.F., Suzuki, J., Thampi, S.M. (eds.) *ACC 2011*. CCIS, vol. 190, pp. 250–259. Springer, Heidelberg (2011)
10. Zheng, G.: Application of Three Design Patterns to Personnel Management System. *Journal of Lanyungang Technical College* 22(1), 6–7 (2009)
11. Zhou, X., Xu, B.: New Research Framework for Automatic Recovery of Design Pattern. *Computer Science* 36(5), 124–128 (2009)
12. Liu, Y., Liu, L.: Research on application of design patterns in information system. *Information Technology* 10, 129–131 (2008)

Ant Colony Optimization with Multi-Agent Evolution for Detecting Functional Modules in Protein-Protein Interaction Networks

Junzhong Ji¹, Zhijun Liu¹, Aidong Zhang², Lang Jiao¹, and Chunnian Liu¹

¹ College of Computer Science and Technology, Beijing University of Technology, Beijing Municipal Key Laboratory of Multimedia and Intelligent Software Technology, Beijing, 100124, China
jjz01@bjut.edu.cn

² Department of Computer Science and Engineering, University at Buffalo, The State University of New York, Buffalo, 14260, USA
azhang@cse.buffalo.edu

Abstract. Functional module identification in a Protein-Protein Interaction (PPI) network is one of the most important and challenging tasks in computational biology. For detecting functional modules, it is difficult to solve the problem directly and always results in a low accuracy and a large discard rate. In this paper, we present a novel algorithm of ant colony optimization with multi-agent evolution for detecting functional modules. The proposed ACO-MAE algorithm enhances the performance of ant colony optimization (ACO) by incorporating multi-agent evolution (MAE). First, the ant colony optimization for solving Traveling Salesman Problems (TSP) is conducted to construct primary clustering results. Then, the multi-agent evolutionary process is performed to move out of local optima. From simulation results, it is shown that the proposed ACO-MAE algorithm has superior performance when compared to other existing algorithms.

Keywords: Protein-Protein Interaction Network, Functional Module Detection, Ant Colony Optimization, Multi-agent Evolution.

1 Introduction

Analysis of the underlying relationships in protein data is a matter of great significance to study the mechanism of human disease, and discover new therapeutic interventions [1]. In particular, protein-protein interactions provide us with a good opportunity to systematically analyze the structure of a large living system and also allow us to use them to understand essential principles. Cellular functions and biochemical events are coordinately carried out by groups of proteins interacting each other in functional modules, and the modular structure of complex networks is critical to function [2]. Therefore, identifying such functional modules in PPI networks is very important for understanding the structure and function of these biological networks. However, experimental approaches for detecting functional modules remained relatively immature, hence

the research of computational approaches for detecting functional modules has become an essential and challenging problem in computational biology [3–11].

This paper is organized as follows. we describe ant colony optimization and multi-agent evolution in Section 2. In Section 3, we present the ACO-MAE algorithm for detecting functional modules. Next, Section 4 shows the simulation results, which are also compared to other existing algorithms to demonstrate the superiority of ACO-MAE algorithm. Conclusions are given in Section 5.

2 Ant Colony Optimization and Multi-agent Evolution

ACO and MAE are population-based search algorithms by maintaining a population of individuals as candidate solutions in solving practice problems. Both algorithms can be embedded with other approaches to speed up the search performance. The basics of ACO and MAE are described below.

2.1 Ant Colony Optimization

Ant colony optimization (ACO) is a new meta-heuristic search algorithm inspired by the observation of real ants looking for food. Ethnologists observed that ants can find the shortest path from their nest to the food source by exploring and exploiting pheromone information, which has been deposited on the path where they traversed. They then can choose routes based on the amount of pheromone. Namely, ants communicate information about food source via pheromone. The larger amount of pheromone is deposited on a route, the greater is the probability of selecting the route by ants. Thus, when one ant finds a good short path from the nest to a food source, other ants are more likely to follow this path, and such a self-strengthening behavior eventually leads all the ants to follow the shortest path. The idea of the ACO is to mimic this behavior with artificial ants walking around the graph representing the problem to solve. ACO was first proposed by Dorigo *et al.* [12]. This algorithm often gives satisfactory results for various optimization problems in a wide range of domains [13], such as data mining, machine learning, bioinformatics and multiple objective optimization problems.

2.2 Multi-agent Evolution

Agent-based computation has also been widely used as a new cooperative search algorithm in the field of computer science [14]. A multi-agent system can be viewed as an evolutionary system where each agent acts on a basis of population ecosystem under the idea of evolution.

In a multi-agent system, an agent, a , is a virtual entity that is able to live and act in the environment, and has some reactive behaviors to the environment. Generally speaking, each agent essentially has two properties: problem knowledge and certain objectives, and is able to communicate, react and cooperate with the environment and other agents. The basic idea of a multi-agent evolution system is: each individual in the traditional evolution algorithm is considered as

an agent who makes use of the evolution mechanism. The agent first exchanges information with the environment and other agents, then performs the cooperation behaviors among agents, eventually achieves the common adaptation among agents and between the agent and the environment. The multi-agent evolution model usually contains the following rules [15]: 1) each agent has an initial energy; 2) each agent is only able to sense and act in its local environment (called neighborhood) which is limited; 3) there are competitions among agents, and the agents with lower energy will die, the behavior is called as survival of the fittest; 4) due to agents died empty some positions, so some alive agents can produce offsprings to substitute, the behavior is called as law of the jungle; and 5) each agent has mating ability, it finds a good mating partner to mate in its neighborhood, and passes its good genes into its offsprings. During the process of interacting with the environment and companions, each agent increases its energy as much as possible, so that the multi-agent evolution can achieve the ultimate purpose of solving the global optimization problem.

3 The ACO-MAE Algorithm

The procedure of the proposed ACO-MAE algorithm is to apply initialization, transformation of problem, ant colony optimization, and multi-agent evolutionary process for detecting functional modules. The last two processes are iterated until a pre-specified termination criterion is satisfied. The proposed algorithm is shown in Fig.1.

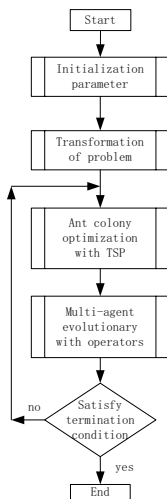


Fig. 1. The flowchart of the proposed algorithm

After clustering proteins in PPI networks is transformed into searching the optimal tour, the ant colony optimization is employed to generate m solutions

in each iteration, and then K solutions with the shortest length are selected and further evolved by using multi-agent evolutionary scheme. Because the methods of problem transformation and ant colony optimization have been explicitly described in [11], thus this paper will focus on the multi-agent evolutionary process in detecting functional modules of PPI networks.

3.1 Agent and Its Energy

In each iteration of the ACO algorithm, a set of solution tours which visit each protein of the PPI network exactly once can be obtained. From these solution tours, we select the Top- k solutions with the shortest lengths as agents to further evolve. The encoding of an agent is denoted as $a = (a_1, a_2, \dots, a_N)$, where N is the number of proteins in the PPI network, a_i denotes the protein connected to i^{th} protein in the corresponding clustering results. Fig.2 shows a solution of an ant, corresponding to the clustering results and the encoding of its agent, where λ denotes a cutoff value. In the paper, let λ be $\gamma \cdot \bar{d}$, where γ is a real parameter and \bar{d} is an average distance of all paths between proteins. By removing these paths whose length are larger than λ in a solution tour, the clusters which include shorter connections between proteins are formed. Such clusters can be viewed as an agent which encodes the connection relationships of the clusters.

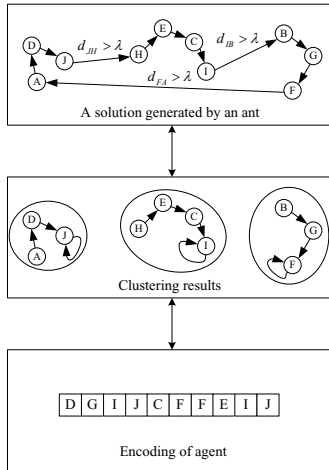


Fig. 2. A solution of an ant, corresponding clustering results and the encoding of its agent

As can be seen in Fig.2, an agent represents a candidate result for detecting functional modules. The value of its energy can be formulated as:

$$Energy(a) = - \sum_{i=1}^n L(C_i), \tag{1}$$

where n is the number of clusters for the agent a , and $L(C_i)$ is the connection length among proteins in the C_i cluster. The purpose of a is to increase its energy as much as possible by sensing and performing some reactive behaviors during its evolution process.

3.2 Multi-agent Evolutionary Environment

In order to realize the local perceptivity of agents, the environment is organized as a lattice structure. All K agents live in such a lattice environment. The size of lattices is $L_k \times L_k$, where L_k is an integer and $L_k = \sqrt{K}$. Each agent is randomly fixed on a lattice-point and it can only interact with its neighbors. The agent lattice can be shown as the one in Fig.3. Each agent, who represents a clustering result, occupies a circle in the evolutionary environment, the data in a circle represents its position in the lattice structure, and two agents can interact with each other if and only if there is a line connecting them.

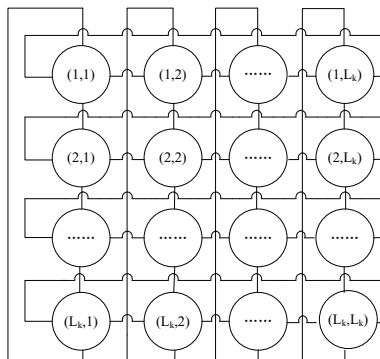


Fig. 3. Evolutionary environment of agent lattices

Suppose that the agent located at (i, j) is $a_{i,j}$, $i, j = 1, 2, \dots, L_k$, then the neighbors of $a_{i,j}$, $neighbor(i, j)$ are defined as follows:

$$neighbor(i, j) = \{a_{i',j}, a_{i,j'}, a_{i'',j}, a_{i,j''}\}, \tag{2}$$

where $i' = \begin{cases} i - 1 & i \neq 1 \\ L_k & i = 1 \end{cases}$,

$$j' = \begin{cases} j - 1 & j \neq 1 \\ L_k & j = 1 \end{cases}, i'' = \begin{cases} i + 1 & i \neq L_k \\ 1 & i = L_k \end{cases}, \text{ and } j'' = \begin{cases} j + 1 & j \neq L_k \\ 1 & j = L_k \end{cases}.$$

3.3 Evolutionary Operators

In the above agent evolutionary environment, agents will compete and cooperate with others so that they can gain more resources. Since each agent can only sense

its local environment, its behaviors of competition and cooperation can only take place between the agent and its neighbors. In such a manner, the information is diffused to the whole agent lattice. Especially, three operators are designed to achieve this purpose for detecting functional modules in PPI networks.

Neighborhood Competition Operator. Suppose that the operator is performed on the agent located at (i, j) , $a_{i,j} = (a_1, a_2, \dots, a_N)$, and $m_{i,j} = (m_1, m_2, \dots, m_N)$ is the agent with maximum energy among the neighbors of $a_{i,j}$, namely, $m_{i,j} \in neighbor(i, j)$ and $\forall a' \in neighbor(i, j)$, then $Energy(a') \leq Energy(m_{i,j})$. If $Energy(a_{i,j}) \geq Energy(m_{i,j})$, $a_{i,j}$ is a winner, so it can still live in the agent lattice; otherwise it is a loser, it will die, and its lattice-point will be occupied by $m_{i,j}$. $m_{i,j}$ has two strategies to occupy the lattice-point, and it selects them with a probability p_o . Let $R(0, 1)$ be a uniform random number generator, If $R(0, 1) < p_o$, occupying strategy 1 is selected; otherwise occupying strategy 2 is selected. In both occupying strategies, $m_{i,j}$ first generates a clone agent $c_{i,j}$, and then $c_{i,j} = (c_1, c_2, \dots, c_N)$ is put on the lattice-point.

Let $d(i, a_i) = al_i$, and $d(i, m_i) = ml_i$, $i = 1, 2, \dots, N$, namely, the connection lengths of $a_{i,j}$ are al_1, al_2, \dots, al_N , and the connection lengths of $m_{i,j}$ are ml_1, ml_2, \dots, ml_N respectively.

Strategy 1. For the connection with the largest length in $m_{i,j}$, $ml_j = \max(ml_1, ml_2, \dots, ml_N)$, c_j is replaced with a_j .

Strategy 2. Each al_i of $a_{i,j}$ is respectively compared with ml_i of $m_{i,j}$. If $al_i < ml_i$, then $c_i = a_i$.

Even $a_{i,j}$ is a loser, it perhaps still has some useful information. Thus, two strategies in this operator play similar roles. The occupying strategy 1 only replaces the worst connection of $m_{i,j}$ with information of $a_{i,j}$, and the strategy 2 is in favor of reserving all advantaged information of a loser.

Neighborhood Crossover Operator. Suppose that two agents are $a = (a_1, a_2, \dots, a_N)$ and $b = (b_1, b_2, \dots, b_N)$ which produce a new agent $e = (e_1, e_2, \dots, e_N)$ by making use of their connection information, and a Mask vector, $\mathbf{M} = (M_1, M_2, \dots, M_N)$, is randomly generated, where each $M_i \in \mathbf{M}$ is either 0 or 1. If $M_i = 1$, then the new agent inherit the connections from a , otherwise it inherit the connections from b . That is,

$$e_i = \begin{cases} a_i & M_i = 1 \\ b_i & M_i = 0 \end{cases} \quad (3)$$

The neighborhood crossover operator has the function of random searching, which is performed on $a_{i,j}$ and its neighbors to achieve the purpose of cooperation in light of a probability p_c .

Mutation Operator. In addition to the behaviors of competition and cooperation, each agent can also increase its energy by using a mutation operator. Based on a mutation probability p_m , an element a_i of an agent $a = (a_1, a_2, \dots, a_N)$ is randomly selected, and then it is replaced with another protein to which i^{th}

protein might connect in the corresponding clustering results. In nature, the mutation operator realizes a local search which only performs a small perturbation on some variables of a .

4 Empirical Study

In this section, we use the protein-protein interaction data downloaded from DIP (Database of Interaction Protein: <http://dip.doe-mbi.ucla.edu/>). We assess the performance of our algorithm, and compare its test results to other classic algorithms on the same benchmark data sets. By large numbers of experiments, the final experimental parameters are confirmed as follows: $N = 100$, $m = 200$, $K = 81$, $Q = 100$, $\rho = 0.5$, $\alpha = 1$, $\beta = 3$, $p_o = 0.2$, $p_c = 0.8$, $p_m = 0.1$.

To evaluate the detected protein modules, the set of real functional modules from [16] is selected as the benchmark. This benchmark set, which consists of 428 protein functional modules, is constructed from three main sources: the MIPS [17], Aloy et al. [18] and SGD database [19] based on Gene Ontology (GO) annotations. In our following experiment, we use core protein interaction data of *saccharomyces cerevisiae* from DIP, which contains 2526 distinct proteins and 5949 highly reliable interactions.

Table 1. Experimental results and comparisons using our algorithm with different γ

γ	Num.	Ave.	Dis.	Pre.	Rec.	F.	Sen.	Pos.	Acc.	$-\log(P.)$
0.70	257	7.84	0.2029	0.3774	0.3691	0.3732	0.5645	0.2951	0.4082	12.74
0.75	213	9.87	0.1685	0.3568	0.2780	0.3125	0.5774	0.2743	0.3979	13.47
0.80	198	10.91	0.1452	0.3737	0.2453	0.2962	0.6038	0.2360	0.3775	12.74
0.85	175	12.60	0.1278	0.2971	0.1799	0.2241	0.6006	0.2355	0.3761	13.69
0.90	162	13.78	0.1167	0.2901	0.1565	0.2034	0.6091	0.2208	0.3667	13.74
0.95	143	15.52	0.1218	0.2657	0.1332	0.1774	0.6072	0.1971	0.3459	12.68
1.0	147	15.21	0.1155	0.3129	0.1495	0.2024	0.6095	0.1915	0.3416	12.12

Though there are many parameters in our algorithm, most of parameters are not sensitive to detection results except for the cut-off parameter γ [11]. To select proper value of γ , we perform many comparisons in light of ten evaluation metrics including the number of modules (Num.), the average size of modules (Ave.), the node discard rate (Dis.), precision (Pre.), recall (Rec.) (where $\omega = 0.2$), F-measure (F.), sensitivity (Sen.), positive predictive value (Pos.), accuracy (Acc.) and P-values (P.), whose meanings and computing formulaes can be found in [20]. Table 1 shows some experimental results. In general, p -value is known as a metric of functional homogeneity. High $-\log(P.)$ (i.e., low p -value) indicates that the module closely corresponds to the function because the network has a lower probability to produce the module by chance, thus we select the value of γ with high $-\log(P.)$. Considering other performance metrics, we set $\gamma = 0.75$ in the following experiments.

To assess the accuracy of each method, we compared our method with several previous state-of-the-art methods. We also computed the statistical p -value using the predicted modules by each algorithm and the reference modules. For each predicted module, we assigned a reference module by finding the lowest p -value, and calculated the average negative $\log(p\text{-value})$ for all predicted modules as the average p-score.

Table 2. Performance comparisons with some clustering algorithms on yeast datasets

methods	Num. of modules	Avg. size of modules	Nodes discard rate (%)	Avg. p-score - $\log(p\text{-value})$
ACO-MAE	213	9.87	16.9	13.47
NACO-FMD	161	10.07	36.1	14.91
Markov Clustering	163	9.79	36.7	8.18
Minimum Cut	114	13.46	35.0	8.36
Neighbor Merging	64	7.91	79.9	9.16
Interconnection Cut	180	10.26	21	8.19

In terms of accuracy (average p-score), the results in Table 2 demonstrate that our method outperforms the other previous methods except for NACO-FMD [12]. Though our method is inferior to NACO-FMD on the functional homogeneity, it can generate much more modules, and it has the least node discard rate among all methods.

5 Conclusions

ACO algorithms have been applied for detecting functional modules of a PPI network in recent years. In this paper, we presented a novel algorithm of ant colony optimization with multi-agent evolution for detecting functional module. In the proposed ACO-MAE algorithm, multi-agent evolution is embedded into ACO algorithm to perform the local search, which can avoid local optima and ameliorate the search performance. From simulation results, it is shown that the proposed algorithm can improve search performance and outperforms other existing algorithms. Our future work includes further studying effectiveness of the ACO method for detecting functional modules, and focusing on the problem of overlapping modules.

Acknowledgments. This work is supported by the Beijing Natural Science Foundation (4102010), and Beijing skeleton teacher program (007000543111511).

References

1. Zhang, A.D.: Protein Interaction Networks: Computational Analysis. Cambridge University Press (2009)
2. Guimera, R., Nunes Amaral, L.: Functional cartography of complex metabolic networks. *Nature* 433, 895–900 (2005)

3. Spirin, V., Mirny, L.A.: Protein complexes and functional modules in molecular networks. *Proceedings of the National Academy of Sciences* 100(21), 12123–12128 (2003)
4. Pei, P., Zhang, A.D.: A "seed-refine" algorithm for detecting protein complexes from protein interaction data. *IEEE Transactions on Nanobioscience* 6(1), 43–50 (2007)
5. Arnau, V., Mars, S., Marin, I.: Iterative cluster analysis of protein interaction data. *Bioinformatics* 21(3), 364–378 (2005)
6. Frey, B.J., Dueck, D.: Clustering by passing messages between data points. *Science* 15(5814), 972–976 (2007)
7. Hwang, W., Cho, Y.R., Zhang, A.D., Ramanathan, M.: CASCADE: a novel quasi all paths-based network analysis algorithm for clustering biological interactions. *BMC Bioinformatics* 9, 64 (2008)
8. Sallim, J., Abdullah, R., Khader, A.T.: ACOPIN: An ACO Algorithm with TSP Approach for Clustering Proteins from Protein Interaction Network. In: *Second UKSIM European Symposium on Computer Modeling and Simulation*, pp. 203–208 (2008)
9. Sallim, J., Abdullah, R., Khader, A.T.: An improved ant colony optimization algorithm for clustering proteins in Protein Interaction Network. In: *International Conference on Software Engineering & Computer Systems, ICSECS 2009* (2009)
10. Wu, S., Lei, X.J., Tian, J.F.: Clustering PPI Network Based on Functional Flow Model through Artificial Bee Colony Algorithm. In: *Seventh International Conference on Natural Computation*, pp. 92–96 (2011)
11. Ji, J.Z., Liu, Z.J., Zhang, A.D., Jiao, L., Liu, C.N.: A New Method based on Ant Colony Optimization for Detecting Functional Modules in Protein-Protein Interaction Networks (in press)
12. Dorigo, M., Maniezzo, V., Colnari, A.: The ant system: Optimization by a colony of cooperating agents. *IEEE Transactions on Systems, Man, and Cybernetics, Part B* 26(1), 29–41 (1996)
13. Angus, D., Woodward, C.: Multiple objective ant colony optimization. *Swarm Intelligence, Special Issue on Ant Colony Optimization* 3(1), 69–85 (2009)
14. Zhong, W.C., Liu, J., Xue, M.Z., Jiao, L.C.: A Multiagent Genetic Algorithm for Global Numerical Optimization. *IEEE Transactions on Systems, Man, and Cybernetics Part B: Cybernetics* 34(2), 1128–1141 (2004)
15. Pan, X.Y., Liu, F., Jiao, L.C.: Density sensitive based multi-agent evolutionary clustering algorithm. *Journal of Software* 21(10), 2420–2431 (2010)
16. Friedel, C.C., Krumsiek, J., Zimmer, R.: Bootstrapping the Interactome: Unsupervised Identification of Protein Complexes in Yeast. In: Vingron, M., Wong, L. (eds.) *RECOMB 2008. LNCS (LNBI)*, vol. 4955, pp. 3–16. Springer, Heidelberg (2008)
17. Mewes, H.W., et al.: MIPS: analysis and annotation of proteins from whole genomes. *Nucleic Acids Res.* 32(Database issue), D41–D44 (2004)
18. Aloy, P., et al.: Structure-based assembly of protein complexes in yeast. *Science* 303(5666), 2026–2029 (2004)
19. Dwight, S.S., et al.: *Saccharomyces Genome Database* provides secondary gene annotation using the Gene Ontology. *Nucleic Acids Research* 30(1), 69–72 (2002)
20. Li, X.L., Wu, M., Kwok, C.K., Ng, S.K.: Computational approaches for detecting protein complexes from protein interaction networks: a survey. *BMC Genomics* 11(suppl. 1), S3 (2010)

Research on Genetic Segmentation and Recognition Algorithms

Zhenjie Hou¹ and Jianhua Zhang²

¹ School of Information Science and Engineering, Changzhou University, Changzhou, China

² College of Computer and Information Engineering, Inner Mongolia Agricultural University, Hohhot, China

hhderek@163.com

Abstract. Proposed an improved genetic segmentation algorithm and recognition method base on PSO. In the genetic algorithm, 2-dimension chromosome coding is adopted; initialization of population with stochastic and symmetrical methods is produced to keep the variety of the population; OTSU is adopted to be as fitness function; a new individual is introduced in updated population. Abstracted three main components from the segmented image, and used neural networks trained by PSO to recognize the blood cells types. Experiments show that good results can be achieved steadily and quickly.

Keywords: Segmentation, Genetic algorithm, PSO.

1 Introduction

In biomedical field, a common task is to distill the features of cells image and classify and count the cells, all of which can be concluded to image segmentation. Segmentation is one of key and basic skills in the image processing and the computer vision, which is to provide the basis for the following classification, recognition and the retrieval after separating the target from background. The methods of segmentation mainly are threshold, edge detection, region tracing method and so on. In which, the method of threshold is widely used in the filed. At present, there are many methods of threshold such as minimum error threshold, OTSU, and the best histogram entropy method.

GA—Genetic Algorithm is a search algorithm which is based on natural selection and heredity theory combining the rule of “survival of the fittest” in biological process with the exchange mechanism of stochastic information of population interior chromosome. Many beneficial explorations have been carried out in the image processing fields [1~6].

The paper combined OTSU with GA and then chooses maximum variance between clusters threshold by using genetic algorithm for the segmentation of blood cells image, and used neural networks trained by PSO to recognize the blood cells types.

2 Designs of Genetic Algorithm Based on OTSU

Two problems should be solved to combine GA with OTSU [7]: the first one is how to code its solution into gene, which is individual coding method; the second one is how to construct fitness function to scale the adapting degree of every chromosome strand.

GA is a kind of adaptive global optimal probability algorithm, which has two remarkable advantages-- implicit parallelism and effective using ability. The former could improve the algorithm speed; the latter could make GA get a better robustness and also could avoid local optimization. Because the process of choosing threshold using OTSU is a process of choosing the best answer, the fast optimization characteristic of GA can be used to optimize it to enhance the efficiency.

2.1 Chromosome Coding

Decision variable is presented as the string structure data by encoding it. Because grey level of image is between 0 and 255, eight binary codes between 00000000-11111111 can be used to present as a segment threshold. Binary coding scheme is adopted in this paper.

2.2 Fitness Function

In the genetic algorithm implementation, many different chromosomes exist at the same time in each generation. Each individual fitness size in the population decides which chromosomes should be inherited to the next generation [7]. The size of fitness can be gained by calculating the value of fitness function and called fitness [8]. The paper takes OTSU as fitness function $F(t)$. Suppose that value t of the threshold divides image into two parts: C_0 and C_1 (object and background), thus C_0 and C_1 separately corresponds to the pixels which gray levels are $\{0,1,\dots,t\}$ and $\{t+1,t+2,\dots,L-1\}$. Suppose that $\sigma_B^2(t)$ represents classes square error when the value of threshold is t in histogram, and then optimum threshold can be obtained by searching maximum of $\sigma_B^2(t)$. The variance between C_0 and C_1 can be calculated as follows:

$$\sigma_B^2(t) = w_1(t) \times w_2(t) \times (\mu_1(t) - \mu_2(t))^2 \tag{1}$$

Where, $w_1(t)$ is the number of pixels in C_0 ; $w_2(t)$ is the number of pixels in C_1 ; $\mu_1(t)$ is the average gray value of all numbers of the pixels; $\mu_2(t)$ is the average gray value of all numbers of the pixels. The value between 0 to $L-1$ should be changed into t , then t which value of δ is the biggest one should be taken as the best threshold T .

2.3 Genetic Operator

Genetic operator mainly consists of 3 types: select, crossover and mutation.

Select operator: choice is made by the method of roulette.

Crossover operator: crossover operator is a significant character for genetic algorithm, which is thus differentiated for other algorithm. It is a main process from which a new unit is produced. The design and application of crossover operator are related closely

to the problems which to be studied. Crossover operator determines the general search capability of genetic algorithm. Generally, it should not disturb the excellent model in individual encoding series, and is able to generate some excellent new individual model. The chromosome in this essay adopts the method of alone point crossing.

Mutation operator: mutation operator is only an auxiliary method for the generation of new unit. It determines the partial search ability of genetic algorithm. In this essay, mutation operation use basic byte mutation operator, that is, within binary coding system, “0” changes to “1”, and “1” changes to “0”.

Introduction of new unit

To keep the variety of group, and to avoid partial pre-mature, a special operator is introduced. Considering the uniqueness of image data, the operator is defined as: when parental generation chromosome produce filial generation group C1 through crossover and mutation, the average value of chromosome in C1 is taken as genetic value, and thus a new unit (X_{new}), is defined as followed:

$$X_{new} = \frac{1}{n} \sum_{i=1}^n X_i \tag{2}$$

$X_i (i = 1, 2 \dots n)$ is the genetic value of units in C_1 .

2.4 Choose Controlling Parameters

Controlling parameters mainly includes the population size and probability of genetic operation etc. In GA, crossover and mutation are genetic operators to produce new individuals where the former plays a main part in maintaining the population diversity and the latter as subsidiary. That means a relatively large crossover probability and small mutation probability should be chosen.

This paper adopts adaptive crossover probability P_c and mutation probability P_m

$$P_c = \begin{cases} P_{c1} & f' < f_{avg} \\ P_{c1} - \frac{(P_{c1} - P_{c2})(f' - f_{avg})}{f_{max} - f_{avg}} & f' \geq f_{avg} \end{cases} \tag{3}$$

$$P_m = \begin{cases} P_{m1} & f < f_{avg} \\ P_{m1} - \frac{(P_{m1} - P_{m2})(f_{max} - f)}{f_{max} - f_{avg}} & f \geq f_{avg} \end{cases} \tag{4}$$

In the formula above, f_{max} is the maximum fitness value in group; f_{avg} is the average fitness value in every generation; f' is the fitness value of the two units that will crossover; f is the fitness value of the unit that will take mutation; P_{c1} is the

fixed maximum crossover rate, 0.9; P_{c2} is the fixed minimum crossover rate, 0.6; p_{m1} is the fixed maximum mutation rate, 0.1; P_{m2} is the fixed rate, 0.001.

2.5 Termination Criteria

As is prescribed, when algorithm operates to the maximum generation number or after 20-generation evolution, the maximum fitness in population has not changed, the algorithm comes to a stop. Then the threshold value referred to by the unit with the maximum fitness is what is needed. In this essay, the maximum generation number is 50.

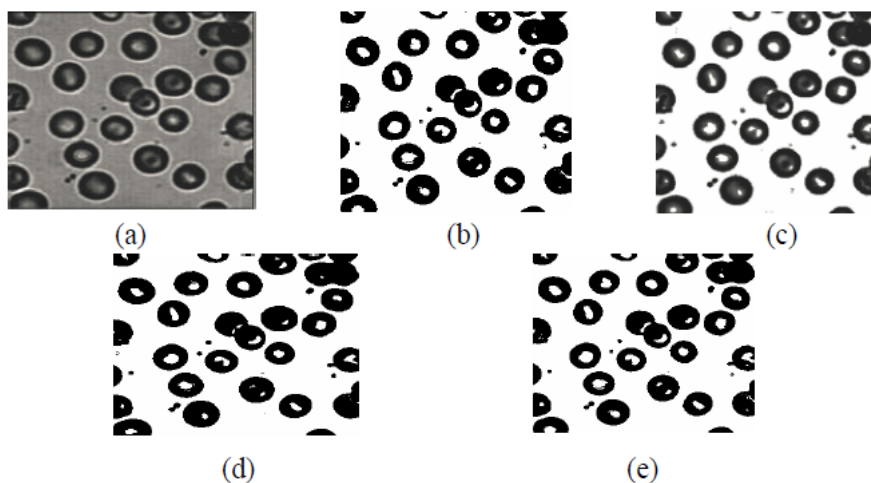


Fig. 1. Contrast results. (a)Original blood (b) Result of optimal threshold segmentation (c) Result of cells image adaptive threshold segmentation (d)Result of Otsu segmentation (e) Result of improved GA segmentation.

3 Experiments of the Segmentation

The algorithm model in this essay and related solution algorithm is realized by Visual C++ program. To segment image by generic algorithm based on Otsu is to calculate Otsu and related threshold value by fast optimization genetic algorithm. The contrasts between the improved GA algorithm and other classic algorithm are shown in Fig.1.

With the consistency, contrast, and calculation cost (time) as evaluation function, quality of image segmentation is tested (shown in table 1).

From the result of experiment, the revised genetic algorithm in this essay has the following characteristics:

Table 1. Results of performance evaluation

Method	Segment threshold	Region homogeneity	Region contrast	Time(second)
Optimal threshold	93	0.9940	0.5240	0.016
adaptive threshold	102	0.9939	0.5060	0.016
otsu	100	0.9939	0.5060	0.016
improved GA	83	0.9937	0.5423	0.015

(1) The revised genetic algorithm based on Otsu is easy to understand, and the algorithm is less complex.

(2) The revised genetic algorithm based on Otsu can separate the target from backdrop and keep details of image, as is shown in table 2.

(3) Partial consistency approaches 1, which indicates the internal elements in the segmentation area of this algorithm are similar, which means the segmentation is highly effective.

(4) The revised genetic algorithm based on Otsu has a relatively high contrast, which means the segmentation method is improved.

(5) The revised genetic algorithm based on Otsu need o calculates less in shorter operation time.

After segmentation, cell's characteristics are acquired from the karyon parts and use the neural networks to classify the types.

4 Identification of Nerve Network Cell Type Based on Particle Swarm Optimization

Cell's characteristic can be acquired by abstracting its three main components. The first main component keeps the comprehensive information of its original image. The second main component reflects the information about the contrast of original image. The third main component signifies information about brightness.

Particle swarm optimization is a swarm intelligence evolutionary computation based on group and fitness. Each individual in group stands for a possible solution. Particle is characterized with speed and space, and the algorithm measures the quality of particle by fitness [9-11].

Supposing in a search space of D-dimension, particles of m number form a group, among with the particle "i" is a vector of D-dimension, which is marked as $\vec{x}_i = (x_{i1}, x_{i2}, \dots, x_{iD}), i = 1, 2, \dots, m$, that is, the position for particle "i" in D-dimension is xi. The position of every particle is a possible solution. The fitness value of Xi can be calculated by taking Xi into a target function built on purpose. Then the quality of Xi can be measured by the fitness value. The speed of particle "i" is also a vector in D-dimension, which is marked as: $\vec{v}_i = (v_{i1}, v_{i2}, \dots, v_{iD})$. The particle I find its optimistic position so far is $\vec{p}_i = (p_{i1}, p_{i2}, \dots, p_{iD})$, also named as pbest; the whole

particle group find their optimistic position so far is $\bar{p}_g = (p_{g1}, p_{g2}, \dots, p_{gd})$, also named as gbest. The formula of PSO is as followed:

$$\begin{aligned} v_{id}^{k+1} &= v_{id}^k + c_1 r_1 (p_{id}^k - x_{id}^k) + c_2 r_2 (p_{gd}^k - x_{id}^k) \\ x_{id}^{k+1} &= x_{id}^k + v_{id}^{k+1} \end{aligned} \tag{5}$$

$i=1,2,\dots,m, d=1,2,\dots,D$; K is iterative number; the study gene c_1 and c_2 is nonnegative constant; r_1 and r_2 is a random number between $[0,1]$; x_{id}^k is the position vector of the particle I ; v_{id}^k is the speed vector of the particle I ; $v_{id}^k \in [-v_{\max}, v_{\max}]$, v_{\max} is constant, settled according to different situation; p_{id}^k is partial optimized value, p_{gd}^k is the general optimized value.

The formula consists of three parts: the first part is the current speed of particle, indicating the current state of particle, the bigger W value at early search stage helping jump out of partial minimum point, while a smaller W value at late search stage helping algorithm convergence; the second part is individual cognitive part, which enables particles with strong general search capability, so as to avoid local minimal; the third part is the social information pool, which helps particles draw experience form other excellent particles, so as to enhance search capability.

For different issues, iterative terminal condition usually selects the maximum iterative number, or the searched optimized position satisfying the minimum fitness threshold value.

The network adopted here contains a hidden level. The input is the three main parts achieved form the former network; the output is coded according to the cell type number. When BP network is trained by PSO, the defined particle's position x_k is the threshold value with proportion and hidden level node under a complete connection structure, with the proportion range between $(-1, 1)$. The network training adopts gradual increasing to increase the number of hidden level nodes to a satisfying degree. Fitness value function is an index for evaluation of nerve network's ability to solve problems. The mean square error of nerve network output works as target function, and its reciprocal works as fitness function. The smaller the error is, the better the function of the referred particle is.

Through experiment, the following parameter is settled: the terminal condition of PSO training BP network is iterative number $K=1000$; $k_1=0.5$; $P_c=0.03$; $P_m=0.01$. the training sample collection consists of 160 cell image of the size 32 by 32; testing samples are 90 cell images, with the average correct rate 91.3%. (Shown in table II).

5 Conclusion

As is concluded, when the revised genetic algorithm is applied to image segmentation based on Otsu, the value of Otsu and related grey level threshold are a quick and stable non-linear solution, thus the image threshold segmentation can be improved in speed and function. Abstracted three main components from the segmented image, and used neural networks trained by PSO to recognize the blood cells types. Experiments show that good results can be achieved steadily and quickly.

References

1. Andrey, P.: Selectionist relaxation: genetic algorithms applied to image segmentation. *Image and Vision Computing* 17(3-4), 175–187 (1999)
2. Van Coillie, F.M.B.: Feature selection by genetic algorithms in object-based classification of IKONOS imagery for forest mapping in Flanders. *Remote Sensing of Environment* 110(4), 476–487 (2007)
3. Tseng, D.-C., Lai, C.-C.: A genetic algorithm for MRF-based segmentation of multi-spectral textured images. *Pattern Recognition Letters* 20(14), 1499–1510 (1999)
4. Angelié, E., de Koning, P.J.H.: Automatic tuning of left ventricular segmentation of MR images using genetic algorithms. *International Congress Series*, vol. 1256, pp. 1102–1107 (June 2003)
5. Kim, E.Y., Park, S.H.: Automatic video segmentation using genetic algorithms. *Pattern Recognition Letters* 27(11), 1252–1265 (2006)
6. Kim, E.Y., Jung, K.: Genetic algorithms for video segmentation. *Pattern Recognition* 38(1), 59–73 (2005)
7. Hou, Z., Ma, S.: Study on segmentation of marrow cells image based on GA. *Journal of Anhui Agricultural University* 32(4), 551–554 (2005)
8. Otsu, N.: A threshold selection method from gray level histogram. *IEEE Trans. Systems Man Cybernet.* 9(1), 62–66 (1979)
9. Hao, M., Ma, S., Hao, X., Ma, L., Wang, L.: Feature selection based on GA and PNN. *Advanced Materials Research*, 1753–1757 (2011)
10. Zhang, P., Yu, Z.: Mechanism of Eggs Classification Based on Machine Vision System. In: *MACE 2010*, Wuhan, China, pp. 5718–5720 (2010)
11. Taghizadeh, M., Gowen, A., O'Donnell, C.P.: Prediction of white button mushroom moisture content using hyperspectral imaging. *Sensing and Instrumentation for Food Quality and Safety* 3, 219–226 (2009)

Blog-Based Distributed Computation

Implementation of Software Verification System

Takayuki Sasajima and Shin-ya Nishizaki

Department of Computer Science, Tokyo Institute of Technology, 2-1-2-W8-69, O-okayama,
Meguro-ku, Tokyo, 152-8552, Tokyo, Japan
nisizaki@cs.titech.ac.jp

Abstract. Nowadays, blogs are regarded as standard text-based communication tools on the internet. In contrast to traditional web pages, the blog has several significant features: it allows authoring via a web browser, and offers automatic backlink requesting, called *trackback*. The latter provides direct communication between blog servers. In this paper, we propose a distributed computation method based on *trackback* communication and present the implementation of a software verification system based on distributed computing. The software verification system consists of ordinary blog systems which are used as frontend interfaces, and verification blog bots which are used as backend inference engines. One of the prominent features is that one can limit the danger of intrusion into vulnerable verifier programs to the private networks where the verification blog bots are operated.

1 Introduction

The word *blog* comes from the term *Web log*, and it is a kind of internet web site. On the Web, articles are displayed in chronological order and we can edit and submit an article to the Web without using an authoring tool such as an HTML editor. Articles on the Web can contain not only text but also multimedia data such as pictures and videos. Today, several kinds of blog publication software are available, for example, *Movable Type* [1] and *Word Press* [2].

Trackback is one of the backlink methods of the Web, and was originally implemented in *Movable Type* [1]. Now it is one of the standard features of the blog. The following is an overview of *trackback*.

Assume that a blog author intends to get a backward link to his/her article from an article in another blog. First, the author specifies a *trackback URL* to his/her referring blog article, which is assigned to the referred blog article.

The blog server A accepts the *trackback URL* from the author and sends a *trackback ping* to the blog server B which carries the referred article.

If the blog server B receives the *trackback ping*, then A makes a *trackback link* from the referred article to the referring article and sends an acknowledgement message to B.

The *trackback protocol* is a de-facto standard for many blog systems.

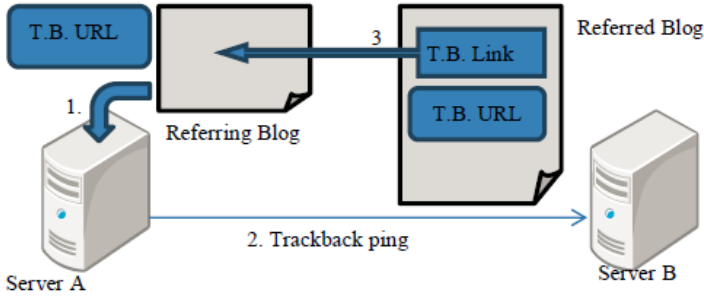


Fig. 1. Overview of Trackback Protocol

Internet bots (or simply bots) are software applications that perform automated tasks. There are several kinds of bots: web crawlers, which browse and analyze web pages systematically, chat-bots which make conversation, disguising themselves as human beings, and gaming bots which participate in network-based games. Bots provide services to others in the network. However, the bots are not categorized as server programs: servers respond to clients according to their requests, but the bots do not wait for requests from others. The public server programs must reply to public requests, which causes security vulnerability. On the other hand, bots can be located on private networks which are connected to the internet via network routers. Thus, bots have far less danger of intrusion caused by vulnerability.

In this paper, because of this advantage with respect to security, we propose a new software system structure in which blog systems act as front-end interfaces and bots provide computation resources as the back-end. The blog systems communicate with each other via the trackback protocol. We call this software structure Blog-based distributed computation. As this instance, we carry out a distributed implementation of software verification system based on the blog-based distributed computation.

2 Blog-Based Distributed Computation

In this section, we give an overview of blog-based distributed computation. Blog-based distributed computation consists of blog systems and bots. The service in this system is open to anonymous clients; it can be restricted to specific ones by using an authentication method if required. The clients are assumed to be ordinary blog systems, which are called client-side blogs. The counterparts of the client-side blogs are also ordinary blog systems, called reception blogs, which receive requests from client-side blogs. Service requests from the client-side blogs to the reception blogs are based on the trackback protocol. These two kinds of blogs are not specifically implemented but powered by ordinary blog publishing software such as Movable Type. The blogs can be hosted by free blog hosting services.

Services in blog-based distributed computation are mainly provided by service bots. A service bot reads the content of a service request published in a reception blog; it executes the service designated by the request; then it publishes the result on

another blog, called a report blog. The service bots can be installed on computers on private networks connected to the internet via network routes, such as home PCs. If a program which is embedded in a service bot and which provides an intended service has some security vulnerability, we can limit the damage inside the private network. This is not surprising because many home PCs connected to the internet are infected with viruses, but public blog hosting services are usually protected from viruses.

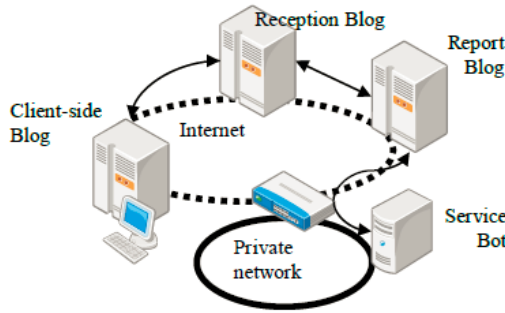


Fig. 2. Overview of Blog-based Distributed Computation

The following is the sequence of procedures from the service request by the client-side blog, to informing by the report blog.

A contributor submits a service request as an article to a client-side blog.(Fig. 3)

The contributor gives an instruction to make a trackback link from the article to an article in a reception blog; consequently, the reception blog sends a trackback ping to the reception blog. The article in the reception blog, called the reception article, is prepared in advance to receive requests for trackback links.(Fig. 4)

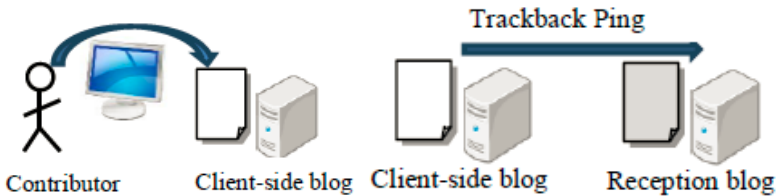


Fig. 3. Procedure 1

Fig. 4. Procedure 2



Fig. 5. Procedure 4(a)(b)

Fig. 6. Procedure 3

Once the reception blog has received the trackback ping, it establishes a trackback link from the article in the reception blog to the article in the client-side blog. (Fig. 6)

A service bot accesses the blogs as follows.

The service bot obtains the article in the client-side blog by following the trackback link in the article in the reception blog. Then the bot applies the service-providing program to the article as its input. If the service is assumed to be software verification, the program is a software verifier and the article in the client-side blog is a code to be verified. (Fig. 5)

The service bot submits the result as an article to the report blog. (Fig. 5)

The service bot sends a trackback ping to the report blog. The ping requests the report blog to make a trackback link to the client-side blog. In other words, the service bot sends the trackback ping, impersonating the report blog

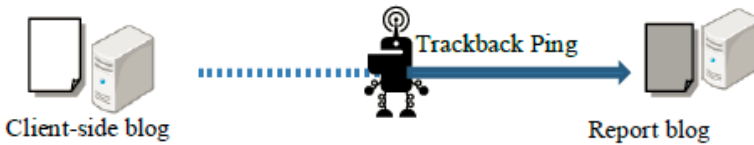


Fig. 7. Procedure 4(c)

Similarly, the service bot sends a trackback ping to the client-side blog. The ping requests the client-side blog to make a trackback link to the report blog.

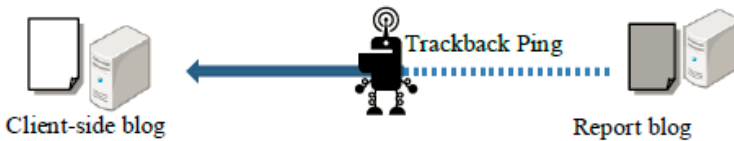


Fig. 8. Procedure 4(d)

Finally, the service bot deletes the trackback link in the reception blog. As a result of these procedures, the blogs reach the following state. (Fig. 9)

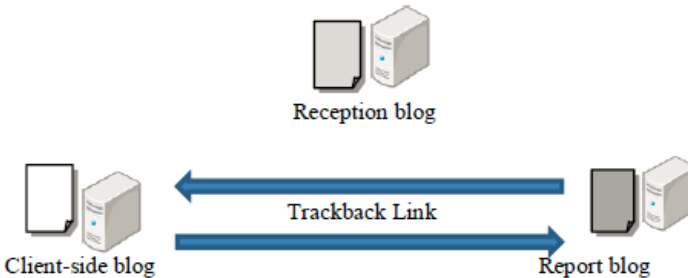


Fig. 9. Result state

In the client-side blog, the article of the service request from the contributor is published.

In the reception blog, the reception article is unchanged.

In the report blog, the article submitted by the service bot is published.

The two trackback links between these two articles in the client-side and report blogs are established.

We discuss the merits and demerits of the blog-based distributed computation below.

Merits. The input interface is implemented as a reception blog and the output interface as a report blog. Most of the system computation is done by service bots. These blogs can be powered by ordinary blog publishing software which provides high security. Programs embedded in service bots may have security vulnerability. Even if they are infected with a computer virus and hijacked, the damage can be limited to the private network in which the bots are located.

Demerits. The reception blogs cannot notify the service requests directly to software bots. The service bots have to wait for requests by polling the reception blogs. Such indirect communication between the report blogs and the service bots causes delays in response. However, if the requirement for responsiveness is not so exacting, for example, in the case of software verification, our approach is sufficiently realistic.

3 Implementation of Software Verification System Based on Blog-Based Distributed Computation

In this section, we implement a software verification system, more specifically, a model checking system, using the blog-based distributed computation introduced in the last section.

If a model of a system to be verified is given, a model checker automatically tests whether this model meets a given specification. If the model contains non-determinism, then the model checker explores exhaustively all possible behaviors of the system. We give an outline of the components introduced in the last section.

Client-side and Report blog. In our implementation, we use a free blog hosting service for the client-side blog, Seesaa Blog (<http://blog.seesaa.jp/>). Similarly to other blog hosting services, one can submit articles to the blog via the XML-RPC protocol.

Reception Blog. The service bot obtains a service request submitted by the contributor. The service bot finds the article of service request by following and analyzing the trackback link to the client-side blog. We use a free web hosting service, FC2 Blog (<http://fc2.com/>), since an id attribute is assigned to an HTML-li entry of each trackback link to the client-side blog, has a certain fixed format in FC2 Blog and it makes easy to follow the link.

Service Bot. The service bot analyzes pages in the reception blog and the client-side blog, submits the results to the report blog and sends trackback pings both to the client-side and report blogs. We chose the PHP scripting language for implementing

the service bot since it has a library for the XML-RPC protocol, its programs are easy to write, and many public web hosting services support PHP scripting.

In this implementation, we choose Seesaa Blog and FC2 Blog. We check whether we can use other public blog hosting services for the servers of our system or not. We investigate nine popular blog hosting services in Japan and show acceptability of trackback ping from the client-side blog to the reception blog (C→Rec.) and from the reception blog to the client-side blogs (Rep. → C), in Table 1.

Service	C→Rec.	Rep.→C	Service	C→Rec.	Rep.→C
Seesaa	OK	OK	goo	OK	OK
FC2	OK	OK	Yahoo!	OK	NG
DTI	OK	NG	Cocolog	NG	OK
livedoor	OK	OK	Excite	OK	NG
Ameba	OK	NG			

The URL of DTI blog, livedoor blog, Ameba blog, goo blog, Yahoo! Japan blog, and Excite blog are <http://www.dtiblog.com/>, <http://blog.livedoor.com/>, <http://ameblo.jp/>, <http://blog.goo.ne.jp/>, <http://blogs.yahoo.co.jp/>, <http://www.cocolog-nifty.com/>, and <http://exblog.jp/>, respectively.

Though our implementation does not work with some of the blog hosting services, it has quite good compatibility with many of them.

4 Conclusion

In this paper, we proposed a blog-based distributed computation in which blog systems are used for input and output interfaces, and the backend for the service computation is provided as bots. We show its feasibility by implementing a software verification system based on the distributed computation.

5 Related Works

A distributed system [3] consists of physically separate computers that communicate through a computer network. The computers communicate with each other in order to archive a common goal. Distributed services such as NFS [4] and DNS [5] are dispensable infrastructure of the internet. Communication protocols used by such distributed network services are provided on Layer-5 (i.e. the session layer) or Layer-7 (i.e. the application layer). The communication protocol of our blog-based distributed computation is defined on trackback protocol which is defined on HTTP defined on Layer-7. Hence, we can say that the communication is exchanged on a much higher layer than Layer-7.

Recently, many researchers have made trials to develop an open infrastructure of computation based on distributed computation. The most famous example is SETI@home [6], which provides an internet-based public volunteer computing in order to search for Extra-Terrestrial Intelligence using PCs at home. We can expect

our blog-based distributed computation to provide us with a more secure and more open infrastructure for volunteer computing.

6 Future Works

The implementation presented in this paper is currently operated on stand-alone PCs. Recently, cloud computing has started to become widely regarded as the infrastructure solution for computationally-scalable requirements. One of the most famous cloud services is Amazon EC2 (Amazon Elastic Compute Cloud). The Amazon EC2 is categorized as IaaS (Infrastructure as a Service), and delivers a platform virtualization environment as a service, along with storage and networking. Such cloud infrastructure can be used for our distributed computation, which is one of the important continuations of this research.

The blog-based distributed computation is designed to be open; in particular, it is assumed that service bots can easily participate in the distributed system. Although it promises to provide us with scalable computational resource, it may cause Denial-of-Service vulnerability. We should therefore investigate resistance against Denial-of-Service attacks [7, 8].

References

1. The Official Web Site of Movable Type, <http://www.movabletype.org/>
2. The Official Web Site of WordPress.com, <http://wordpress.org/>
3. Andrews, G.R.: Foundations of Multithreaded, Parallel and Distributed Programming. Addison-Wesley (2000)
4. RFC 3530, Network File System (NFS) Version 4 Protocol, <http://www.ietf.org/rfc/rfc3530.txt>
5. RFC 1034, Domain Names – Concept and Facilities, <http://www.ietf.org/rfc/rfc1034.txt>
6. SETI@home, <http://setiathome.berkeley.edu/>
7. Tomioka, D., Nishizaki, S., Ikeda, R.: A Cost Estimation Calculus for Analyzing the Resistance to Denial-of-Service Attack. In: Futatsugi, K., Mizoguchi, F., Yonezaki, N. (eds.) ISSS 2003. LNCS, vol. 3233, pp. 25–44. Springer, Heidelberg (2004)
8. Ikeda, R., Narita, K., Nishizaki, S.: Cooperation of Model Checking and Network Simulation for Cost Analysis of Distributed Systems. International Journal of Computers and Applications 33(4), 323–329 (2011)
9. Homepage of XML-RPC, <http://xmlrpc.scripting.com/default.html>

Model Transformation Method for Compensation Events and Tasks from Business Process Model to Flowchart

Jian Deng^{1,*}, Bo Chen², and Jiazhi Zeng²

¹ School of Information & Software Engineering, University of Electronic Science & Technology of China, Chengdu, China

² School of Computer Science & Engineering, University of Electronic Science & Technology of China, Chengdu, China
18980891251@189.cn, jockeydj@gmail.com,
{bochen, jzzeng}@uestc.edu.cn

Abstract. To satisfy the requirement of developing software rapidly for business processes with compensation, a special kind of computation independent model with compensation named well-structured associated business process model was researched. The model mapping method for compensation events and tasks was also proposed. The computation independent model was designed using XML Process Definition Language 2.1 version. It can be directly mapped to flowchart model in Windows Workflow Foundation 4.0 version, and the tool for model conversion and testing was implemented on the Windows platform. Experiment results shows that this method satisfies the functional requirement of processes compensation, and improves the speed of software development for business processes.

Keywords: Compensation modeling, business process model, flowchart model, XPDL, WF.

1 Introduction

Due to market fluctuations, policy changes, technology upgrades and other factors, users need business processes developed quickly to meet the changing requirements. This has made software development face enormous challenges, so the Model Driven Architecture (MDA) [1] has become the research focus in software engineering [2-4].

The computation independent model of business processes can be illustrated by Business Process Modeling Notation (BPMN)[5], and designed using XML Process Definition Language (XPDL)[6] proposed by Workflow Management Coalition (WfMC). In this paper, XPDL model means BPMN model in the file format of XPDL. In the situation of a business process having completed but needing to eliminate its effects, we can model this situation with compensation.

* Corresponding author.

The flowchart in Windows Workflow Foundation (WF) [7] 4.0 can be used to design the platform specific model. It supports designing and executing a process model with compensation on .NET platform.

One of the key technologies of MDA [8] is automatic model conversion. To our knowledge, there has been no published conversion method for models with compensation from XPDL to flowchart in WF. To solve the problem, this paper presents a method to transform XPDL to WF. The business process with compensation is modeled using XPDL 2.1 version, and converted to flowchart in WF 4.0 directly.

The rest of the paper is organized as the followings: Section 2 introduces the basic definitions in business process modeling. Section 3 presents the method to convert models from XPDL to WF. Section 4 gives insight into how the conversion technique can be applied to a typical case. Section 5 describes related work in this area. Finally Section 6 draws the conclusion and previews future research content.

2 Business Process Model

XPDL provides a set of business process graphical symbols that are readily understandable by business users, system analysts and technical developers. To facilitate the model transformation, the definition of process in reference [9] can be extended as the followings.

Definition 1. A process in a business process model is an eight-tuple $P = (E, T, S, G, O, D, C, \varphi)$, where

E is a set of events which can be partitioned into disjoint sets of start events E_S , intermediate events E_I and end events E_E ,

T is a set of tasks,

S is a set of sub-processes that can be mapped to processes,

G is a set of gateways used to model the convergence and divergence in a process,

O is a set of pools,

D is a set of data objects,

C is a set of connections which can be partitioned into disjoint sets of sequence flow C_S , association flow C_A and message flow C_M ,

φ is a set of function mappings from sub-processes to processes.

Compensation events and association flows in Definition 1 will be used to model business process with compensation.

Windows Workflow Foundation 4.0 was released by Microsoft in 2010. Developers can use flowcharts to model the business process. According the features of the flowchart in WF, it can be described with Definition 2.

Definition 2. A flowchart in Workflow Foundation is a directed graph $F = \{N, C\}$, where N is a set of nodes, C is a set of connections, $C \subseteq N \times N$.

There are three types of nodes in the flowchart: FlowStep, FlowDecision and FlowSwitch. FlowStep is used to perform a specified activity. The Action attribute of

it can be an instance of a basic activity, sequential activity, flowchart activity, or any class derived from `System.Activities.Activity`. The `Next` attribute of `FlowStep` is used to specify an unconditional successor node.

`FlowDecision` has at most two conditional successor nodes. It evaluates the expression in its `Condition` attribute, and then selects to execute the node specified in the `True` or `False` attribute.

There may be many conditional successor nodes in `FlowSwitch`. `FlowSwitch` uses `Expression` attribute to select which node in its `Default` and `Cases` attributes to execute.

3 Mapping XPDL to WF

This section discusses the methods to map data objects, compensation events, tasks with compensation, parallel gateways and exclusive gateways in XPDL to flowchart in WF.

3.1 Mapping of Data Objects

For the purpose of modeling the scripts in tasks and expressions in gateways, data objects can be defined in XPDL model. In this paper, the type of data object is a string stored in its `Description` attribute. A data object in XPDL can be mapped to a variable in WF. Table 1 show the data types supported currently.

Table 1. Mapping of data type

XPDL model	WF model
INTEGER	int
FLOAT	float
BOOLEAN	Boolean
STRING	String
DATE	DateTime
TIME	DateTime
DATETIME	DateTime
CompensationToken	CompensationToken

In a flowchart model, all variables should be added into the *Variables* attribute of the flowchart so that they can be accessed by the activities in the model.

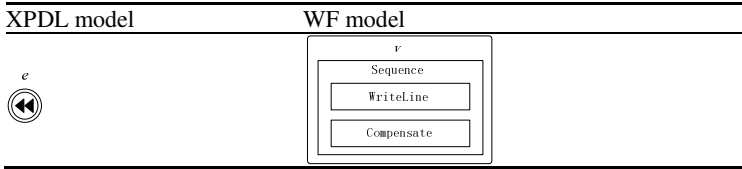
3.2 Mapping of Compensation Events

In XPDL model, compensation events are used to indicate that compensation is necessary. A compensation event can be thrown or caught. To throw a compensation event, the event may be an intermediate event or end event. To catch a compensation event, the event must be an intermediate event at the boundary of the task which needs to be compensated.

In WF model, the activities which relate to compensation are *Compensate* and *CompensableActivity*. *Compensate* is the activity that will call the *CompensationHandler* of *CompensableActivity*. A variable type of *CompensationToken* is used to establish the correspondence of *Compensate* and *CompensableActivity*. Table 2 shows the mapping method for throwing compensation event.

For any other type of events, it will be mapped to a sequence consisting of only *WriteLine* activity currently.

Table 2. Mapping of throw compensation event



3.3 Mapping of Tasks with Compensation

Table 3 shows part of a process with compensation in XPDL model and the corresponding WF model. In XPDL model, there should be an intermediate event *e1* at the boundary of task *A0* which needs to be compensated. The *Target* attribute of *e1* is the ID of *A0*. Task *A1* is used to compensate task *A0*, and it is connected with event *e1* by association connection *C1*. When task *A0* has completed and event *e1* is caught, task *A1* should be executed. The *A0*, *e1*, *C1*, *A1* in XPDL model can be implemented by a compensable activity *v* in WF model. The *Body* attribute of *v* is used to model the action of task *A0*, *CompensationHandler* attribute of *v* is used to convert *e1* and *A1*.

To model the function of task *A0*, some script text can be stored in *Description* or *Implementation*. *TaskScript* attribute of it. For example, there are two lines of VisualBasic scripts in Table 4. It can be converted to activity *v* which contains a sequence of one *WriteLine* activity and two *Assign* activities in WF model.

When there is no script text in task *A0*, it will be mapped to a sequence consisting of only *WriteLine* activity currently.

Table 3. Mapping of task with compensation

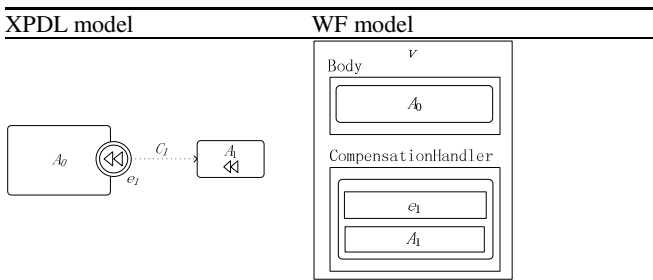
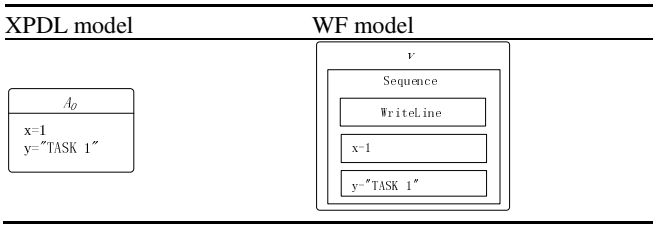


Table 4. Mapping of normal task



3.4 Mapping of Parallel Gateways

Parallel gateways are used to model the creation and synchronization of parallel flows in XPDL model. In this paper, we assume that there is only one task in each branch of parallel flow so that it can be mapped to WF model easily.

Table 5. Mapping of Parallel gateways

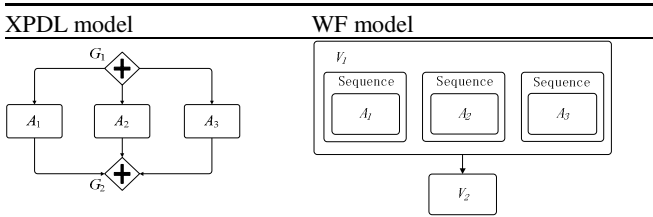
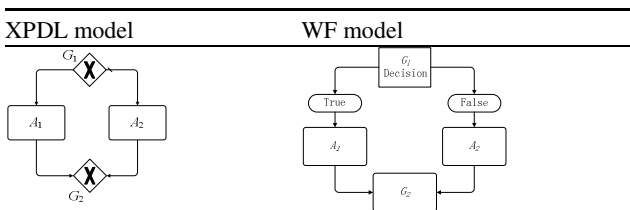


Table 5 shows that the divergent gateway G_1 can be mapped to a Parallel activity V_1 in WF. The Branches attribute of V_1 is used to map each branch of G_1 . Convergent gateway G_2 can be mapped to V_2 which contains one WriteLine activity in WF.

3.5 Mapping of Exclusive Gateways

In XPDL model, Exclusive divergent gateways are used when only one path can be taken from two or more alternative paths. When there are only two alternative paths, the exclusive divergent gateway in XPDL model can be mapped to a *FlowDecision* in WF model. The exclusive convergent gateway can be mapped to a *FlowNode* in WF.

Table 6. Mapping of Exclusive gateways



4 Case Study

In this paper, we use TIBCO Business Studio [10] 3.1 version to design XPD model. By extending the workflow designer provided by reference [11], we have implemented a lab tool named XPD2WF4. The tool converts XPD model to WF model directly. Users can use this tool to view and execute the result WF model.

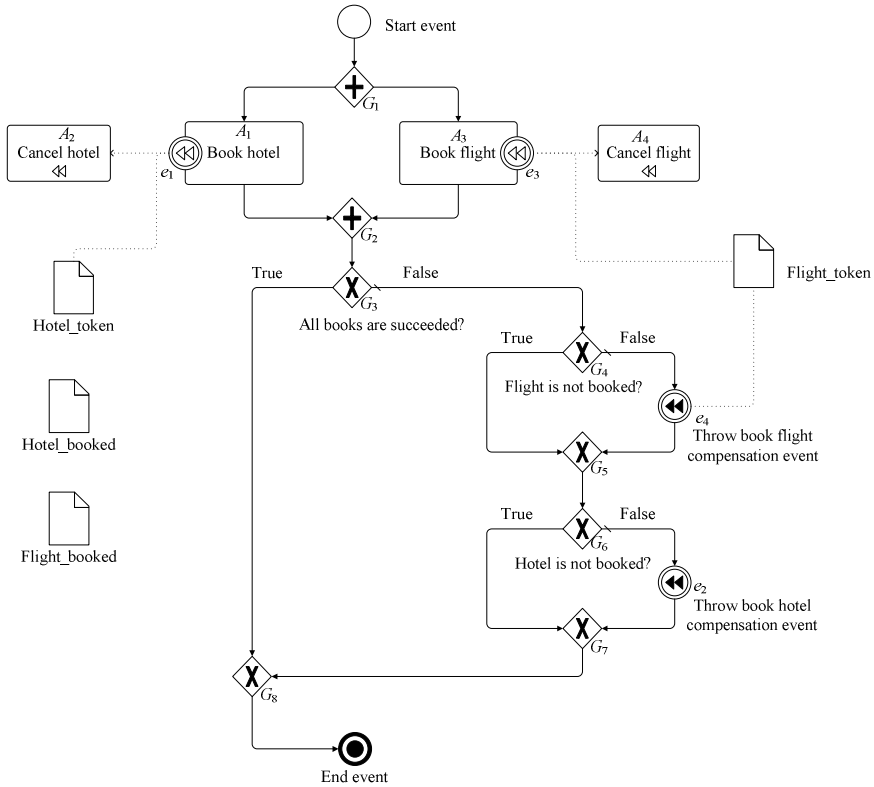


Fig. 1. Process model with compensation

Fig. 1 shows a business process model for a travel agency. It contains two activities with catching compensation event. The compensation tokens are modeled by data objects. To model the correspondence of throwing and catching compensation events, we can use association connection or store the name of token into Description attribute of event. For example, e3 and e4 are both connected with Flight_token. Event e1 is connected with data object Hotel_token, and the string “Hotel_token” is store in the Description attribute of event e2.

In order to execute the flowchart model, two data objects Hotel_booked and Flight_booked are designed, and the values of Description attribute for elements A1, A3, G3, G4 and G6 are assigned according to Table 7.

Table 7. Description attribute of some elements in Fig. 1.

Name of element	Description attribute of element
A ₁	Hotel_booked = true
A ₃	Flight_booked = false
G ₃	Hotel_booked = true and Flight_booked = true
G ₄	Flight_booked = false
G ₆	Hotel_booked = false

Fig. 2 shows the workflow model converted from Fig 1. We can execute the workflow model in our XPDL2WF tool. When the body of task A1 and A3 are completed, the WF 4 execute engine will set the results into Hotel_token, Flight_token respectively, the task A2 “Cancel hotel” will be executed because the Target attribute of e2 is Hotel_token.

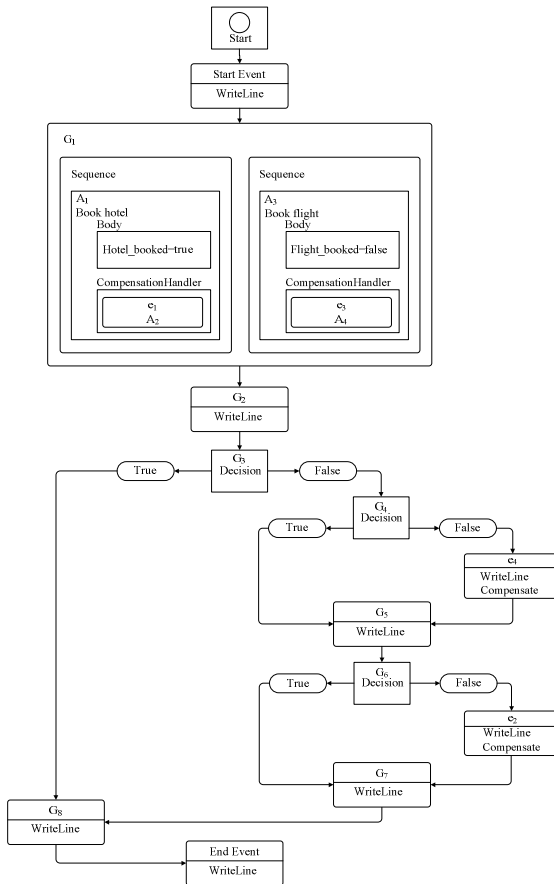


Fig. 2. The flowchart model mapped from Fig.1

5 Related Works

In the domain of business process modeling, BPMN can be used to design computation independent models (CIM); Business Process Execute Language (BPEL) [12] can be used to design platform independent models (PIM); and Microsoft's Windows Workflow Foundation (WF) can be used to design platform specific models (PSM).

Since there is no published method to convert BPEL models to WF models, the conversion from XPDL to BPEL and to WF does not exist. To our knowledge, the graphical layout information of elements was not specified in BPEL standard. In order to preserve the graphical layout information in PSM model so that it can be checked visually, we consider converting XPDL to WF directly.

According to the research works of reference [13,14], the XPDL model can be mapped to WF model. Lin [15] has developed a method to transform the XPDL model into sequential model of WF 3.5.

The concept of compensation was proposed by Gray[16] and widely adopted by standards such as BPMN, XPDL, BPEL and WF. Since there are so many modeling elements in the specification of XPDL, Lin [15] didn't discuss how to transform business process model with compensation.

To implement the model transformation from XPDL to WF, we consider using WF 4.0 because the framework of workflow has been totally reconstructed. There are many features that are improved or added. For example, the WriteLine, Flowchart, FlowDecision, FlowSwitch activities are introduced in the new version of Workflow Foundation. The technology of WF 3.5 is not updated any more.

6 Conclusion

This paper presents a model transformation method for business process with compensations. The innovation is that it converts the business process model with compensation in XPDL 2.1 version to the flowchart model in Windows Workflow Foundation 4.0. Details of the transformation methods for data objects, compensation events and tasks, and gateways are discussed in this paper. The lab tool is implemented and can be used to view and execute the business processes.

Future researches can include: (1) business processes combined with transactions and compensations; (2) the conversion of message links. Further studies may also include converting from WF to XPDL, so that we can convert the models bi-directionally.

Acknowledgments. This paper is supported by the project "Enterprise resource cooperation and research of key technologies" (No. 2007AA040801) and the project "International trade and regional economic cooperation platform and research of key technologies" (No.2009BAH46B03). The authors are extremely grateful to Professor Sun Linfu, Associate Professor Wang Shuying and Han Min for their help.

References

1. Soley, R.: Model Driven Architecture, <http://www.omg.org/cgi-bin/doc?omg/00-11-05.pdf>
2. Meservy, T., Fenstermacher, K.: Transforming software development: An MDA road map. *Computer* 38(9), 52–58 (2005)
3. Braganca, A., Machado, R.J.: A model-driven approach for the derivation of architectural requirements of software product lines. *Innovations in Systems and Software Engineering* 5(1), 65–78 (2009)
4. Jindan, F., Dechen, Z., Lanshun, N., et al.: Modeling business object platform independent model and its completeness. *Jisuanji Jicheng Zhizao Xitong/Computer Integrated Manufacturing Systems* 17(6), 1308–1316 (2011)
5. Object Management Group: DTC/2008-01-17 Business Process Modeling Notation, <http://www.omg.org/spec/BPMN/1.1/PDF>
6. Workflow Management Coalition: WFMC-TC-1025 Process Definition Interface – XML Process Definition Language, <http://www.wfmc.org/xpdl.html>
7. Microsoft: Windows Workflow Foundation, <http://www.microsoft.com/visualstudio/en-us/products/2010-editions>
8. Sendall, S., Kozaczynski, W.: Model transformation: The heart and soul of model-driven software development. *IEEE Software* 20(5), 42–45 (2003)
9. Jian, D., Zhi, C., Jiazhi, Z.: Formal verification of business process models using Petri nets. *Jisuanji Jicheng Zhizao Xitong/Computer Integrated Manufacturing Systems* 17(5), 1110–1119 (2011)
10. TIBCO: TIBCO Business Studio, http://developer.tibco.com/business_studio/
11. Bruce, B.: *Pro WF: Windows Workflow in .NET 4*. Apress, New York (2010)
12. Jordan, D., Evdemon, J.: Web Services Business Process Execution Language Version 2.0, <http://docs.oasis-open.org/wsbpel/2.0/OS/wsbpel-v2.0-OS.pdf>
13. Weqqing, L., Jian, W.: Analysis of process control pattern of XPDL 2.0. *Jisuanji Jicheng Zhizao Xitong/Computer Integrated Manufacturing Systems* 13(9), 1839–1846 (2007)
14. Zapletal, M., Aalst, W.M.P., Russell, N., et al.: An Analysis of Windows Workflow's Control-Flow Expressiveness. In: Eshuis, R., Grefen, P., Papadopoulos, G. (eds.) *Proceedings of the 7th IEEE European Conference on Web Services*, pp. 200–209. IEEE Computer Society, Los Alamitos (2009)
15. Lin, M., Tao, J., Jianmin, W.: Transformation Technology from XPDL Model to WWF Model. *Jisuanji Yanjiu yu Fazhan/Journal of Computer Research and Development* 46(suppl.), 165–171 (2009)
16. Gray, J., Reuter, A.: *Transaction Processing: Concepts and Techniques*. Morgan Kaufmann, San Francisco (1993)

Towards Efficient Replication of Documents in Chord: Case (r,s) Erasure Codes^{*}

Rafał Kapelko

Institute of Mathematics and Computer Science,
Faculty of Fundamental Problems of Technology,
Wrocław University of Technology, Poland
rafal.kapelko@pwr.wroc.pl

Abstract. In this article we investigate the replication mechanism, called Global Policy with (r, s) erasure codes for Chord peer-to-peer system. Each document is divided into s fragments and then is encoded into $s + r$ fragments that can tolerate r failures. The independent hash functions: H_1, \dots, H_{r+s} are used for placing fragments of documents into the Chord. In our paper we show analytical formulas describing basic probabilistic properties of the proposed scheme.

Keywords: peer-to-peer network, Chord, replication of documents, erasure codes.

1 Introduction

Peer-to-Peer systems (see e.g. [16], [14], [13], [5]) and others have been popular and powerful networks for successful sharing of certain resources. The replication mechanism is used to increase data availability (see e.g. [4], [7], [9]). Instead of replicating the whole file we propose to replicate a portion of the file.

Erasure Codes (r, s) in replication in P2P systems were first studied in [18]. We fix two parameters s and r . Each data item is divided into s fragments and then is encoded into $r+s$ fragments. We call $r/(r+s)$ the *rate* of encoding. The *storage cost* is increased by a factor $(r + s)/s$. The key property of erasure codes is that the original data item can be reconstructed from any s fragments.

In this paper we focus on presenting the analytical formulas describing basic probabilistic properties of the (r, s) erasure codes for replication in Chord P2P system. We consider Global Policy replication mechanism (see [6], [3], [4]) to store $s + r$ fragments of the file via a family H_1, \dots, H_{r+s} of independent hash functions.

There are several previous works related to the subject [12], [8], [10]. In [12] authors present the availability analysis of whole-file replication versus erasure codes replication. The effectiveness of erasure codes replication with different parameters is discussed. In [8] an analytical optimization theory for benchmarking the performance of replication algorithms that employ erasure codes is developed. The performance of erasure codes for different models in P2P storage systems was evaluated in [10]. The

^{*} Supported by grant nr 2011/S10026 of the Institute of Mathematics and Computers Science of the Wrocław University of Technology.

comparative analysis is based on that node session length follows exponential distribution, Pareto distribution and Weibull distribution.

The remainder of this article is organised as follows. In Section 2 we describe Chord P2P system with the basic facts and notation. Section 3 presents the structure called Chord with (r, s) -replicas of documents. The analytical formulas describing basic probabilistic properties of the replicas scheme are presented in Section 4. Then, in Section 5 we discuss the main results. Finally, Section 6 summaries and concludes our work.

2 Basic Facts and Notation

The classical Chord protocol defined in [16] and developed in [11], and many other papers from the formal point of view may be described as a structure

$$\text{Chord} = (\{0, 1\}^{160}, H, H_1),$$

where H is a hash function assigning position to each node and H_1 is a hash function assigning position of descriptors of documents. The space $\{0, 1\}^{160}$ is identified with the set $\{0, 1, \dots, 2^{160} - 1\}$ considered as the circular space with the ordering $0 < 1 < \dots < 2^{160} - 1 < 0 < \dots$. Each new node X obtains a position $H(Id)$ (where Id is an identifier of the node) in the space $\{0, 1\}^{160}$ and is responsible for the interval starting at point $H(Id)$ and ending at the next point from the set $\{H(Id') : Id' \neq Id\}$. This node is called the successor of the node X . Each document with a descriptor doc is placed at point $H_1(doc)$ in the space $\{0, 1\}^{160}$ and the information about this document is stored by the node which is responsible for the interval into which $H_1(doc)$ falls.

We shall identify the space $\{0, 1\}^{160}$, after a proper scaling, with the unit interval $[0, 1)$ and we shall interpret positions of nodes of a Chord as elements of $[0, 1)$ (see e.g. [4], [15]). Moreover, we may assume that one node is at point 0. The random sets corresponding to nodes of Chord are generated in the following way: we generate independently n random points X_1, \dots, X_n from $[0, 1)$ using the uniform distribution on $[0, 1)$ and then we sort them in increasing order and obtain a sequence $X_{1:n} \leq \dots \leq X_{n:n}$. This construction will be used as a formal model of the Chord protocol with $n + 1$ nodes. We call the segment $[X_{i:n}, X_{i+1:n})$ the interval controlled by the node i .

We will use several times the Eulerian function

$$B(a, b) = \int_0^1 x^{a-1}(1-x)^{b-1} dx,$$

which are defined for all complex numbers a, b such as $\Re(a) > 0$ and $\Re(b) > 0$.

We will use the following basic identity $B(a, b) = \Gamma(a)\Gamma(b)/\Gamma(a+b)$, where $\Gamma(z)$ denotes the standard generalization of the factorial function. The following identities hold: $n! = \Gamma(n+1)$, $\Gamma(z+1) = z\Gamma(z)$.

Let X_1, \dots, X_n be independent random variables with the uniform density on $[0, 1)$. The order statistics $X_{1:n}, \dots, X_{n:n}$ are the random variables obtained from X_1, \dots, X_n by sorting each of their realizations in the increasing order. The probabilistic density $f_{k:n}(x)$ of the variable $X_{k:n}$ equals

$$f_{k:n}(x) = \frac{1}{B(k, n-k+1)} x^{k-1}(1-x)^{n-k} \tag{1}$$

(see e.g. [11]). Let us recall that these kinds of probabilistic distributions are called Beta distributions.

Let X be a random variable. We denote its expected value and variance by $\mathbf{E}[X]$ and $\mathbf{var}[X]$ respectively.

3 Chord with (r, s) –Replicas of Documents

The Chord with (r, s) -replicas of documents is the structure

$$r_{(r,s)}\text{-Chord} = (\Omega, H, \{H_1, \dots, H_{r+s}\}) ,$$

where H, H_1, \dots, H_{r+s} are independent hash functions. H is used for putting nodes into the system. Each document is divided into $r + s$ fragments that can tolerate r failures. The independent hash functions: H_1, \dots, H_{r+s} , are used for placing fragments of documents into the system. If we put d documents into the structure $r_{(r,s)}\text{-Chord}$, then the total number of „information items” in the system is $d \cdot (r + s)$.

4 Replication of Documents in Chord

In this section we investigate resistance to the loss of documents in Chord with (r, s) –replicas of documents. We present the analytical formulas for expected value and variance of random variable describing the number of nodes which must be removed in order to lose some information from the structure.

Theorem 1. *Let $n+1$ denote the number of nodes and let d be the number of documents put into this structure $r_{(r,s)}\text{-Chord}$. Let $K_{n:d}^{(r,s)}$ denote the number of nodes which must be removed in order to lose some information from the system. Then*

$$\mathbf{E} \left[K_{n:d}^{(r,s)} \right] = 1 + n \int_0^1 (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l} \right)^d dx.$$

Proof. Assume that the number of nodes in a $r_{(r,s)}\text{-Chord}$ is $n + 1$. Let X_1, \dots, X_n denote the family of independent random numbers from the interval $[0, 1]$ with uniform distributions, which corresponds to nodes in Chord. Let d denotes the number of documents put into the system. Let Y_i^j denotes the position of i -th fragments of j -th document. We treat $(Y_i^j)_{i=1\dots s+r, j=1\dots d}$ as a family of independent random variables with values in the interval $[0, 1]$ with the uniform distribution.

Let A be a subset of nodes. Let the total area controlled by nodes from the set A be L_A . Let S_A^j denote the event that after removing the set A of nodes we do not lose any information stored in the j -th document. Let S_A denote the event that after removing the set A of nodes we do not lose any information stored in the system.

There are $r + 1$ reasons why removing the set A of nodes we do not lose any information stored in the j -th document. Let $l \in \{0, 1, \dots, r\}$ The l -th reason is that exactly l random variables from the set $\{Y_1^j, Y_2^j, \dots, Y_{s+r}^j\}$ belong to the set L_A . Therefore,

$$\Pr[S_A^j | L_A = x] = (1-x)^s \sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l}.$$

Hence

$$\Pr[S_A|L_A = x] = \Pr\left[\bigwedge_{j=1}^d S_A^j|L_A = x\right] = \prod_{j=1}^d \Pr[S_A^j|L_A = x] = (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l}\right)^d. \tag{2}$$

Let $K_{n:d}^{(r,s)}$ denotes the number of nodes which must be removed in order to lose some information from the system. From (1) and (2) we obtain

$$\Pr[K_{n:d}^{(r,s)} > k] = \frac{1}{B(k, n-k+1)} \int_0^1 (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l}\right)^d x^{k-1} (1-x)^{n-k} dx.$$

Notice that, $\Pr[K_{n:d}^{(r,s)} > 0] = 1$ and $\Pr[K_{n:d}^{(r,s)} > n+1] = 0$. Therefore,

$$\begin{aligned} \mathbf{E} \left[K_{n:d}^{(r,s)} \right] &= \sum_{k \geq 0} \Pr[K_{n:d}^{(r,s)} > k] = 1 + \sum_{k=1}^n \Pr[K_{n:d}^{(r,s)} > k] \\ &= 1 + \sum_{k=1}^n k \binom{n}{k} \int_0^1 (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l}\right)^d x^{k-1} (1-x)^{n-k} dx = \\ &= 1 + \int_0^1 (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l}\right)^d \left(\sum_{k=1}^n k \binom{n}{k} x^{k-1} (1-x)^{n-k}\right) dx. \end{aligned}$$

Finally, from the identity $\sum_{k=1}^n k \binom{n}{k} x^{k-1} (1-x)^{n-k} = n$ we get

$$\mathbf{E} \left[K_{n:d}^{(r,s)} \right] = 1 + n \int_0^1 (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l}\right)^d dx,$$

hence the proof is finished. □

Now we calculate the second factorial moment of random variable $K_{n:d}^{(r,s)}$.

Theorem 2

$$\begin{aligned} \mathbf{E} \left[K_{n:d}^{(r,s)} \left(K_{n:d}^{(r,s)} - 1 \right) \right] &= 2n \int_0^1 (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l}\right)^d dx \\ &+ 2n(n-1) \int_0^1 (1-x)^{sd} x \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l}\right)^d dx. \end{aligned}$$

Proof. Firstly, we calculate

$$\begin{aligned} & \sum_{k=1}^n \Pr[K_{n:d}^{(r,s)} > k](k-1) = \\ & \sum_{k=1}^n k(k-1) \binom{n}{k} \int_0^1 (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l} \right)^d x^{k-1} (1-x)^{n-k} dx = \\ & \int_0^1 (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l} \right)^d \left(\sum_{k=1}^n k(k-1) \binom{n}{k} x^{k-1} (1-x)^{n-k} \right) dx . \end{aligned}$$

From the identity $\sum_{k=1}^n k(k-1) \binom{n}{k} x^{k-1} (1-x)^{n-k} = xn(n-1)$ we get

$$\begin{aligned} & \sum_{k=1}^n \Pr[K_{n:d}^{(r,s)} > k](k-1) = \\ & n(n-1) \int_0^1 (1-x)^{sd} x \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l} \right)^d dx . \end{aligned}$$

Let us recall that for a discrete nonnegative random variable X we have

$\mathbf{E}[X(X-1)] = 2\sum_{k \geq 0} \Pr[X > k]k$ (see e.g. [17]). Therefore, from Theorem 1 we get

$$\begin{aligned} \mathbf{E} \left[K_{n:d}^{(r,s)} \left(K_{n:d}^{(r,s)} - 1 \right) \right] &= \mathbf{E} \left[K_{n:d}^{(r,s)} \right] - 2 + 2 \sum_{k=1}^n \Pr[K_{n:d}^{(r,s)} > k](k-1) = \\ & 2n \int_0^1 (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l} \right)^d dx \\ & + 2n(n-1) \int_0^1 (1-x)^{sd} x \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l} \right)^d dx . \end{aligned}$$

hence the proof is finished. □

The following results about variance of the random variable $K_{n:d}^{(r,s)}$ follows directly from Theorem 1 and Theorem 2

Theorem 3

$$\begin{aligned} \mathbf{var} \left[K_{n:d}^{(r,s)} \right] &= n \int_0^1 (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l} \right)^d dx \\ & - n^2 \left(\int_0^1 (1-x)^{sd} \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l} \right)^d dx \right)^2 \\ & + 2n(n-1) \int_0^1 (1-x)^{sd} x \left(\sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l} \right)^d dx . \end{aligned}$$

5 Discussion

Let us assume that $s = 1$ and consider $(r, 1)$ erasure codes. From Theorem 1 the identity $\sum_{l=0}^r \binom{r+1}{l} x^l (1-x)^{r-l} = \frac{1-x^{r+1}}{1-x}$ and the definite integral expressible in terms of the Gamma function (see e.g. [2] page 10)

$$\int_0^1 (1-x^{r+1})^d dx = \frac{\Gamma(1+d)\Gamma(1+\frac{1}{r+1})}{\Gamma(1+d+\frac{1}{r+1})} \tag{3}$$

we get

$$\mathbf{E} \left[K_{n;d}^{(r,1)} \right] = 1 + n \frac{\Gamma(1+d)\Gamma(1+\frac{1}{r+1})}{\Gamma(1+d+\frac{1}{r+1})}. \tag{4}$$

The formula (4) was also obtained in [4] with more complex techniques. However, due to the limited size of the paper the proof of (4) was not presented.

Now, for fixed r we compare erasure codes $(r, 1), (r, 2), \dots, (r, s)$. Notice that, for all $x \in [0, 1]$ we have

$$\begin{aligned} (1-x)^s \sum_{l=0}^r \binom{r+s}{l} x^l (1-x)^{r-l} - (1-x)^{s+1} \sum_{l=0}^r \binom{r+s+1}{l} x^l (1-x)^{r-l} = \\ x^{1+r} (1-x)^s \binom{s+r}{s} \geq 0 . \end{aligned}$$

From this we deduce the following result:

Corollary 1

$$\mathbf{E} \left[K_{n;d}^{(r,s)} \right] < \mathbf{E} \left[K_{n;d}^{(r,1)} \right].$$

Therefore, for fixed r and d erasure code $(r,1)$ with the biggest storage cost has the biggest resistance to the loss of documents in the set $\{(r, 1), (r, 2), \dots, (r, s)\}$ of erasure codes.

From (3) and (4) we also deduce that for fixed d the expected value $\mathbf{E} \left[K_{n;d}^{(r,1)} \right]$ as the function of parameter r is increasing.

6 Conclusion

In this paper, we study replication scheme, called Global Policy with (r, s) erasure codes for Chord P2P system. We obtain the analytical formulas for expected value and variance of random variable $K_{n;d}^{(r,s)}$ describing the number of nodes which must be removed in order to lose some information from the structure Chord with (r, s) -replicas of documents. By comparing erasure codes $(r, 1), (r, 2), \dots, (r, s)$ we prove that $(r, 1)$ has the biggest resistance to the loss of documents.

References

1. Arnold, B., Balakrishnan, N., Nagaraja, H.: *A First Course in Order Statistics*. John Wiley & Sons, New York (1992)
2. Bateman, H.: *Higher Transcendental Functions*. McGraw-Hill Book Company, Inc., Malabar (1953)
3. Caron, S., Giroire, F., Mazaauric, D., Monteiro, J., Pérennes, S.: Data Life Time for Different Placement Policies in P2P Storage Systems. In: Hameurlain, A., Morvan, F., Tjoa, A.M. (eds.) *Globe 2010*. LNCS, vol. 6265, pp. 75–88. Springer, Heidelberg (2010)
4. Cichoń, J., Jasiński, A., Kapelko, R., Zawada, M.: How to Improve the Reliability of Chord? In: Meersman, R., Tari, Z., Herrero, P. (eds.) *OTM 2008 Workshops*. LNCS, vol. 5333, pp. 904–913. Springer, Heidelberg (2008)
5. Rowstron, A., Druschel, P.: Pastry: Scalable, Decentralized Object Location, and Routing for Large-Scale Peer-to-Peer Systems. In: Guerraoui, R. (ed.) *Middleware 2001*. LNCS, vol. 2218, pp. 329–350. Springer, Heidelberg (2001)
6. Giroire, F., Monteiro, J., Pérennes, S.: P2P Storage Systems: How Much Locality Can They Tolerate? In: *34th Annual IEEE Conference on Local Computer Networks (LCN)*, Zurich, Switzerland, pp. 320–323 (2009)
7. Park, G., Kim, S., Cho, Y., Kook, J., Hong, J.: Chordet: An Efficient and Transparent Replication for Improving Availability of Peer-to-Peer Networked Systems. In: *2010 ACM Symposium on Applied Computing (SAC)*, Sierre, Switzerland, pp. 221–225 (2010)
8. Kangasharju, J., Ross, K.W., Turner, D.A.: Optimizing File Availability in Peer-to-peer Content Distribution. In: *26th IEEE International Conference on Computer Communications (INFOCOM)*, Anchorage, Alaska, USA, pp. 1973–1981 (2007)
9. Kapelko, R.: Improving Data Availability in Chord p2p System. In: Liu, B., Chai, C. (eds.) *ICICA 2011*. LNCS, vol. 7030, pp. 208–215. Springer, Heidelberg (2011)
10. Li, J., Xu, G., Zhang, H.: Performance Comparison of Erasure Codes for Different Churn Models in P2P Storage Systems. In: Huang, D., Zhang, X., Reyes Garcia, C.A., Zhang, L. (eds.) *ICIC 2010*. LNCS, vol. 6216, pp. 410–417. Springer, Heidelberg (2010)
11. Liben-Nowell, D., Balakrishnan, H., Karger, D.: Analysis of the Evolution of Peer-to-Peer Systems. In: *ACM Conference on Principles of Distributed Computing*, Monterey, California, USA, pp. 233–242 (2002)
12. Lin, W.K., Chiu, D.M., Lee, Y.B.: Erasure Code Replication Revisited. In: *2004 IEEE Fourth International Conference on Peer-to-Peer Computing, Peer-to-Peer Computing (P2P)*, Zurich, Switzerland, pp. 90–97 (2004)
13. Maymounkov, P., Mazières, D.: Kademlia: A Peer-to-Peer Information System Based on the XOR Metric. In: Druschel, P., Kaashoek, F., Rowstron, A. (eds.) *IPTPS 2002*. LNCS, vol. 2429, pp. 53–65. Springer, Heidelberg (2002)
14. Ratnasamy, S., Francis, P., Handley, M., Karp, R., Shenker, S.: A Scalable Content-addressable Network. In: *SIGCOMM 2001*, San Diego, California, USA, pp. 161–172 (2001)
15. Rumín, R.C., Uruña, M., Banchs, A.: Routing Fairness in Chord: Analysis and Enhancement. In: *28th IEEE International Conference on Computer Communications (INFOCOM)*, Rio de Janeiro, Brazil, pp. 1449–1457 (2009)
16. Stoica, I., Morris, R., Karger, D., Kaashoek, M.F., Balakrishnan, H.: Chord: A Scalable Peer-to-Peer Lookup Service for Internet Applications. In: *SIGCOMM 2001*, San Diego, California, USA, pp. 149–160 (2001)
17. Szpankowski, W.: *Average Case Analysis of Algorithms on Sequences*. A Wiley-Interscience Publication, New York (2001)
18. Weatherspoon, H., Kubiatowicz, J.D.: Erasure Coding Vs. Replication: A Quantitative Comparison. In: Druschel, P., Kaashoek, F., Rowstron, A. (eds.) *IPTPS 2002*. LNCS, vol. 2429, pp. 328–338. Springer, Heidelberg (2002)

Path Planning for Crawler Crane Using RRT*

Yuanshan Lin¹, Di Wu¹, Xin Wang², Xiukun Wang¹, and Shunde Gao²

¹School of Computer Science, Dalian University of Technology, Dalian, Liaoning, P.R. China

²School of Mechanical Engineering, Dalian University of Technology, Dalian, Liaoning, P.R. China

Linyuanshan2008@126.com, {wudi,Jsjwxk}@dlut.edu.cn,
{Wangxbd21,gaoshunde}@163.com

Abstract. Lift path planning is a key subtask in preparing lift plan. Many researchers have concerned the topic of automated path planning for crane lifts. Previous studies provided some exact approaches or optimization methods to solve the problem. However, these approaches failed when they handle the problem with high DOFs. Due to this, this paper presents an approach for automated path planning of crawler crane using RRT*. The experimental case indicates that the proposed algorithm can handle lift path planning problem with high DOFs and effectively find a safe path in complex lifting environment.

Keywords: Path planning, Crane, RRT, High DOFs, lift plan.

1 Introduction

The lift path planning of a crane is an important subtask within the overall heavy lift planning process. In a heavy lift, a detailed lift plan is required to be developed prior to actual lifting. As the core of lift plan, path planning for a crane lift is generating an action sequence of crane. And during the series of actions are performed, the crane must be collision-free and not be overload. The quality of lift path determines whether the lift is successful. But planning the lift path of a crane is a challenging task since it requires considering the degrees of freedom of the crane, the lifting capacity, and potential obstacles on the site. Particularly, on retrofit and expansion construction projects the site will be congested and finding a suitable lift path can be a difficult and tedious task. It is very necessary to develop a path-planning tool to aid the planner to rapidly generate lift paths from the picking pose to the placing pose. Further more, an alternate path can be regenerated using the tool when some last-minute changes to the constraints happen.

Actually, the lift path planning has drawn much concern of some researchers. Koshy Varghese etc. used search methods, A^* , hill climbing and GA, one after another to accomplish path planning of single/two crane lifting, besides analyzed and compared these algorithms[1, 2]. Later, they made a smooth path with 4th order trigonometric splines-Based trajectory planning [3]. Mohamed Al-Hussein provided a method for single crane path planning in 2010, which segmented the lifting process to make control easier [4]. Erashima Kazuhiko etc. plan and control the motion of single

bridge crane to achieve the safety and high efficiency of fast lifting [5]. In 2011, Wangxin etc. proposed path planning method for single crane based on ant colony algorithm [6]. These researchers have made many positive attempts and provided their solutions to the path planning problem of crane lift. However, they mainly focus on the crane lift without considering mobility of crane and these approaches become very inefficient in practice when handling the problem with high DOFs. Path planning considering mobility of crane is a computationally challenging problem because of the high DOFs of the system. No attempt has been made to automate the path planning of crawler crane lift.

Rapidly-exploring Random Trees (RRTs) is one of the most popular path planning approaches, which is specifically designed to handle nonholonomic constraints (including dynamics) and high degrees of freedom [7]. Its idea is to incrementally grow a space-filling tree by sampling the space at random and connecting the nearest point in the tree to the new random sample. In the meanwhile, the RRT has been successfully used in numerous robotic systems, including DARPA Urban Challenge vehicles, military UAVs, humanoid robots, soccer-playing robots and so on [8-12]. RRT* is a variant of RRT and is proved that it is globally asymptotically optimal for path planning problems without differential constraints [14]. So, this paper tries to apply RRT* to solve path planning problem considering the mobility of crawler crane. Finally, we verify the feasibility and validity of this algorithm by a simulation case of construction lifting.

2 Problem Formulation

The problem of path planning for crawler crane is: given picking configuration and placing configuration of crawler crane, finding a safe action sequence in construction site with obstacles and making sure that the crawler crane is under its capacity and is collision-free along with the action sequence. The path found must meet two requirements: 1) the crawler crane must be not overload; 2) and there is no collision among the crawler crane, the object lifted and obstacles. Furthermore, we expect that the trajectory of object lifted is shortest. Therefore, the problem can be formulated as following:

$$P = (S, Obs, q_{init}, q_{goal}, U, f_{col}) \quad (1)$$

Where,

S : State space of crawler crane, a subset of the set R^n ;

Obs : Representation of Obstacles in the construction site;

q_{init} : Picking state of crawler crane, the starting point of lift path;

q_{goal} : Placing state of crawler crane, the end point of lift path;

U : Input set of crawler crane, containing releasing hook, rising hook, luffing up boom, luffing down boom, slewing left, slewing right, traveling forward, traveling backward, turning left, turning right and combinations of above motions.

f_{col} : Function of collision detection, no collision among the crawler crane, the object lifted and obstacles in construction site.

3 Based-on RRT* Path Planning for Crawler Crane

3.1 RRT* Algorithm

RRT*[13] is a novel variant of the RRT algorithm, which is proven that it is a sampling-based algorithm with the asymptotic optimality property. The pseudo-code of this approach is shown as Algorithm 1 and Algorithm 2. RRT* introduces the concept of near neighbors, the nodes within a certain radius of a node. Once a new node is added to the tree, the RRT* checks the near neighbors whether can be rewired through the new node with lower cost.

Algorithm 1:

```

program GrowRRT* (  $q_{init}$  )
   $V \leftarrow \{q_{init}\}$ ;  $E \leftarrow \Phi$ ;  $i \leftarrow 0$ ;
  While  $i < N$  do
     $G \leftarrow (V, E)$ ;
     $q_{rand} \leftarrow \text{Sample}(i)$ ;  $i \leftarrow i + 1$ ;

     $(V, E) \leftarrow \text{ExtendRRT}^*(G, q_{rand})$ 

```

Algorithm 2:

```

program ExtendRRT* (  $G, q_{rand}$  )
   $V' \leftarrow V$ ;  $E' \leftarrow E$ ;
   $q_{nearest} \leftarrow \text{Nearest}(G, q_{rand})$ ;
   $q_{new} \leftarrow \text{Steer}(q_{nearest}, q_{rand})$ ;
  if (CollisionFree( $q_{nearest}, q_{new}$ ))
     $V' \leftarrow V' \cup \{q_{new}\}$ ;
     $q_{min} \leftarrow q_{nearest}$ ;
     $Q_{near} \leftarrow \text{Near}(G, q_{new})$ ;
    For all  $q_{near} \in Q_{near}$  do
      if (CollisionFree( $q_{near}, q_{new}$ ))
         $c \cdot \text{Cost}(q_{near}) + C(\text{Line}(q_{near}, q_{new}))$ ;
        if ( $c < \text{Cost}(q_{new})$ )
           $q_{min} \leftarrow q_{near}$ ;
     $E' \leftarrow E' \cup \{(q_{min}, q_{new})\}$ ;
    For all  $q_{near} \in Q_{near} \setminus \{q_{min}\}$  do
      if (CollisionFree( $q_{new}, q_{near}$ )
        and  $\text{Cost}(q_{near}) > \text{Cost}(q_{new}) + C(\text{Line}(q_{new}, q_{near}))$ )
         $q_{parent} \leftarrow \text{Parent}(q_{near})$ ;
         $E' \leftarrow E' \setminus \{q_{parent}, q_{near}\}$ ;
         $E' \leftarrow E' \cup \{q_{parent}, q_{near}\}$ ;
  Return  $G' = (V', E')$ ;

```

3.2 State-Space Representation of Crawler Crane

Crawler crane, a kind of mobile crane widely used in the construction sites, can travel with a load on the hook. It consists of bottom chassis, superstructure, boom system and lift system. The bottom chassis contains crawlers and chassis frame; Superstructure, as a supporting body of upper-structure, contains power devices and mechanisms devices; Boom system is composed of main boom, mast and its attachment; and lift system includes lifting ropes and the hook. Slewing ring connects bottom chassis and superstructure, which allows superstructure to slew in 360° . The boom system is mounted to the top of the superstructure, some distance away from the slewing centerline. The lift system is hinged on the top of boom. So, the crawler crane has five motions: traveling, turning, slewing, luffing, hoisting. In some situations, manpower is required to rotate the lifted object to pass narrow spaces. We add a motion “Hook Rotating” for the crawler crane to imitate such action. The components and motions of crawler crane are shown in Fig. 1.

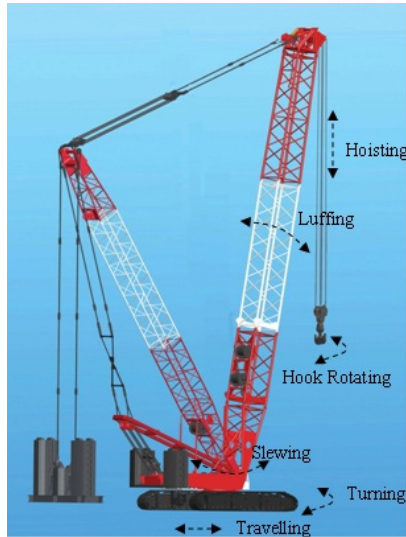


Fig. 1. Components and motions of a crawler crane

It is worth mentioning that the deformation of components of crawler crane and off-lead of lifting rope are neglected. Thus, the state of the crawler crane is described with a 7 dimensions tuple $(x, z, \alpha, \beta, \gamma, h, \omega)$. Where, (x, z) is the Cartesian coordinates of the crane, α is the heading angle, β is the slewing angle, γ is the boom angle, h is the length of hoisting rope, and ω is the rotating angle of hook. All states of the crawler crane compose of state-space of the crawler crane.

3.3 Length of Path

We define that the distance metric between two states is the length of object-lifted trajectory generated by crane’s input applied and path length is the sum of the Euclidean distances between consecutive poses in the solution trajectory. Let $q_i = (x_i, z_i, \alpha_i, \beta_i, \gamma_i, h_i, \omega_i)$ and $q_{i+1} = (x_{i+1}, z_{i+1}, \alpha_{i+1}, \beta_{i+1}, \gamma_{i+1}, h_{i+1}, \omega_{i+1})$ are two elements in state-space of crawler crane, then:

$$d(q_i, q_{i+1}) = |x_{i+1} - x_i| + |z_{i+1} - z_i| + |r \times (\alpha_{i+1} - \alpha_i)| + |l_z \times (\beta_{i+1} - \beta_i)| + |h_{i+1} - h_i| + |l_s \times (\omega_{i+1} - \omega_i)| \tag{2}$$

Where, $d(q_i, q_{i+1})$ is the distance between q_i and q_{i+1} , r is the working radius, l_z is the boom length, l_s is the length of lifted object. Thus, the whole length of the path L can be expressed as: $L = \sum_{i=0}^{n-1} d(q_i, q_{i+1})$. Where, n is the number of state in the path.

4 Case Study

To verify the effectiveness of the proposed approach for crawler crane lift, an actual lift is tested in this section. This case is a petrochemical refinery lift project. The electric desalting tank lifted is 27m long, 3.8m diameter, and 37 tons weight. At the beginning lifting, the object lifted (electric desalting tank) is placed at (-170.0, 0.0, -71.0), the angle between the centerline of object lifted and the X axis is 0 degree, and the installation position is (-220.7, 15.4, -11.8) with the angle of the object centerline and X axis $\frac{\pi}{2}$.

According to the length, diameter, weight parameters of the lifted object and the work environment, a crawler crane is selected to accomplish such lift task. Its parameters are shown in Table 1.

Table 1. Crane parameters

Crane Parameters	Value
Crane Type	LR1400-2
Configuration	S
Boom Length(m)	49
Counter Weight (t)	135
Center Weight(t)	43
Hook Weight(t)	8

In this case, our task is finding a safe lift path from the picking state (-170.0 57.0 - 3.14 1.57 1.34 45 0.0) to the placing state (-204.5 56.7 -3.1415 -0.8987 1.023 43.3

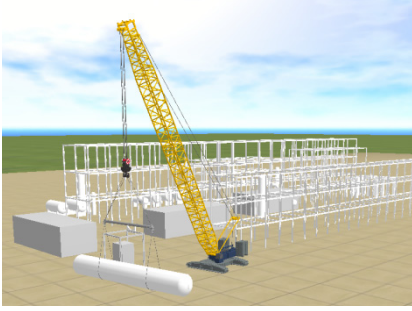


Fig. 2. Picking state

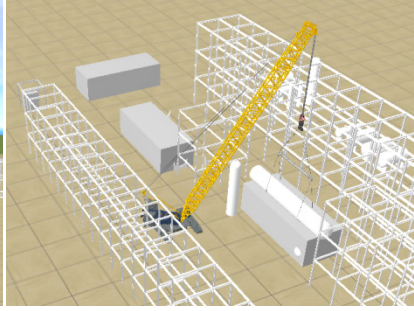


Fig. 3. Placing state

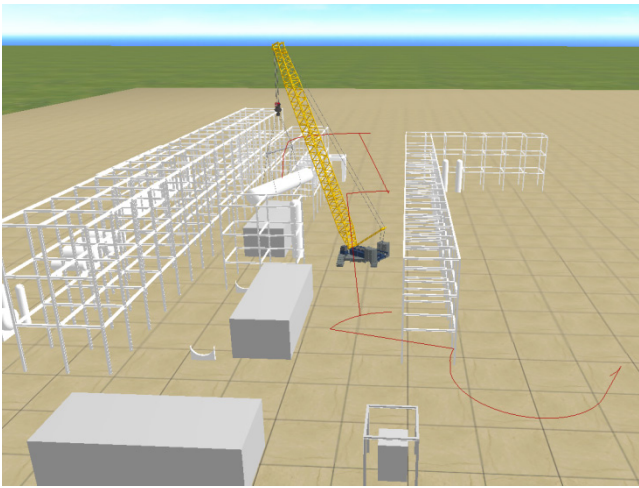


Fig. 3. Path generated by the proposed approach (the red curve is the trajectory of object-lifted)

0.8987) of crawler crane in construction site with obstacles. The path which is generated using the proposed approach is shown as Fig. 4. Fig. 5 illustrates that the crawler crane operates along with lift path.

The algorithm is realized with C++, whose collision detection library is PQP. And we make Kd-trees as secondary data structure to raise efficiency of selecting nearest neighbor(s). The test is performed on Thinkpad T60 2.0GHz computer with 1GB of memory. We run the approach 100 trials and observe that it takes about 30.68s for attaining the first path. The trajectory of lifted object is not very optimal and smooth, but the proposed approach can generate a sub-optimal path. It shows that the proposed approach is capable of effectively handle the path planning problem of a crawler crane.

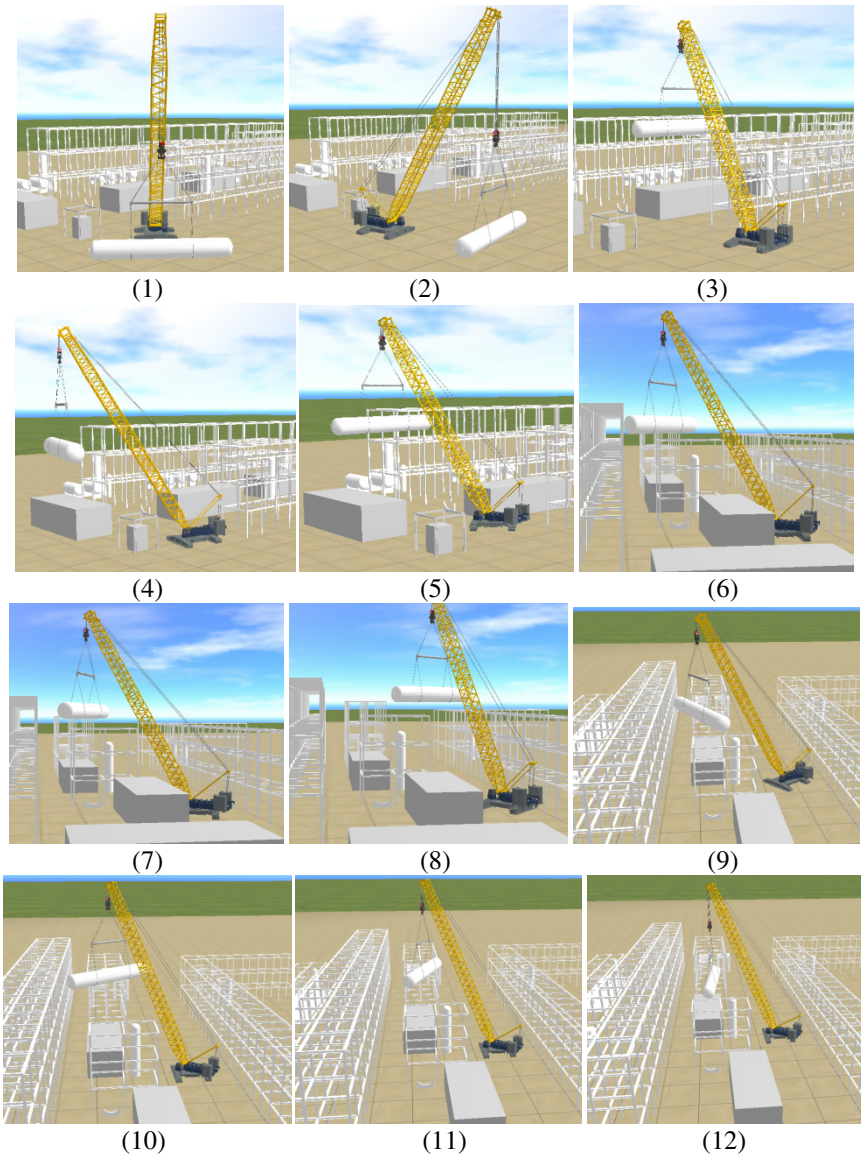


Fig. 4. Simulation snapshots of crawler crane operates along with lift path

5 Conclusion

This paper presents a path planning approach for crawler crane lifts based-on RRT*. The mobility of crawler cranes is considered in the approach. Firstly, the path planning problem for crawler crane lifts is formulated according to the characteristics of RRT. And then the state-space and path length are defined. Finally, this paper

verifies the effectiveness of the approach using an actual case. The case indicates that the proposed algorithm can effectively find a safe path in complex lifting environment.

References

1. Sivakumar, P.L., Varghese, K., Babu, N.R.: Automated path planning of cooperative crane lifts using heuristic search. *Journal of Computing in Civil Engineering* 17(3), 197–207 (2003)
2. Deen Ali, M.S.A., Babu, N.R., Varghese, K.: Collision free path planning of cooperative crane manipulators using genetic algorithm. *Journal of Computing in Civil Engineering* 19(2), 182–193 (2005)
3. Bhaskar, S.V., Babu, N.R., Varghese, K.: Spline Based Trajectory Planning for Cooperative Crane Lifts. In: *Proceedings of the 23rd ISARC, Tokyo, Japan*, pp. 418–423 (2006)
4. Olearczyk, J., et al.: Spatial trajectory analysis for cranes operations on construction sites. In: *Construction Research Congress 2010: Innovation for Reshaping Construction Practice - Proceedings of the 2010 Construction Research Congress, Banff, AB, Canada*, pp. 359–368 (2010)
5. Kaneshige, A., Miyoshi, T., Terashima, K.: The development of an autonomous mobile overhead crane system for the liquid tank transfer. In: *IEEE/ASME International Conference on Advanced Intelligent Mechatronics, AIM, Singapore*, pp. 630–635 (2009)
6. Wang, X., et al.: Collision-Free Path Planning for Mobile Cranes Based on Ant Colony Algorithm. *Key Engineering Materials* 467, 1108–1115 (2011)
7. Lavalle, S.M.: *Rapidly-Exploring Random Trees: A New Tool for Path Planning* (1998)
8. Kuffner Jr., J.J., et al.: Dynamically-stable motion planning for humanoid robots. *Autonomous Robots* 12(1), 105–118 (2002)
9. Kuwata, Y., et al.: Real-Time Motion Planning With Applications to Autonomous Urban Driving. *IEEE Transactions on Control Systems Technology* 17(5), 1105–1118 (2009)
10. Shkolnik, A., et al.: Bounding on rough terrain with the LittleDog robot. *The International Journal of Robotics Research* 30(7), 846–894 (2011)
11. Teller, S., et al.: A voice-commandable robotic forklift working alongside humans in minimally-prepared outdoor environments. In: *2010 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 526–533 (2010)
12. Berenson, D., Kuffner, J., Choset, H.: An optimization approach to planning for mobile manipulation. In: *IEEE International Conference on Robotics and Automation, ICRA 2008*, pp. 1187–1192 (2008)
13. Karaman, S., Frazzoli, E.: *Incremental Sampling-based Algorithms for Optimal Motion Planning* (2010)

Study on Data Preprocessing for Daylight Climate Data

Ping Guo¹, Shuai-Shuai Chen¹, and Ying He²

¹ School of Computer Science, Chongqing University, Chongqing, 400044, China

² College of Architecture and Urban Planning, Chongqing University, Chongqing, China
guoping@cqu.edu.cn

Abstract. It is well known that the real-world data tend to exist many data quality problems such as incompleteness and noisy data. Data preprocessing technology can improve data quality effectively and provide more reliable data for the next step. A data preprocessing approach for daylight climate data is presented in this paper according to the characteristics of this data. Then this approach is applied to the real-world data and the experimental results show that the approach can enhance the data quality effectively. Besides, the integration of the domain knowledge into data preprocessing is emphasized in this paper in order to make data preprocessing more effective and more targeted.

Keywords: Data preprocessing, Data quality, Knowledge base, Intelligent information processing, Daylight climate data.

1 Introduction

With the development of the technology of computer software and hardware and the maturity of the data acquisition method, data which we collect in every field are increasing in exponentially speed. The equipment failure, the environmental noise interference and error caused by manual input lead to many data quality problems in data acquisition process. Data preprocessing can settle these problems effectively and output clean, accurate and concise data.

With the development of science and technology and the improvement of living standards for the people, demands for energy are rising. But the energy such as oil and coal in the earth is limited. For saving the limited energy, people become more and more interested in utilization of natural light. It is very important for human to utilize natural light in order to reducing energy consumption for lighting. The key factor that the building use natural light for lighting is sky luminance distribution, so the law of sky luminance distribution is the main object which we study on [1].

Daylight climate data is used as important reference for studying the law of sky luminance distribution [11]. For obtaining better results, it requires higher quality climate data. In this paper, we are concerned with the application of data preprocessing in daylight climate data. The rest of this paper is structured as follows. In the next section, we give a brief description of data quality and data preprocessing particularly. Section 3 presents a data preprocessing approach for daylight climate data. Section 4 the simulation experiment flow chart and results are displayed and the conclusions in the last section are stated.

2 Data Quality and Data Preprocessing

Real world has amounts of data. More data, more questions. We divide the questions into three categories approximately: incomplete data (lacking attribute values), noisy (containing errors, or outlier that deviate from the expected) and inconsistent (e.g., containing discrepancies in the department codes used to categorize items) [2,12]. It is well known that garbage in, garbage out. If the data set is low-quality, we can't get a high-quality result. Data preprocessing is applied to do a series of processing steps such as cleaning, integration, transformation and reduction in original data. Simply speaking, data preprocessing is a process that turn the dirty data into clean data. It can help improve the quality of the original data. Besides it can also improve the efficiency and accuracy of the next process. Quality decisions must be based on quality data, thereby the importance of data preprocessing is obviously.

2.1 Data Quality Index

Data quality reflects the value of the data. The higher quality of the data is, the greater value of the data is. The next application of the data must be based on high quality data and data preprocessing is to output the high quality data. How can we measure the data quality? In different fields, the standards of measure are different. In literature [3], the data quality is defined as the satisfaction of the consistency, accuracy, completeness, minimality. The literature [4] proposes "being appropriate for use" as the preliminary standard. In general, the standards of measure data quality are as follows [5]:

- (1) Accuracy: the consistent extent of correct data value and the data value from data source, in other word, it requires that the "noisy" data is as little as possible.
- (2) Completeness: require no missing data or missing attribute value in data source.
- (3) Uniqueness: require no duplicate data in data source.
- (4) Validity: require data from data source to meet the conditions defined by user or be within a certain range of threshold.

2.2 Basic Methods of Data Preprocessing

Nowadays, there are a number of data preprocessing methods and the frequently-used ones of them are in the following [2]:

- (1) Data Cleaning: Because the original data are incomplete and noisy, the tasks of data cleaning include eliminating the noisy data and irrelevant data, filling in missing data, identifying or removing outliers, and resolving the data inconsistency problem.
- (2) Data Integration: As the data may from many different data stores, it emerges problems when translate the heterogeneous data. The main task of data integration is to fuse the heterogeneous data and combine them into coherent data storage.
- (3) Data Transformation: According to the original data's character representation, data transformation converts the data into the specific format which meet the demand of the data processing.

(4) Data Reduction: The definition of data reduction is compressing the data under the premise of ensuring the data's completeness and correctness by the methods of data cluster, dimension reduction and so on.

The methods of data preprocessing are not incompatible but correlated. For example, the data smoothing of the noisy data is not only a form of data cleaning but also a form of data transformation. In practical applications, more than one data preprocessing methods will commonly be used.

3 Data Preprocessing Model

By integrating corresponding domain knowledge and business requirements, this paper improves the data quality further in the data preprocessing of daylight climate. Besides, it can also make the amounts of data get more reasonable reduction.

3.1 The Analysis of the Original Data

The data in this paper come from the "International Lighting Observation Year" project, launched by the International Commission on Illumination (CIE) and the International Meteorological Organization. It includes the data more than two and a half year. In the observation, the data acquisition systems controlled by computer were used to observe the solar radiation, luminance, sun elevation angle and zenith luminance etc. totally 15 features between sunrise and sunset every day and every minute. The observation equipment is manufactured by national professional factories of our country. The measurement methods are satisfied CIE's requirements. All devices are measured annually by the authority of the national identification departments to ensure the accuracy of observed data.

Before determining how to improve the data quality, it has to make sure that what concrete problems of data quality may exist in the original data. In the observation, there are errors which probably cause by the failures of data acquisition systems or operation carelessness. The main quality problems include accuracy, completeness, uniqueness, effectiveness.

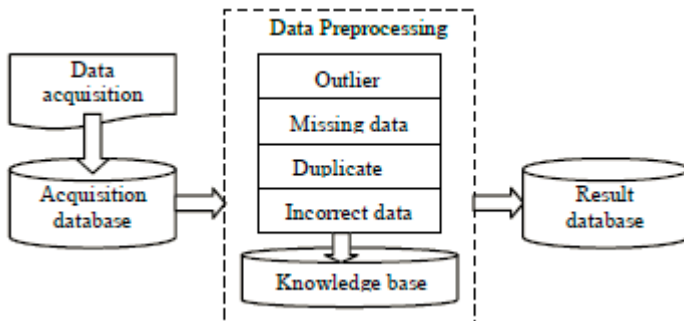


Fig. 1. The model chart of data preprocessing of daylight climate

By analysis of the original data above and integration the data next processing requirements, this article proposes a data preprocessing model appropriate for the daylight climate data, the model focuses on four aspects: outlier, missing data, duplicate data and incorrect data, as Fig.1.

3.2 The Data Preprocessing Method of Daylight Climate

Data Preprocessing and the Rules. Through analysis of the daylight climate data quality, we can conclude that there exist data quality problems in the original data.

1) Outlier

There are a certain amount of exception values, generally termed as outliers in data source. It maybe results from measurement or execution error. Alternatively, outlier may be the result of inherent data variability. Many algorithms try to minimize the influence of outliers or eliminate them all together. Taking the characteristics of daylight climate data into consideration, this article use the outlier detection approach mentioned in literature [6]. We must define the conformity of an attribute value firstly and then detect outlier according to it. If the conformity of an attribute value is out of range, it is called an outlier. For daylight climate data, the conformity of an attribute value by the following membership function u_R :

$$u_R(V_{ij}) = 2 / (1 + e^{\frac{D_i}{TV_i \times N_{ij}}}) \tag{1}$$

Where V_{ij} is the j -th value of the attribute A_i ; γ is the shape factor, representing the user attitude towards conformity of rare values; D_i is the total number of records where $A_i \neq null$; TV_i is the total number of distinct values taken by the attribute A_i ; and N_{ij} is the number of occurrences of the value V_{ij} . It can be easily verified that the Equation (1) agrees with the definition of fuzzy measure (see[7]).The conformity becomes close to zero, when the number of value occurrences is much smaller than the average number of records per value (given by D_i/TV_i). On the other hand, if the data set is very small (D_i close to zero), u_R approaches one, which means that even a single occurrence of a value is not considered an outlier. The subset of outlier in each attribute is found by defining an α -cut of all the attribute values :{ $V_{ij}: u_R < \alpha$ }. In real application, we define the conformity of a record as the average of the sum of the conformity of every attribute value.

Rule R1: IF $(u_R(V_{0j}) + \dots + u_R(V_{nj})) / n < \alpha$ THEN delete record j .

2) Missing data

The incomplete data is an important factor that causes the problem of data quality. Because of equipment malfunction and worker carelessness in inputting data, these factors lead to the missing problem in the original data. As a rule, there are usually two solutions to solve the missing data: delete or fill. The missing data includes two cases:

Attribute value missing

Calculate the ratio (A_i) of missing attribute values for each record, which is equaled to the number of missing attribute values divided by the number of values with the

whole record. Set a threshold β , if the value of A_i is smaller than β , then fill the missing value; otherwise, delete this record.

Rule R2: IF $A_i \geq \beta$ THEN delete record i .

There are many methods in filling the missing values [8, 9], such as attribute average substitution, attribute common value substitution, and so on. To ensure the accuracy of the data and minimize the error, we use the weighted average value of attributes to fill the missing value. Order by the record-time, we find k records which are the nearest to the missing value, then use the weighted average value which corresponds to k records to fill the missing field value. The record is nearer to the missing value, the weight is bigger. It's worth noting that it needs to judge that the field values of the k records are missed or not. The missing ratio M_i is equal to the number of missing field values divided by k . Set a threshold γ , if the value of M_i is smaller than γ , fill the missing value; otherwise, delete this record.

Rule R3: IF $A_i < \beta$ and $M_i \geq \gamma$ THEN delete the record.

Rule R4: IF $A_i < \beta$ and $M_i < \gamma$ THEN fill the missing value by the weighted average of attributes value.

Record missing

The record missing ratio (D_i) is the number difference of expected records and actual records divided by the number of expected records. Set a threshold δ , if the value of D_i is smaller than δ , fill this record; otherwise, delete all records which belong to this date.

Rule R5: IF $D_i \geq \delta$ THEN delete all records of the date D_i .

Because the record is measured day by day and minute by minute, the base of judging the missing record is that the attribute *time* has to be continuous. Before filling the missing record, we need to judge whether exists continuous missing in missing time. Set a threshold ε of continuous missing time point, if the continuous missing time point is smaller than ε , then fill this record; otherwise, delete all records which belong to this date.

Rule R6: IF $D_i < \delta$ and the number of continuous missing record $\geq \varepsilon$ THEN delete all records of the date D_i .

Rule R7: IF $D_i > \delta$ and the number of continuous missing record $< \varepsilon$ THEN fill the missing value by the weighted average of attributes value mentioned above.

3) Duplicate data

Because the data is measured day by day and minute by minute, it is not allowed to exist no duplicate time point record in database. Find the duplicate records (assumed record A and record B), then we can find k records near to the duplicate records according to the attribute time. Calculate the distance between the duplicate records A , B and k records separately. Stay the shortest distance, and remove the others. The distance function is as follows:

$$d(X_a, Y) = \sum_{i=1}^k \sum_{j=1}^n |x_{a_i} - y_{ij}| \quad (2)$$

Rule R8: IF $d(X_A, Y) < d(X_B, Y)$ THEN delete the record B .

4) Incorrect data

The range of every attribute

Every attribute in data source has its range. Out of its range are called the incorrect data, it must delete. Set V_i stands an attribute value.

Rule R9: IF $V_i \notin$ its range THEN delete this record.

The form of every attribute

Every attribute in data source has its specific form. The data with error form are considered the incorrect data. For example, the attribute *time* form is "***:**:00".

Rule R10: IF the attribute value doesn't meet the attribute form THEN modify this attribute value.

Data Normalization. Since the data in the data source is usually composed of multiple attributes, and the unit of measurement is different to the diverse attributes. The smaller unit of measurement used, the larger the range of the attributes value is, which has a greater impact on calculation results. In order to avoid the dependence on the choice of the measurement units, data should be normalized before we do the calculation. Therefore, data normalization is an important auxiliary in the data preprocessing. The minimum-maximum normalization is used in this paper. This method is a linear transformation of the original data, so that each attribute is in a specific interval. The advantage of this method is that the relationship among the original data is still remained after the normalization. The formula is as follows:

$$v' = \frac{new_max - new_min}{max_A - min_A} (v - min_A) + new_min \quad (3)$$

Where min_A and max_A represent the minimum value and maximum value of attribute A respectively, the value of the original data of attribute A mapped into the interval $[new_min, new_max]$.

4 The Experiment of the Data Preprocessing

Based on the model of the data preprocessing of daylight climate data, we develop the data preprocessing system for daylight climate data as Fig.2.

Data normalization is to normalize the data in acquisition database and then store the data into result database.

Rules R1 to R10 are composed of the knowledge base of the model of the data preprocessing. The structure of rules in the knowledge base is (rule ID, rule name, the condition of the rule, the process function), where the condition is a logical expression, for example, Rule R3, the condition is $A_i \geq \beta$, the process function is `doDeleteRow()`.

We take a certain strategies when choose the correspondent rule to match the known facts in database during data preprocessing. If match successfully, then use it

to handle the data and use the result data to replace the original data. We choose the rules as follows:

R9→R10→R2→R3→R4→R8→R1→R5→R6→R7

After preprocessing, the changes of data quality are shown in Table 1.

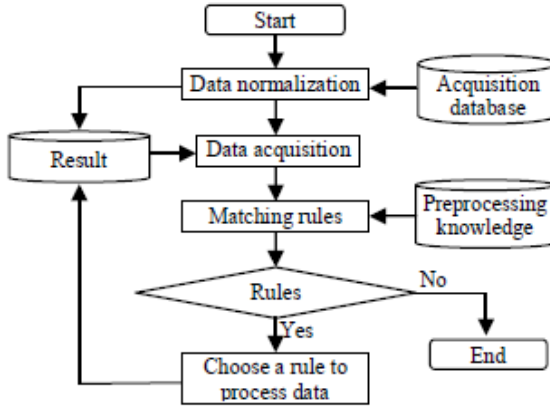


Fig. 2. The flow of the data preprocessing model

Table 1. The result data after data preprocessing

The situation of the original data	The number of records before data preprocessing	The number of records after data preprocessing
total records	367798	309250
incorrect value	852	delete records:852
incorrect form	472	modify records: 472
duplicate data	408	delete records:204
missing value	653	delete records:501, fill records: 152
missing record	310	delete records:56991, fill records: 27

5 Conclusion

There are a variety of quality problems existing in the real-world data. However, the high quality decision depends on the high quality data. Generally speaking, the data preprocessing can fill the missing data, correct the incorrect data, remove the duplicate data, modify the date format into the required, and also can eliminate the irrelevant attributes, so as to make sure the data consistency, simplicity and accuracy. In this paper, data preprocessing technology is applied to handle the light climate data and improve the data quality effectively.

Acknowledgments. This work was supported by the National Natural Science Foundation of China-Youth Fund (Grant No. 1010200220090070).

References

1. Zhang, N., Lu, W.F.: An Efficient Data Preprocessing Method for Mining Customer Survey Data. In: Proceedings of the 5th IEEE International Conference on Industrial Informatics, vol. 1, pp. 573–578. IEEE Press, Vienna (2007)
2. Han, J.W.: Data Mining: Concepts and Techniques. Higher Education Press, Beijing (2006)
3. Aebi, D., Perrochon, L.: Towards improving data quality. In: Proceedings of the International Conference on Information Systems and Management of Data, pp. 273–281. Institution of Engineers, Delhi (1993)
4. Yu, H.R.: The key technologies research for Data quality and data cleaning. Msc thesis, Fudan University (2002)
5. Tayi, G.K., Ballou, D.P.: Examining data quality. Communications of the ACM 41(2), 54–57 (1998)
6. Last, M., Kandel, A.: Automated detection of outliers in real-world data. In: Proceedings of the 2nd International Conference on Intelligent Technologies, Bangkok, Thailand, pp. 292–301 (2001)
7. Klir, G.J., Yuan, B.: Fuzzy Sets and Fuzzy Logic: Theory and Applications. Prentice Hall, Upper Saddle River (1995)
8. Grzymala-Busse, J.W., Hu, M.: A Comparison of Several Approaches to Missing Attribute Values in Data Mining. In: Ziarko, W.P., Yao, Y. (eds.) RSCTC 2000. LNCS (LNAI), vol. 2005, pp. 378–385. Springer, Heidelberg (2001)
9. Gustavo, E.A., Batista, P.A., Monard, M.C.: An analysis of four missing data treatment methods for supervised learning. Applied Artificial Intelligence 17(5-6), 519–533 (2003)
10. Zou, Y., An, A.J., Huang, X.J.: Evaluation and automatic selection of methods for handling missing data. In: IEEE International Conference on Granular Computing, vol. 2, pp. 728–733. IEEE Press, Beijing (2005)
11. He, Y., Guo, P., Lin, Y.: Study on the sky luminance distribution of information methods by ant colony systems. Applied Mechanics and Materials 48-49, 1202–1207 (2011)
12. Guo, Z.M., Zhou, A.Y.: Research on Data Quality and Data Cleaning: a Survey. Journal of Software 13(11), 2076–2082 (2002)
13. Ordonez, C.: Data set preprocessing and transformation in a database system. Intelligent Data Analysis 15(4), 613–631 (2011)
14. Heinrich, J., Elter, T., Ulrich, J.: Data Preprocessing of In Situ Laser-Backscattering Measurements. Chemical Engineering and Technology 34(6), 977–984 (2011)
15. Elmagarmid, A.K., Ipeirotis, P.G., Verykios, V.S.: Duplicate Record Detection: A Survey. IEEE Transactions on Knowledge and Data Engineering 19(1), 1–16 (2007)
16. Outrata, J.: Boolean factor analysis for data preprocessing in machine learning. In: Proceedings of the 9th International Conference on Machine Learning and Applications, pp. 899–902. IEEE Computer Society, Washington, D.C. (2010); IEEE Transactions on Knowledge and Data Engineering
17. Nick, J.P.: Fuzzy quartile encoding as a preprocessing method for biomedical pattern classification. Theoretical Computer Science 412(42), 5909–5925 (2011)

Scalable Technique to Discover Items Support from Trie Data Structure

A. Noraziah¹, Zailani Abdullah², Tutut Herawan¹, and Mustafa Mat Deris³

¹ Faculty of Computer Systems and Software Engineering,
Universiti Malaysia Pahang Lebuhraya Tun Razak, 26300 Kuantan Pahang, Malaysia

² Department of Computer Science,
Universiti Malaysia Terengganu, 21030 Kuala Terengganu, Terengganu, Malaysia

³ Faculty of Computer Science and Information Technology,
Universiti Tun Hussein Onn Malaysia, Parit Raja, Batu Pahat 86400, Johor, Malaysia
{noraziah, tutut}@ump.edu.my, zailania@umt.edu.my,
mmustafa@uthm.edu.my

Abstract. One of the popular and compact trie data structure to represent frequent patterns is via frequent pattern tree (FP-Tree). There are two scanning processes involved in the original database before the FP-Tree can be constructed. One of them is to determine the items support (items and their support) that fulfill minimum support threshold by scanning the entire database. However, if the changes are suddenly occurred in the database, this process must be repeated all over again. In this paper, we introduce a technique called Fast Determination of Item Support Technique (F-DIST) to capture the items support from our proposed Disorder Support Trie Itemset (DOSTrieIT) data structure. Experiments through three UCI benchmark datasets show that the computational time to capture the items support using F-DIST from DOSTrieIT is significantly outperformed the classical FP-Tree technique about 3 orders of magnitude, thus verify its scalability.

Keywords: Frequent Pattern Tree, Trie Data Structure, Fast Technique.

1 Introduction

Since the introduction of Apriori by Agrawal et al. [4], frequent pattern mining has received a great deal of attentions and explorations [1,2,3,20]. Frequent pattern mining tries to find interesting patterns from database such as association rules, correlation, causality, sequences, clusters etc. Until know, more than hundreds of research papers have been published based on developing or enhancing the techniques and data structures. Generally, the main problem in frequent patterns mining is how to manage the huge data in computer's memory in efficient manner. One of the remarkable solutions is to store the data in trie data structure.

Frequent pattern tree (FP-Tree) [1] is one of them and already became a benchmarked trie data structure in frequent pattern mining. Since that, several variations of constructing or updating the FP-Tree have been proposed and discussed

in the literature [1,5,6,7,8,9,10,11]. However, there are still two major drawbacks encountered. First, when the existing database is updated, all items that fulfill the minimum support must be recounted again. Second, the current FP-Tree is probably not valid anymore and must be rebuilt all over again due to the latest items support. Thus in order to improve the overall computational time, the first part of the problem is worth to consider. Therefore, we proposed a scalable technique called Fast Determination of Item Support Technique (F-DIST) and enhanced trie namely Disorder Support of Trie Itemset (DOSTrieIT). Three datasets from Frequent Itemset Mining Datasets from Repository [21] were employed in the experiments. The performance analysis between F-DIST and the benchmarked FP-Tree technique was performed in determining the items support.

In summary, there are three main contributions from this work. First, we propose a novel, complete and incremental pattern tree data structure, DOSTrieIT that can keep the entire transactional database. Second, we embed a feature called Single Item Without Support (SIWE) in DOSTrieIT to speed up the process of capturing the items support. Lastly, we do the experiments with three benchmarked dataset from UCI Repository [21] and show the scalability of F-DIST.

The paper structure is organized as follows. Section 2 explains the related works. In section 3, the basic concept and terminology in association rules is discussed. Section 4 elaborates the proposed methods. Detail discussions of the experiments are reported in section 5. Finally, Section 5 concludes the paper.

2 Related Works

Since the introduction by Agrawal et al. in 1993 [4], more than hundreds of papers have been intensively published in an attempt to increase its efficiencies and scalabilities. In general, the algorithms for mining the frequent itemset could be classified into three; Apriori-like algorithms, frequent pattern-based algorithms and algorithms that use the vertical data format.

Due to the problem of two nontrivial costs in Apriori [12], the frequent pattern based technique without candidate itemsets generation has been proposed. This technique constructs a compact trie data structure known as FP-Tree [1] from the original database. Typically, before constructing the FP-Tree, the items that satisfy the minimum support threshold will be captured first from database. Then, the items from each transaction will be extracted and sorted in support descending order before they can be transformed into FP-Tree.

Since the introduction of FP-Tree, there are abundant researches have been put forward such as H-Mine Algorithms [13], PatriciaMine [14], FPgrowth* [15], SOTrieIT [16], AFOPF [5], AFPIM [6] and EFPIM [7], CATS-Tree [8], CanTree [17], FUFPT [9], CP-Tree [18], BSM [19], BIT [11]. However, most of the mentioned above approaches need to employ the original database to determine the current items support. Moreover, to the best of our knowledge, none of the research has been conducted so far to efficiently determine the items support from trie data structure. In summary, most of them are still depend on the original database.

3 Association Rules

Throughout this section the set $I = \{i_1, i_2, \dots, i_{|A|}\}$, for $|A| > 0$ refers to the set of literals called set of items and the set $D = \{t_1, t_2, \dots, t_{|U|}\}$, for $|U| > 0$ refers to the data set of transactions, where each transaction $t \in D$ is a list of distinct items $t = \{i_1, i_2, \dots, i_{|M|}\}$, $1 \leq |M| \leq |A|$ and each transaction can be identified by a distinct identifier TID. A set $X \subseteq I$ is called an itemset. An itemset with k-items is called a k-itemset. The support of an itemset $X \subseteq I$, denoted $\text{supp}(X)$ is defined as a number of transactions contain X. Let $X, Y \subseteq I$ be itemset. An association rule between sets X and Y is an implication of the form $X \Rightarrow Y$, where $X \cap Y = \emptyset$. The sets X and Y are called antecedent and consequent, respectively. The support for an association rule $X \Rightarrow Y$, denoted $\text{supp}(X \Rightarrow Y)$, is defined as a number of transactions in D contain $X \cup Y$. The confidence for an association rule $X \Rightarrow Y$, denoted $\text{conf}(X \Rightarrow Y)$ is defined as a ratio of the numbers of transactions in D contain $X \cup Y$ to the number of transactions in D contain X. Thus $\text{conf}(X \Rightarrow Y) = \text{supp}(X \Rightarrow Y) / \text{supp}(X)$. An itemset X is called frequent item if $\text{supp}(X) > \beta$, where β is the minimum support. The set of frequent item will be denoted as Frequent Items and Frequent Item = $\{X \subset I \mid \text{supp}(X) > \beta\}$. Frequent pattern mining has been received a great deal of attentions from data mining researchers [22-31].

4 Proposed Model

4.1 Definition

In order to easily comprehend the whole process that is performed in DOSTrieIT, some required definitions and a sample transactional data are presented.

Definition 1. Disorder Support Trie Itemset (DOSTrieIT) is defined as a complete tree data structure in canonical order of itemsets. The order of itemset is not based on the support descending order. DOSTrieIT contains n-levels of tree nodes (items) and their support. Moreover, DOSTrieIT is constructed in online manner and for the purpose of incremental pattern mining.

Example 1

Let $T = \{\{1, 2, 5\}, \{2, 4\}, \{2, 3\}, \{1, 2, 4\}, \{1, 3\}, \{2, 3, 6\}, \{1, 3\}, \{1, 2, 3, 5\}, \{1, 2, 3\}\}$.

A step by step to construct DOSTrieIT is explained in the next section. Graphically, an item is represented as a node and its support is appeared nearby to the respective node. A complete structure of DOSTrieIT is shown as in Fig. 1.

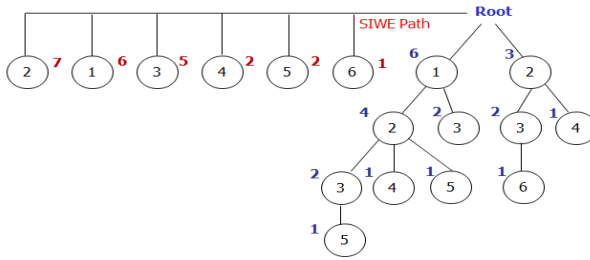


Fig. 1. DOSTrieIT and SIWE path arranged in support descending order

Definition 2. Single Item without Extension (SIWE) is a prefix path in the tree that contains only one item or node. SIWE is constructed upon receiving a new transaction and as a mechanism for fast searching of single item support. It will be employed during tree transformation process but it will not be physically transferred into the others tree.

Example 2. From Example 1, the transactions have 6 unique items and it is not sorted in any order. In Fig. 2, SIWE for DOSTrieIT i.e. $SIWE = \{2,1,3,4,5,6\}$

Proposition 1. (Instant Support of Single Items Property). For any item a_i , the items support is instantly obtained from the 1-level of DOSTrieIT. All these items or nodes have no extension or also known as SIWE.

Justification. Let single item a_1, a_2, \dots, a_n , from Definition 8, Single Item without Extension (SIWE) a_1, a_2, \dots, a_n is a prefix path in the tree that contains only one item or node. In this case a_1, a_2, \dots, a_n is constructed upon receiving a new transaction. To this we can accelerate the process of updating and/or searching a support of a_1, a_2, \dots, a_n . The trie-traversal for examining a support of a_1, a_2, \dots, a_n is truncated once it reaches at the last of single items without extension. It will be employed during tree transformation process but it will not be physically transferred into the FP-Tree.

Example 3. Let examine the sample transaction from Example 1. The items support of 2, 1, 3, 4, 5, 6 is 7, 6, 5, 2, 2, 1 respectively. For CanTree and CATS-Tree, the support information is only can be captured after scanning 9 lines of transactions. However, the similar information can be easily determined from DOSTrieIT via SIWE as shown in Fig. 1. The items support is obtained in DOSTrieIT by traversal in the trie and immediately stopped after no more single items without extension is found.

4.2 Activity Diagrams

Activity diagram is employed in visualizing the details processes of constructing DOSTrieIT. It is one of the prominent diagrams in Unified Modeling Language (UML) to graphically represent the workflows of stepwise with support of choice (condition), iteration (loop) and concurrency. Fig. 2 and Fig. 3 show the activity diagrams for constructing DOSTrieIT data structure and F-DIST, respectively.

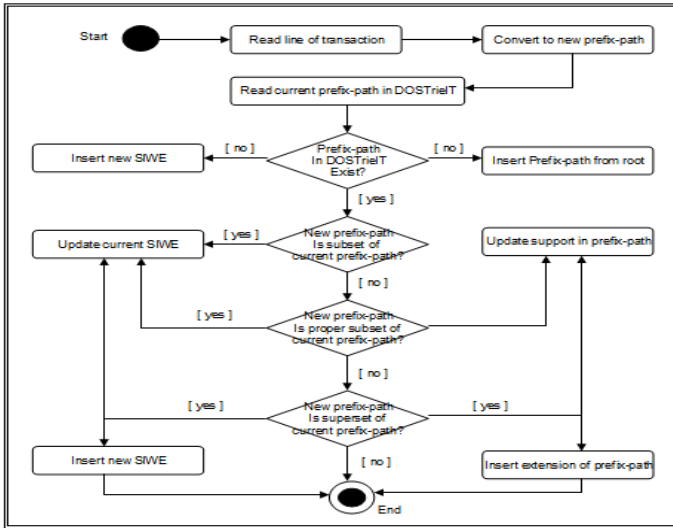


Fig. 2. An activity diagram for DOSTrieT

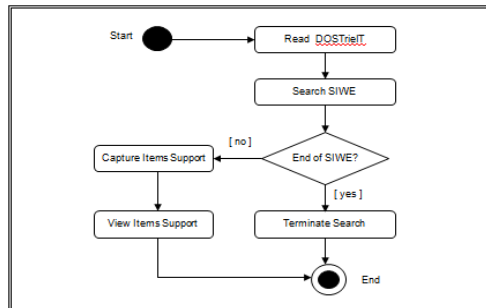


Fig. 3. An activity diagram for F-DIST

5 Experimental Setup

In this section, we do comparison tests between F-DIST and benchmarked FP-Tree technique. The performance analysis was carried out by comparing the computational time required to determine the items support from the respective data source. In term of system specifications, the experiment was conducted on Intel® Core™ 2 Quad CPU at 2.33GHz speed with 4GB main memory, running on Microsoft Windows Vista. All coding have been developed using C# as a programming language.

Three benchmarked datasets from Frequent Itemset Mining Dataset Repository [21] were employed in the experiment. The first dataset was synthetic dataset T10I4D100K was used. It is a sparse dataset. In this dataset, the frequent itemsets are short but they are not abundant. For the second experiment, the dataset was Retail and

it contains the retail market basket data from an anonymous Belgian retail store. The third benchmarked dataset was Pumsb. The dataset contains census data for population and housing. It is one of the difficult data sets to mine because of the long average record width coupled with the high number of popular attribute-value pairs. The fundamental characteristics of the datasets are depicted in Table 1.

Fig. 4 shows the comparison between F-DIST and FP-Tree technique in term of duration taken (or processing time) to capture the items support. In overall, duration to determine the items support using F-DIST was less than FP-Tree technique. For T10I4D100K and Pumsb datasets, processing time via F-DIST technique was 3 orders of magnitude faster than FP-Tree technique. Based on Retail dataset, the performance of F-DIST was 2 orders of magnitude better than FP-Tree technique. In summary, the average duration to determine the items support by F-DIST was nearly 3 orders of magnitude better than FP-Tree technique.

Table 1. Fundamental characteristics of datasets

Datasets	Size	#Trans	#Items	Average length
T10I4D100K	3.83 MB	100,000	1,000	10
Retail	4.15 MB	88,136	16,471	10
Pumsb	16.59 MB	49,046	2,113	74

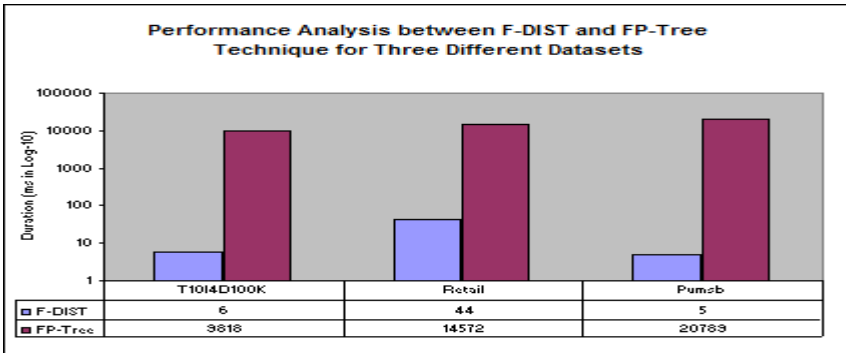


Fig. 4. Performance analysis for determining the items support using three different datasets

6 Conclusion

FP-Tree is a variation of the trie data structure for storing the frequent patterns in compressed manner. Typically, two scanning processes are performed before FP-Tree can be completely built. One of them is to capture the list of items that fulfils the minimum support. Only these items will be captured from each transaction in the database during the next scanning process. However when the current content in database is changed, the FP-Tree must be rebuilt again. Prior to this, the first scanning process must be repeated all over again in order to obtain the valid list of items. At the moment, most of the trie-based techniques are still depend on the original dataset to capture these items rather than their own trie data structure. In order to resolve this

problem, we proposed a scalable technique called F-DIST to capture the items and their support using our proposed DOSTrieIT data structure. We do the experiment with three benchmarked datasets from Frequent Itemset Mining Dataset Repository [21] datasets and found that our proposed technique is outperformed at 3 order of magnitude faster than benchmarked FP-Tree technique.

Acknowledgement. This research is supported by Fundamental Research Grant Scheme (FRGS) from Ministry of Higher Education of Malaysia No. Vote RDU 100109.

References

1. Han, J., Pei, H., Yin, Y.: Mining Frequent Patterns without Candidate Generation. In: Proceeding of the 2000 ACM SIGMOD, pp. 1–12 (2000)
2. Zheng, Z., Kohavi, R., Mason, L.: Real World Performance of Association Rule Algorithms. In: Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pp. 401–406. ACM Press (August 2001)
3. Han, J., Pei, J.: Mining Frequent Pattern without Candidate Itemset Generation: A Frequent Pattern Tree Approach. *Data Mining and Knowledge Discovery* 8, 53–87 (2004)
4. Agrawal, R., Imielinski, T., Swami, A.: Database Mining: A Performance Perspective. *IEEE Transactions on Knowledge and Data Engineering* 5(6), 914–925 (1993)
5. Liu, G., Lu, H., Lou, W., Xu, Yu, J.X.: Efficient Mining of Frequent Patterns using Ascending Frequency Ordered Prefix-Tree. *Data Mining and Knowledge Discovery* 9, 249–274 (2004)
6. Koh, J.-L., Shieh, S.-F.: An Efficient Approach for Maintaining Association Rules Based on Adjusting FP-Tree Structures. In: Lee, Y., Li, J., Whang, K.-Y., Lee, D. (eds.) DASFAA 2004. LNCS, vol. 2973, pp. 417–424. Springer, Heidelberg (2004)
7. Li, X., Deng, Z.-H., Tang, S.-W.: A Fast Algorithm for Maintenance of Association Rules in Incremental Databases. In: Li, X., Zaïane, O.R., Li, Z.-H. (eds.) ADMA 2006. LNCS (LNAD), vol. 4093, pp. 56–63. Springer, Heidelberg (2006)
8. Cheung, W., Zaïane, O.R.: Incremental Mining of Frequent Patterns without Candidate Generation of Support Constraint. In: Proceeding of the 7th International Database Engineering and Applications Symposium, IDEAS 2003 (2003)
9. Hong, T.-P., Lin, J.-W., We, Y.-L.: Incrementally Fast Updated Frequent Pattern Trees. *An International Journal of Expert Systems with Applications* 34(4), 2424–2435 (2008)
10. Tanbeer, S.K., Ahmed, C.F., Jeong, B.S., Lee, Y.K.: Efficient Single-Pass Frequent Pattern Mining Using a Prefix-Tree. *Information Science* 279, 559–583 (2009)
11. Totad, S.G., Geeta, R.B., Reddy, P.P.: Batch Processing for Incremental FP-Tree Construction. *International Journal of Computer Applications* 5(5), 28–32 (2010)
12. Agrawal, R., Shafer, J.: Parallel Mining of Association Rules: Design, Implementation, and Experience. *IEEE Transaction Knowledge and Data Engineering* 8, 962–969 (1996)
13. Pei, J., Han, J., Lu, H., Nishio, S., Tang, S., Yang, D.: Hmine: Hyper-Structure Mining of Frequent Patterns in Large Databases. In: The Proceedings of IEEE International Conference on Data Mining, pp. 441–448 (2001)
14. Pietracaprina, Zandolin, D.: Mining Frequent Item sets Using Patricia Tries. In: The Proceedings of the ICDM 2003 (2003)
15. Grahne, G., Zhu, J.: Efficiently using prefix-trees in mining frequent itemsets. In: Proceeding of FIMI 2003 (2003)

16. Woon, Y.K., Ng, W.K., Lim, E.P.: A Support Order Trie for Fast Frequent Itemset Discovery. *IEEE Transactions on Knowledge and Data Engineering* 16(7), 875–879 (2004)
17. Leung, C.K.-S., Khan, Q.I., Li, Z., Hoque, T.: CanTree: A Canonical-Order Tree for Incremental Frequent-Pattern Mining. *Knowledge Information System* 11(3), 287–311 (2007)
18. Tanbeer, S.K., Ahmed, C.F., Jeong, B.-S., Lee, Y.-K.: CP-Tree: A Tree Structure for Single-Pass Frequent Pattern Mining. In: Washio, T., Suzuki, E., Ting, K.M., Inokuchi, A. (eds.) *PAKDD 2008*. LNCS (LNAI), vol. 5012, pp. 1022–1027. Springer, Heidelberg (2008)
19. Tanbeer, S.K., Ahmed, C.F., Jeong, B.-S., Lee, Y.-K.: Sliding Window-based Frequent Pattern Mining Over Data Streams. *Information Sciences* 179, 3843–3865 (2009)
20. Ivancsy, R., Vajk, I.: Fast Discovery of Frequent Itemsets: a Cubic Structure-Based Approach. *Informatica (Slovenia)* 29(1), 71–78 (2005)
21. Frequent Itemset Mining Dataset Repository, <http://fimi.ua.ac.be/data/>
22. Abdullah, Z., Herawan, T., Deris, M.M.: An Alternative Measure for Mining Weighted Least Association Rule and Its Framework. In: Zain, J.M., Wan Mohd, W.M.B., El-Qawasmeh, E. (eds.) *ICSECS 2011, Part II*. CCIS, vol. 180, pp. 480–494. Springer, Heidelberg (2011)
23. Herawan, T., Yanto, I.T.R., Deris, M.M.: Soft Set Approach for Maximal Association Rules Mining. In: Ślęzak, D., Kim, T.-H., Zhang, Y., Ma, J., Chung, K.-I. (eds.) *DTA 2009*. CCIS, vol. 64, pp. 163–170. Springer, Heidelberg (2009)
24. Abdullah, Z., Herawan, T., Deris, M.M.: Mining Significant Least Association Rules Using Fast SLP-Growth Algorithm. In: Kim, T.-H., Adeli, H. (eds.) *AST/UCMA/ISA/ACN 2010*. LNCS, vol. 6059, pp. 324–336. Springer, Heidelberg (2010)
25. Herawan, T., Deris, M.M.: A soft set approach for association rules mining. *Knowledge Based Systems* 24(1), 186–195 (2011)
26. Abdullah, Z., Herawan, T., Noraziah, A., Deris, M.M.: Extracting Highly Positive Association Rules from Students' Enrollment Data. *Procedia Social and Behavioral Sciences* 28, 107–111 (2011)
27. Abdullah, Z., Herawan, T., Noraziah, A., Deris, M.M.: Mining Significant Association Rules from Educational Data using Critical Relative Support Approach. *Procedia Social and Behavioral Sciences* 28, 97–101 (2011)
28. Herawan, T., Vitasari, P., Abdullah, Z.: Mining Interesting Association Rules of Student Suffering Mathematics Anxiety. In: Zain, J.M., Wan Mohd, W.M.B., El-Qawasmeh, E. (eds.) *ICSECS 2011, Part II*. CCIS, vol. 180, pp. 495–508. Springer, Heidelberg (2011)
29. Herawan, T., Vitasari, P., Abdullah, Z.: Mining Interesting Association Rules on Student Suffering Study Anxieties using SLP-Growth Algorithm. *International Journal of Knowledge and Systems Science* 3(2), 24–41 (2012)
30. Herawan, T., Yanto, I.T.R., Deris, M.M.: SMARViz: Soft Maximal Association Rules Visualization. In: Badioze Zaman, H., Robinson, P., Petrou, M., Olivier, P., Schröder, H., Shih, T.K. (eds.) *IVIC 2009*. LNCS, vol. 5857, pp. 664–674. Springer, Heidelberg (2009)
31. Abdullah, Z., Herawan, T., Deris, M.M.: Visualizing the Construction of Incremental Disorder Trie Itemset Data Structure (DOSTriEIT) for Frequent Pattern Tree (FP-Tree). In: Badioze Zaman, H., Robinson, P., Petrou, M., Olivier, P., Shih, T.K., Velastin, S., Nyström, I. (eds.) *IVIC 2011, Part I*. LNCS, vol. 7066, pp. 183–195. Springer, Heidelberg (2011)

On Soft Partition Attribute Selection

Rabiei Mamat¹, Tutut Herawan², Noraziah Ahmad², and Mustafa Mat Deris³

¹ Department of Computer Science, Universiti Malaysia Terengganu

² Faculty of Computer System and Software Engineering, Universiti Malaysia Pahang

³ Faculty of Computer Science and Information Technology, Universiti Tun Hussein Onn
Malaysia

rab@umt.edu.my, {tutut,noraziah}@ump.edu.my, mmustafa@uthm.edu.my

Abstract. Rough set theory provides a methodology for data analysis based on the approximation of information systems. It revolves around the notion of discernibility i.e. the ability to distinguish between objects based on their attributes value. It allows inferring data dependencies that are useful in the fields of feature selection and decision model construction. Since it is proven that every rough set is a soft set, therefore, within the context of soft sets theory, we present a soft set-based framework for partition attribute selection. The paper unifies existing work in this direction, and introduces the concepts of maximum attribute relative to determine and rank the attribute in the multi-valued information system. Experimental results demonstrate the potentiality of the proposed technique to discover the attribute subsets, leading to partition selection models which better coverage and achieve lower computational time than that the baseline techniques.

Keywords: Soft set theory, Partition attribute, Attribute relative, Complexity.

1 Introduction

It is no-doubt that rough set theory [1,2] is a well developed approach to solve problems regarding to the investigation of structural relationship among data especially when involving uncertainties. For example, Mazlack et al. [3] exploiting the relationship between attributes by the term Total Roughness (TR) to determine the partition attribute. Later, Parmar et al. [4] proposed a technique called Min-Min Roughness (MMR) to improved the works of Mazlack et al. It is followed by Herawan et al. [5] that proposed a technique called Maximum Dependency Attributes (MDA) to determine the partition attribute by exploiting the dependency relationship among attributes. Although Herawan et al. used a different approach; the technique is still based on the rough sets theory. Recently, a new way for managing uncertain data called soft set theory have been proposed by Molodtsov [6]. Unlike theories before this, soft set theory is free from the inadequacy of the parameterization tools. The usage of parameterization sets make it very convenient and easy to apply for problem solving as shown by Molodtsov in various applications such as the smoothness of function, game theory, operations research, Riemann integration,

Perron integration and measurement theory. Presently, great progresses of study on soft set theory have been made such as in [7,8,9]. At the same time, the great progress on practical applications of soft set theory in various fields also has been achieved such as in [10,11,12]. However, most of the research describes above is heavily rely on Boolean-valued information system and only scarce researchers deals with multi-valued information system. In addition, analysis of attributes relationship in multi-value information system is very important in data mining. In this paper, we proposed an alternative technique to partition attribute selection using soft set theory. The rest of the paper is organized as follows. Section 2 described the soft set theory. In section 3, the proposed technique is described. An experiment and analysis is described in section 4. Finally, conclusion of this work described in section 5.

2 Soft Set Theory

In 1999, Molodtsov [6] proposed soft set theory as a new mathematical tool for dealing with vagueness and uncertainties. At present, work on the soft set theory is progressing rapidly and many important results have been achieved, including the works of [13]. In this section, we review some definitions and properties with regard to soft sets. Throughout this section, U is refers to an initial universe, E is a set of parameters describing objects in U , $P(U)$ is the power set of U and $A \subseteq E$.

Definition 1: (See [6].) A pair (F, A) is called a soft set over U , where F is a mapping given by

$$F : A \rightarrow P(U) \tag{1}$$

In other words, a soft set over U is a parameterized family of the universe U . For $\varepsilon \in A$, $F(\varepsilon)$ may be considered as the set of ε -elements of the soft set $F(A)$ or as the set ε - approximate elements of the soft set. Clearly, a soft set is not a (crisp) set.

Consider the following as an example. Let a universe $U = \{c_1, c_2, c_3, c_4, c_5, c_6, c_7, c_8, c_9, c_{10}\}$ be a set of candidates, a set of parameters $E = \{e_1, e_2, e_3, e_4, e_5, e_6, e_7\}$ be a set of soft skills which stand for the parameters “communicative”, “critical thinking”, “team work”, “information management”, “entrepreneurship”, “leadership” and “moral”, respectively. Consider F be a mapping of E into the set of all subsets of the set U as $F(e_1) = \{c_1, c_2, c_4, c_5\}$, $F(e_2) = \{c_3, c_8, c_9\}$, $F(e_3) = \{c_6, c_9, c_{10}\}$, $F(e_4) = \{c_2, c_3, c_4, c_5, c_8\}$, $F(e_5) = \{c_2, c_5, c_6, c_7, c_8, c_9, c_{10}\}$, $F(e_6) = \{c_6, c_9, c_{10}\}$ and $F(e_7) = \{c_6, c_9, c_{10}\}$. Now consider a soft set (F, E) which describes the “capabilities of the candidates for hire”. According to the data collected, the soft set (F, E) is given by

$$(F, E) = \left\{ \begin{array}{l} e_1 = \{c_1, c_2, c_4, c_5\}, \\ e_2 = \{c_3, c_8, c_9\}, \\ e_3 = \{c_6, c_9, c_{10}\}, \\ e_4 = \{c_2, c_3, c_4, c_5, c_8\}, \\ e_5 = \{c_2, c_5, c_6, c_7, c_8, c_9, c_{10}\}, \\ e_6 = \{c_6, c_9, c_{10}\}, \\ e_7 = \{c_6, c_9, c_{10}\} \end{array} \right\}$$

Definition 2: A quaternion $S = (U, A, V, f)$ is called a information system where U is a nonempty finite set of objects called universe of discourse. A is nonempty finite set of attribute, $V = \bigcup V_r$ where V_r is called the value domain of attribute r , and f is a information function specifying the attributes-value for each object and denoted by $f : U \times A \rightarrow V$.

Definition 3: Let $S = (U, A, V_{\{0,1\}}, f)$ Abe an information system. If $V_a = \{0,1\}$, for every $a \in A$, then $S = (U, A, V_{\{0,1\}}, f)$ is called a Boolean-valued information system.

Proposition 1: Each soft set may be considered as a Boolean-valued information system.

Proof. Let (F, E) be a soft set over the universe U , $S = (U, A, V, f)$ be a information system. Obviously, the universe U in (F, E) may be considered the universe U , the parameter set E may be considered the attributes A . The information function f is defined by

$$f = \begin{cases} 1, & h \in F(e), \\ 0, & h \notin F(e), \end{cases}$$

That is, when $h_i \in F(e_j)$, where $h_i \in U$ and $e_j \in E$, then $f(h_i, e_j) = 1$, otherwise $f(h_i, e_j) = 0$. To this, we have $V = \{0,1\}$. Therefore, a soft set (F, E) may be considered as a Boolean-valued information system $S = (U, A, V_{\{0,1\}}, f)$. For multi-valued information system, it needs to be converted to multi soft-sets [27]. It is based on the notion of a decomposition of a multi valued information system. Let $S = (U, A, V, f)$ be a multi valued information system and $S^i = (U, a_i, V_{\{0,1\}}, f)$, $i = 1, 2, \dots, |A|$ be the $|A|$ binary valued information system. From proposition 1,

$$\begin{aligned}
 S = (U, A, V, f) &= \begin{cases} S^1 = (U, a_i, V_{\{0,1\}}, f) \Leftrightarrow (F, a_1) \\ \vdots \qquad \qquad \qquad \vdots \qquad \qquad \vdots \\ S^{|A|} = (U, a_{|A|}, V_{\{0,1\}}, f) \Leftrightarrow (F, a_{|A|}) \end{cases} \\
 &= ((F, a_i), (F, a_2) .. (F, a_{|A|}))
 \end{aligned}$$

We define $(F, E) = ((F, a_1), (F, a_2) .. (F, a_{|A|}))$ as a multi soft-set over universe U representing a multi valued information system $S = (U, A, V, f)$.

3 The Proposed Technique

Throughout this section, a pair (F, A) refers to multi-soft sets over the universe U representing a categorical valued information system $S = (U, A, V, f)$.

Definition 4: Let (F, A) be a multi soft-sets over the universe U , where $(F, a_i), \dots, (F, a_{|A|}) \subseteq (F, A)$ and $(F, a_i), \dots, (F, a_{|A|}) \subseteq (F, a_i)$. Support of (F, a_i) by (F, a_k) denoted $Sup_{(F, a_k)}(F, a_i)$ is defined as

$$Sup_{(F, a_k)}(F, a_i) = \frac{|(F, a_i) \cap (F, a_k)|}{|(F, a_k)|} \tag{2}$$

Definition 5: Maximum support is a summation of all support which value equal to 1. For each soft set (F, a_i) , maximum support is denoted by $MaxSup_{(F, a_i)}$ and is defined as

$$MaxSup_{(F, a_i)} = \sum Sup_{(F, a_k)}(F, a_i) = 1 \tag{3}$$

Definition 6: Minimum support is a summation of all support with value less than 1. For each soft set (F, a_i) minimum support which denoted by $MinSup_{(F, a_i)}$ is define as

$$MinSup_{(F, a_i)} = \sum Sup_{(F, a_k)}(F, a_i) \neq 1 \tag{3}$$

Proposition 2: If $Mode\{MaxSup_{(F, a_{i_1})}, \dots, MaxSup_{(F, a_{|A|})}\} = 1$, then (F, a_i) is a partition attribute.

Proof. Let (F, a_i) and (F, a_j) be a two soft sets over the universe U , if $a_i = a_j$, then the support value of (F, a_i) is equal to 1, therefore it is said that a_i is relative to a_j . Therefore, a_i can be used to describe a_j and vice versa. If $a_i \neq a_j$, then there exist (F, a_k) where $a_i = a_k$ and $a_j \neq a_k$ then, support value of (F, a_i) is greater than 1. Based on Definition 5, clearly that (F, a_i) is selected as a partition attribute.

Corollary 3: If $Mode\left(\left(\text{Maxsup}_{(F, a_{i_j})}, \dots, \text{Maxsup}_{(F, a_{m_{bl}})}\right) = \text{Max}\right) > 1$ then, $\text{Max}\left(\text{Minsup}_{(F, a_{i_j})}, \dots, \text{Minsup}_{(F, a_{m_{bl}})}\right)$ is a partition attribute.

Proof: The proof is clear from Definition 6 and Proposition 1.

3.1 Complexity

Suppose that in an information system, there are n objects, m attributes and l is the maximum distinct values of each attribute. Computational cost to determining elementary set of all attributes is nm . The proposed technique needs $ml(ml-1)$ times to determine the support for each category. Thus the computational complexity of the proposed technique is the polynomial $O(ml(ml-1)+nm+1)$. The algorithm for the proposed technique is as shown in Fig. 1.

```

Input: Categorical-valued dataset
Output: A partition attribute


---


Begin
1. Builds the multi-soft set approximation
2. Calculate Support, MaxSup and MinSup
   For i = all categories
     For j=all categories
       Intersection=Data(i) And Data(j)
       Sup(i,j)=Intersection / Data(j)
       If Sup = 1 then
         MaxSup(i) = MaxSup(i) + Sup
       Else
         MinSup(i) = MinSup(i) + sup
     End
   End
3. If Mode (MaxSup(data(1)..data(i)))=1 then
   Partition attribute = Maxsup(Data(1)..Data(n))
4. else
   Partition Attribute = Minsup(Data(1)..Data(n))
End

```

Fig. 1. Algorithm for the proposed method

4 Experiment Results

This section explains and discusses experiments that were done using the proposed technique. These experiments are focused on performance measurement using execution time as a parameter. Two previous techniques called Min-Min Roughness (MMR) and Maximum Dependency Attribute (MDA) are used as the comparison factors. Two dataset from UCI repository are used in this experiment. The first dataset is lung cancer dataset purposely to analyst the performance when involving high number of attribute. Second is mushroom dataset purposely used to analyst the performance when involving scalability. Both two datasets have been modified by removing instances with missing values and attribute with single categorical value. All experiments are implemented using Ruby programming language version 1.9 under Windows 7 Home Edition operating system powered by Intel i5 processor with 4 GB memory.

4.1 Modified Lung Cancer Dataset

Originally, lung cancer dataset is comprised of 32 instances with 56 attributes. A modified lung cancer dataset is obtained by removing an attribute with single value and deletes five objects that have incomplete attribute. New modified lung cancer dataset is a collection of 27 instances with 55 attributes. Result of the experiment shows that MMR technique needs an approximately 0.25 seconds of execution times, while MDA technique consumes 0.15 seconds of execution time. Compared to two above techniques, MAR technique is a bit faster where it only requires 0.12 seconds to accomplish the process. Fig. 1 above clearly shows that there is an improvement of execution time given by MAR technique. But, the percentage of improvement apparently is not very significant since the process is involving a small instance of dataset with a large number of attribute.

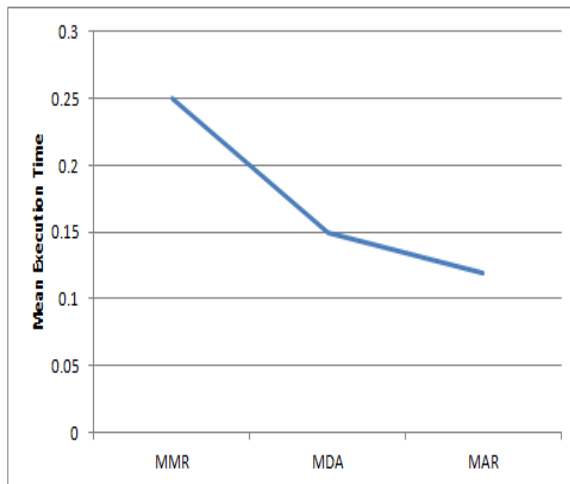


Fig. 2. Comparison of mean execution times using lung cancer dataset

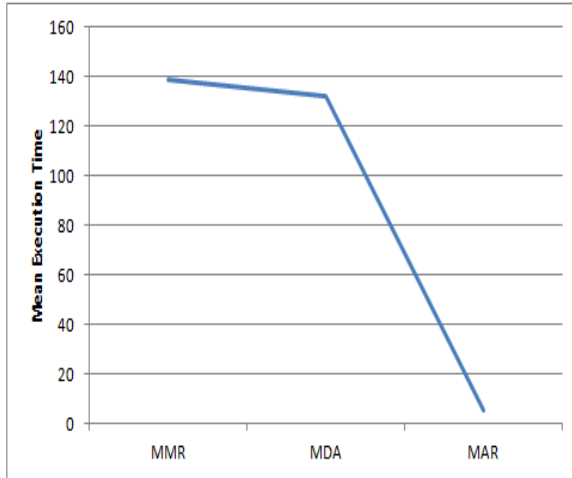


Fig. 3. Comparison of mean execution times using mushroom dataset

4.2 Modified Mushroom Dataset

Originally, mushroom dataset consist of 8124 instances with 22 attributes. Modified version of mushroom dataset acquired after removing an attribute with the missing values for this experiment and lefts only 21 attributes but number of records still the same. Results of the experiment show that MMR technique needs an approximately 139 seconds of execution times, while MDA technique requires 132 seconds of execution time. Using this dataset, MAR technique only needs 5.5 seconds to accomplish the task. As shows in Fig. 2, it is clear that MAR technique improves the execution time almost up to 95% when involving large dataset.

5 Conclusion

In this paper, a soft set based technique for partition attribute selection in multi-valued information systems has been proposed. This alternative technique, called Maximum Attribute Relativity, exploits the value of relativity relationship among categorical attributes value of an information system to determine a partition attribute. Using comparisons at the attribute level is firstly computed before the representative for each attributes is compared to select a partition attribute. Using this technique, execution time is slightly improves especially when involving large categorical data sets as compared to rough set-based techniques.

Acknowledgement. This paper is supported by Universiti Tun Hussein Onn Malaysia under the FRGS grant number 0758.

References

1. Pawlak, Z.: Rough sets. *International Journal of Computer and Information Science* 11, 341–356 (1982)
2. Pawlak, Z.: *Rough Sets - Theoretical Aspect of Reasoning about Data*. Kluwer Academic Publisher, Boston (1991)
3. Mazlack, L.J., He, A., Zhu, Y., Coppock, S.: A Rough Sets Approach in Choosing Partitioning Attributes. In: *Proceeding of ICSA 13th International Conference, CAINE 2000*, pp. 1–6 (2000)
4. Parmar, D., Wu, T., Blackhurst, J.: MMR: An Algorithm for Clustering Categorical Data using Rough Set Theory. *Data and Knowledge Discovery* 63, 879–893 (2007)
5. Herawan, T., Deris, M.M., Abawajy, J.H.: A Rough Set Approach for Selecting Clustering Attribute. *Knowledge Based System* 23, 220–231 (2010)
6. Molodtsov, D.: Soft set theory - First results. *Computer and Mathematics with Applications* 37, 19–31 (1999)
7. Maji, P.K., Biswas, R., Roy, A.R.: Fuzzy soft sets. *Journal of Fuzzy Mathematics* 9, 589–602 (2001)
8. Maji, P.K., Biswas, R., Roy, A.R.: Soft Set Theory. *Computer and Mathematics with Applications* 45, 555–562 (2003)
9. Ali, M.I., Feng, F., Liu, X., Min, W.K., Shabira, M.: On some new operation in soft sets theory. *Computer and Mathematics with Applications* 57, 1547–1553 (2009)
10. Maji, P.K., Roy, A.R., Biswas, R.: An application of soft sets in a decision making problem. *Computer and Mathematics with Applications* 44, 1077–1083 (2002)
11. Roy, A.R., Maji, P.K.: A fuzzy soft set theoretic approach to decision making problems. *Journal of Computational and Applied Mathematics* 203, 412–418 (2007)
12. Kong, Z., Gao, L., Wang, L.: Comment on “A fuzzy soft set theoretic approach to decision making problems”. *Journal of Computational and Applied Mathematics* 223, 540–542 (2009)
13. Herawan, T., Rose, A.N.M., Mat Deris, M.: Soft Set Theoretic Approach for Dimensionality Reduction. In: Ślęzak, D., Kim, T.-H., Zhang, Y., Ma, J., Chung, K.-I. (eds.) *DTA 2009. CCIS*, vol. 64, pp. 171–178. Springer, Heidelberg (2009)

Effect of Vocabulary Preparation on Students Vocabulary and Listening Comprehension

Lijun Li¹, Kaida He^{2,4}, and Qiudong He³

^{1,3}School of Foreign Languages, Hubei Engineering University, 272 Jiaotong Road,
Xiaogan City, China

²Graduate School of Translation and Interpretation, Beijing Foreign Studies University,
No.2 West Third Ring North Road, Haidian District, Beijing City, China

⁴College of Foreign Languages and Literature, Wuhan University, Luojiashan Wuhan City,
China

heqiudong@163.com

Abstract. This paper investigates and analyzes the effects of vocabulary preparation on students' vocabulary and listening comprehension in local engineering universities. There are three groups in this study. The three groups were given three different preparation times to learn new vocabulary that would be heard in a listening text. The students' performances of vocabulary and listening comprehension were also tested and analyzed statistically. The instruments include a vocabulary test, a listening comprehension test, and a questionnaire to investigate students' confidence levels. The results show that significant differences between groups were found only in the vocabulary test but not in the listening comprehension test. It is concluded that allowing students to prepare vocabulary before a listening test could improve their vocabulary knowledge and confidence level but not their listening comprehension. The students' listening comprehension performance mainly depends on their level of listening comprehension and listening strategies rather than vocabulary preparation. Another finding is that the group with 30-min preparation time showed the higher levels of confidence. According to the results, some strategies were proposed to improve the listening comprehension of students in local engineering universities.

Keywords: vocabulary preparation, vocabulary, listening comprehension, confidence, local engineering universities.

1 Introduction

It is not easy for the students' to listen and understand English completely in local engineering universities, on one hand because their first language plays an important role in most of their communication; on the other hand learning how to listen is largely through formal instruction in the classroom and with limited exposure to English outside formal study. The students are usually confronted with a number of difficulties, such as limited vocabulary, unfamiliar topic, fast speech, unfamiliar dialects, and listening only once and so on. Among these factors, lack of vocabulary

knowledge seems to be the most difficulty to EFL learners[1-2]. So, lots of students have become used to prepare vocabulary before a listening test.

Pre-task activities have been demonstrated to be helpful to fluency and complexity in the performance of oral narratives[3-4]and with writing skills[5]. Three major types of pre-task activities were proposed by Skehan[6]: teaching, consciousness-raising, and planning. Teaching is about the introduction of new language to the interlanguage. Consciousness-raising activities involve pre-task discussion and exposure to material relevant to the task. Planning involves the issue of time. Several studies suggest that planning time influences output to a great degree in terms of fluency, complexity and accuracy[3]. Whether it is very helpful to receptive skills, e.g. listening is what should be explored in this study.

2 Literature Review

Nowadays there is no research on the effect of vocabulary preparation on vocabulary performance and listening comprehension of the students in local engineering universities. However, there are a few studies focusing on pre-teaching vocabulary. Berne[7] looked at how vocabulary-preview and question-preview affected her Spanish learners. The results show that the question-preview group scored higher than the vocabulary-preview group, but the difference was not significant. Elkhafaifi[8] replicated Berne's study and the results were found comparable. It was apparent that studying vocabulary before listening did not have a significant facilitative effect on the students' comprehension, whereas additional exposure to the input seemed to be very helpful in enhancing listening comprehension.

Chang and Read[2] found that neither high nor low level learners seemed to have benefited from the vocabulary preparation they received immediately before the test. Chang[9] ever found that students had no time or very little time to practice the vocabulary before a test. Buck[10] notes "when second language learners learn some new element of a language, at first they have to pay conscious attention and think about it; that takes time and their use of it is slow." Based on the findings of these studies, this study aims to investigate whether L2 students will perform differently when vocabulary lists are given with different preparation times. The study will solve the following questions:

1. Does varying preparation time make a difference to the students' performance of vocabulary or listening comprehension?
2. Does different preparation time make a difference to the students' confidence?

3 Research Procedure

3.1 Research Participants

A total of 120 Hubei local university students, aging from 18 to 22, took part in the study. The participants were from three intact classes with English listening courses.

At the beginning of the new academic year they were all given a listening test, the test results revealed no significant differences across the three classes. In addition, these three classes have the same teacher and the same teaching methods. At the time the research was being conducted, student numbers divided into three Groups: (Group A, Group B, Group C).

3.2 Research Design

Give three groups three different times to prepare vocabulary. Group A was given a list of vocabulary a week before the listening comprehension test, and Group B, one day before. Group C was not given any materials for home, but they studied the vocabulary in classroom for 30min and then discussion was allowed. The participants were told that vocabulary on the list would occur in the listening text of the test. The research design is set out in Table 1.

Table 1. Research design

Experimental groups	Group A	Group B	Group C
Preparation time	1 week	1 day	30 min
No. of participants	40	40	40

3.3 The Listening Text

The listening text used in this study was a 790-word short story selected from Penguin Young Readers, level 4, written from a 1300-word wordlist. This story was chosen because it is not a popular one and it is very unlikely that the participants would know the story, so the background knowledge would not be a factor affecting the listening performance. The listening text was read by a native speaker of English.

3.4 Instruments

The instruments used in this study are a vocabulary test, a listening comprehension test, and a questionnaire. The details of each instrument are as follows:

1. A vocabulary test: Based on the vocabulary list, 25 words were selected to test whether participants had learned the words.
2. A listening comprehension test: The whole test contained three tests: multiple-choice, gap-filling, and short-answer questions. There are 8 items to test the students' listening comprehension of the story.
3. A post-test questionnaire: A post-test questionnaire containing 10 items was used to investigate the students' confidence level with different vocabulary preparation. Five-scale Likert was adopted for scoring.

3.5 Preparation Materials

The vocabulary list contained 25 words from the story; 7 nouns, 10 verbs, 3 proper nouns and 5 adjectives. To avoid variations about the difficulty of the words between the researcher and test-takers and between low and high level students, these words were chosen by 6 students who volunteered at the preliminary stage of the research, but they were not students in this study.

3.6 Research Procedure

Group A was given the vocabulary list a week before the test, and Group B, a day before the test. Group C was given the vocabulary list 30min before the test, and then the list was collected by the teacher. The participants were first given the vocabulary test, then the listening comprehension test. When the test was finished, they filled out the questionnaire.

3.7 Analysis

SPSS Statistics 17.0 was employed for the data analysis. A univariate analysis using a generalized linear model (GLM) was conducted to evaluate the effect of vocabulary preparation on the students' vocabulary knowledge, listening comprehension, confidence level. The scores obtained from the tests, the questionnaire are the dependent variables; the two independent variables are Groups A–C.

4 Results

4.1 Results of Test

Table 2 lists the test results and it shows that students scored higher on the vocabulary than the listening comprehension test. In the vocabulary test, Group A answered 76% correctly, higher than Group B with 65%, and Group C at 63%; however, the difference between Groups B and C was only two percent. In the test of listening comprehension, Group A scored highest, 59%, Group B, 56% and Group C, 54%. Although the differences among the three groups were minimal, the scores show the longer the preparation time, the higher the scores they achieved.

Table 2 shows the impact of vocabulary preparation on the students' performances of vocabulary and listening comprehension. A GLM univariate analysis of variance was performed to uncover what factors contributed to differences and whether the differences are significant. The result shows that the scores in Vocabulary are higher than those of Listening Comprehension. Table 3 is a summary of the GLM analysis on the effects of preparation time on the students' vocabulary performance. This table presents that different preparation times have significant main effects on the vocabulary performance of the groups, $p=.000<.05$.

Table 2. Results of vocabulary test and listening test

Experimental groups	Group A	Group B	Group C
Vocabulary Test	77.6	66.4	64.7
Listening Test	59.7	57.0	56.6

Table 3. Results of the effects of preparation time on the students' vocabulary

Source	SSIII	df	MS	F	Sig.	Effect size(eta squared)
Model	586435	42	13963	254.4	.000	.993
Group	2073	39	53	.97	.533	.326
Preparation Time	3897	2	1949	35.5	.000	.477
Error	4281	78	55			
Total	590716	120				

Table 4 presents a summary of the GLM analysis on the effects of preparation time on the students' listening comprehension. The independent variables were the same as in the previous section, but the dependent variable was the score of listening comprehension. This table shows that different preparation times have less significant effects on the listening comprehension, $P=.082>.05$. The results suggest that The students' listening comprehension relies more on their listening proficiency than on the amount of preparation time with vocabulary, and listening comprehension cannot be improved only by preparing vocabulary.

Table 4. Results of the effects of preparation time on listening comprehension

Source	SSIII	df	MS	F	Sig.	Effect size(eta squared)
Model	402304	42	9579	214.4	.000	.991
Group	1866	39	48	1.07	.390	.349
Preparation Time	231	2	115	2.58	.082	.062
Error	3484	78	45			
Total	405788	120				

4.2 Results of Questionnaire

This questionnaire is to find out students' confidence levels for different lengths of preparation time. Table 5 tells us that Group C had higher scores than Groups A or B. The data shows that the students with 30min preparation time have higher confidence level than the students with a week or a day of preparation time.

Table 5. Descriptive statistics of the questionnaire results

Experimental groups	Group A	Group B	Group C
No. of the participants	40	40	40
Confidence level	2.88	2.60	3.30

Table 6 presents the statistics of the students' listening confidence affected by preparation time. From this table, preparation times of vocabulary have significant effect on the listening confidence, $P=.037<.05$; While preparation times of vocabulary have less significant effect on the listening performances, $P=.996>.05$.

Table 6. Results of the effects of preparation time on listening confidence

Source	SSIII	df	MS	F	Sig.	Effect size(eta squared)
Model	1039	12	86.5	75.08	.000	.893
Preparation Time	7.86	2	3.9	3.41	.037	.059
Scores of listening Comprehension	1.89	9	.21	.182	.996	.015
Error	124.5	108	1.15			
Total	1163	120				

5 Suggestions for Teachers

5.1 Providing Time for Preparation

This study shows that vocabulary preparation had some impact on the students' confidence. From the questionnaire, a lot of the students think that studying vocabulary was very helpful to their comprehension and many of them were able to use the given words to predict some content. In addition to vocabulary preparation, there are many other forms of preparation, such as including visual aids or allowing multiple listening [2]. Language teachers may choose some preparations which fit for their students' language level. However, it is very important to allow sufficient time for the students to prepare in a classroom listening test, providing students with sufficient time for preparation may help them familiarize the pronunciation, link the lexical items to the aural text. Underwood [11] argued it is unfair to plunge students straight into the listening text. So before listening, students should be given time so they know what to expect. We have to remember that our goal is to give the students the experience of success and to help them understand spoken English.

5.2 Adopting Oral Repetition Strategy

Knowing the pronunciation of a word includes being able to recognize the word when it is heard and also being able to produce the spoken form[12]. Many students said that they could not match the sound and the written form of a word. So doing listening practice, students should be encouraged to speak the words which are not familiar to them. Gu and Johnson[13] claim that the use of oral repetition strategy is positively correlated with general language proficiency, whereas visual repetition was the strongest negative predictor of general language proficiency. In this study, very few students in Groups A and B enunciate the vocabulary. This influenced the effect of listening comprehension through vocabulary. On the contrary, most students in Group C said that they had tried to articulate every word and discuss. So although Group C had the least time for preparing the vocabulary, they achieved a score comparable

with the other two groups. Gathercole and Baddeley^[14] think keeping a word in your phonological short-term memory is an important factor influencing vocabulary learning.

5.3 Encouraging Discussion in Class

In listening classroom, students often feel anxious, even frightened because of fearing they can't understand what heard. So it is very important to care about student's affection. Wang Chuming[15] claimed affection is the booster to remain learning. In the class teachers should give students more opportunities to discuss to remove their anxieties, and to enhance their confidence and the level of listening performance. This study also shows that providing more time to study can't guarantee better comprehension or significantly higher levels of confidence. But encouraging discussion in classroom can produce as good effect as giving students a large amount time to prepare for a test. Dansereau[16] proposed discussion in classroom and cooperative learning can promote active processing of information and enhance motivation in the participants. Therefore, encouraging discussion in classroom and encouraging group study may be a good approach to improve listening comprehension.

6 Conclusion

This study examines the effect of vocabulary preparation time on the students' performances of vocabulary and listening comprehension, and how this preparation affected the students' confidence. Vocabulary preparation can be quite useful in a classroom test because it enhances the students' confidence and increases their willingness to complete the task. But the results of this study show that providing students with the vocabulary of an aural text did not greatly enhance the students' listening comprehension, because the students' listening comprehension cannot be improved only by familiarizing the vocabulary, but mainly by the students' listening strategy use and confidence. So in listening teaching, teachers should try to develop students' abilities to use listening strategies, and try to improve students' confidence.

Acknowledgments. The editors would like to acknowledge the assistance provided by Professor Zhang Weiyu in Central China Normal University in China. Without his splendid help, this paper would not have been possible.

References

1. Boyle, J.P.: Factors affecting listening comprehension. *English Language Teaching Journal* 5, 34–38 (1984)
2. Chang: Read: The effects of listening support on the listening performance of EFL learners. *TESOL Quarterly* 4, 375–397 (2006)

3. Ellis, R.: Inter-language variability in narrative discourse: style shifting and use of the past tense. *Studies in Second Language Acquisition* 9, 19–20 (1987)
4. Crookes, G.: Planning and inter-language variation. *Studies in Second Language Acquisition* 1, 367–383 (1989)
5. Kroll, B.: What does time buy? ESL student performance on home versus class compositions. Cambridge University Press, New York (1990)
6. Skehan, P.: *A Cognitive Approach to Language Learning*. Oxford University Press, Oxford (1998)
7. Berne, J.E.: How does varying pre-listening activities affect second language listening comprehension? *Hispania* 3, 316–329 (1995)
8. Elkhafaifi, H.: The effect of pre-listening activities on listening comprehension in Arabic learners. *Foreign Language Annals* 3, 505–513 (2005)
9. Chang, C.: The effect of lexical support in EFL listening comprehension tests. In: *The Proceedings of the 22nd International Conference on English Teaching and Learning in the Republic of China*, Crane, Taipei, Taiwan (2005)
10. Buck, G.: *Assessing Listening*. Cambridge University Press, Cambridge (2001)
11. Underwood, M.: *Teaching Listening*. Longman, London (1989)
12. Nation, I.S.P.: *Range and Frequency: Computer Programme for Windows based PCs*. University of Wellington, Victoria (2001)
13. Gu, Y., Johnson, R.: Vocabulary learning strategies and language learning outcomes. *Language Learning* 4, 643–679 (1996)
14. Gathercole, S., Baddeley, A.D.: Valuation of the role of phonological STM in the development of vocabulary in children: a longitudinal study. *Journal of Memory and Language* 8, 200–213 (1989)
15. Wang, C.: Two major factors influencing L2 learning and their effect on L2 teaching. *Foreign Language World* 6, 8–9 (2001)
16. Dansereau, D.F.: *Cooperative learning strategies*, pp. 103–120. Academic Press, New York (1988)

3D Parametric Design on Trough Type Liquid Distributor Based on AutoCAD VBA

Pengfei Zhang and Luojia Wan

Chemical Engineering Institute, Tianjin University, Tianjin, China
zhangpf@tju.edu.cn, xufeidiewan@yahoo.cn

Abstract. To reduce the duplicated work and enhance the drawing efficiency in the tower internals design process, this paper introduces a 3D parametric design method by taking the example of trough type liquid distributor. This new method utilizes the inlaid development language Visual Basic of Applications (VBA) to realize the secondary development of AutoCAD, and then establishes a parametric design program. By changing parameters, this program can automatically generate the 3D models of trough type liquid distributor with different sizes. These 3D models are more intuitive and convenient than the common used 2D graphics. At present, there is a great demand for 3D parametric design system in modern industries and design enterprises.

Keywords: VBA, AutoCAD, 3D parametric design, trough type liquid distributor.

1 Introduction

With the rapid development of computer technology and the emergence of the corresponding software, Computer-Aided Design (CAD) technology has been widely used. Currently, the most popular CAD software is AutoCAD which is developed and sold by American Autodesk Corporation. AutoCAD provides an ideal re-development platform as well as the powerful drawing function for the designers, which makes both 2D and 3D parametric design possible. In order to reduce duplicated work in the drawing process and at the same time enhance the drawing quality and efficiency [1], the re-development language must be utilized to develop AutoCAD.

In AutoCAD model space drawing environment, size-driven principles are used to realize the size changes automatically correspond with the geometry bound changes [2]. This will do a big favor to the 3D parametric design when objects with the same structure but with different sizes are requested to be designed. If the 3D parametric design software package of the object is tailored by taking advantage of the second development language, a 3D drawing will be automatically generated after parameter inputting. With the benefits such as lower product development costs and a greatly shortened design cycle, this 3D parametric design method will be widely used in modern industries and design enterprises [3].

This article takes trough type liquid distributor as an example, illustrates in detail how to use the VBA language to develop the automatic 3D drawing procedure.

2 The Technical Approaching of Parametric Design

What follows is a foundation for understanding the 3D parametric design based on AutoCAD VBA. Here, we use AutoCAD2008 as medium or drawing support software and the Visual Basic of Applications (VBA) as the second development language [4].

2.1 Overview of VBA Technology Built into AutoCAD

Although AutoCAD graphic processing ability is very strong, but when we solve some quite specialized questions, especially data processing, synthesize design calculation and so on, the effect of directly using the interactive function in AutoCAD is poor. Therefore, it's particularly important to use a second development language to develop AutoCAD and simplify the 3D model design in engineering. At present, on second development languages there are Visual Lisp, AutoLISP, Visual Basic, ARX, VBA and so on. VBA is an object-oriented programming language, and it has well inherited the easy-to-use features of Visual Basic (VB). What's more, VBA has the newest development technology and the formidable development function [5]. It is resident in main application procedure, the code operating efficiency of VBA is quite high [6].

VBA is an implementation of programming language VB 6.0 and its associated Integrated Development Environment (IDE). It's built into some applications such as AutoCAD. When a new project is added, VBA IDE can be used to edit the form, program code and references of the project as well as compile or run the project.

AutoCAD ActiveX Automation technology of VBA provides users the programming mechanism from inside or outside of the AutoCAD [7], and it has been also designed to have a capacity of transferring data to other windows applications [8]. This enables users to operate another windows application program (such as Microsoft Word, Excel, Access, Visio and so on) to call on the VBA program conveniently through ActiveX technology (Figure 1).

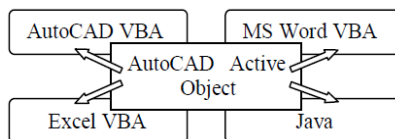


Fig. 1. AutoCAD ActiveX Automation

To realize 3D parametric design, there are three key factors: AutoCAD, ActiveX Automation and VBA IDE. The designer must have a capability of using these three tools proficiently.

2.2 The Object Model of AutoCAD ActiveX

Through AutoCAD ActiveX Automation, lots of functions of AutoCAD are encapsulated in ActiveX objects described by the AutoCAD Object Model with their

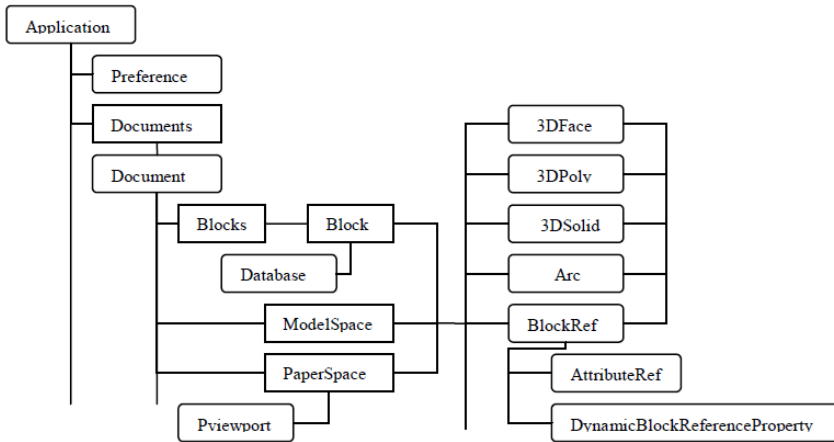


Fig. 2. The AutoCAD ActiveX Object Model

properties, methods and events. All the ActiveX objects can be exposed in one way for further programming, and they are composed by hierarchies. The AutoCAD Object Model is very large so that only a small part is shown here (figure 2).

3 Program Design Process of Trough Type Liquid Distributor Based on AutoCAD VBA

In chemical engineering, tower internals design is a very important constituent of the process package design. If chemical engineers establish the 3D parametric drawing program of the tower internals, it will be greatly helpful for workers to understand the assembly relationship and the structure of the tower internals. Liquid distributor is a kind of packing tower internals. Besides the packing itself, a good liquid distributor plays a greatly important role in a well operated packing tower as well [9], it can efficiently improve the liquid initial distribution and the mass transfer efficiency. The following illustration will show the procedures of 3D parametric design of the trough type liquid distributor.

3.1 Create a New Project and Open the VBA IDE

The steps of creating a new project is: firstly, open AutoCAD2008 and choose “tools”; then click “Macros” and VBA Manger Dialog Box will appear; finally, click “new” to create a project. If the project has been created, we just need to load it by click “load” in the VBA Manger Dialog Box.

There are three ways can be used to open the VBA IDE. First, type “VBAIDE” in the command line. Second, choose “tools” and click “Macros” and then choose Visual Basic Editor. Third, press “Alt + F11”.

3.2 Insert a UserForm and the Related Controls

To add a UserForm to the project, we just need to choose “insert” and then choose “UserForm”. The related controls are added via Toolbox, here we insert two Frame named as Input and Output. In the “Input” frame, we insert five Label named as The 3D Trough Type Liquid Distributor Design Program, Column Diameter, Manhole Diameter, Distribution Trough Width and Distance Between Distribution Trough, then insert three Text Box and a Combo Box named as txtbox1, txtbox2, txtbox3 and trough distance. In the “Output” frame, we insert four Label named as Trough NO., Long side length, Short side length and Hole NO., then insert four Text Box named as txtbox4, txtbox5, txtbox6 and txtbox7. At last insert one Label named as The 3D Trough Type Liquid Distributor Design Program and two Command Button named as Exit and Draw (figure 3).

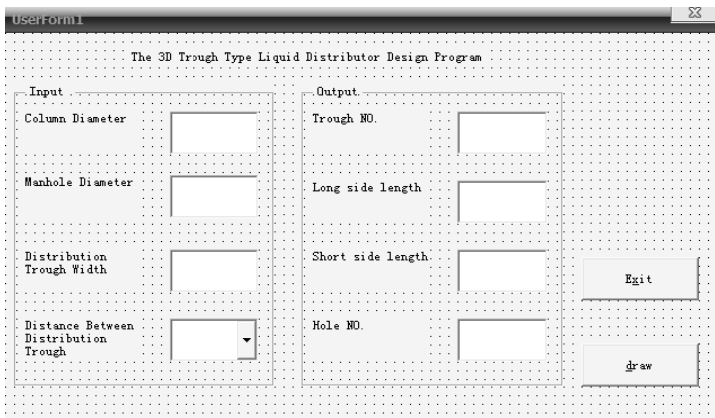


Fig. 3. The 3D trough type liquid distribution design program

This UserForm is a kind of dialog box, which is created for inputting parameters. The user can input the different parameters for different requirements. When the parameters are obtained by the program, the program itself will analyze and process the data, and then all the results of data process together with the macro will be transferred to AutoCAD by ActiveX Automation technology.

3.3 Add Codes for UserForm

Because the code of the whole project is too long, here we select part of the program to illustrate. Double-click the UserForm to open the code window and enter the following codes.

Create User Coordinate System in the Current Graphics.

```

Public Sub NewUCS()
    Dim ucsObj As AcadUCS
    Dim origin(0 To 2) As Double
    Dim xAxisPnt(0 To 2) As Double
    Dim yAxisPnt(0 To 2) As Double
    'Define user coordinate system
    origin(0) = 500#: origin(1) = 500#: origin(2) =
500#
    xAxisPnt(0) = 600#: xAxisPnt(1) = 500#: xAxisPnt(2) =
500#
    yAxisPnt(0) = 500#: yAxisPnt(1) = 600#: yAxisPnt(2) =
500#
    'Add "UCS" to UserCoordinatesSystem set
    Set ucsObj =
ThisDrawing.UserCoordinateSystems.Add(origin, xAxisPnt,
yAxisPnt, "New_UCS")
    MsgBox ucsObj.Name & " System has added." & vbCrLf &
"base point: " & ucsObj.origin(0) & ", " &
ucsObj.origin(1) & ", " & ucsObj.origin(2)
    'Display UCS icon
    ThisDrawing.ActiveViewport.UCSIconAtOrigin = True
    'Set the current "UCS"
    ThisDrawing.ActiveUCS = ucsObj
End Sub

```

Some Code of Creating Trough Type Liquid Distributor

```

Sub create_model()
    .....
    Dim polyline3(0) As Acad3DPolyline
    Dim path As AcadLine
    Dim p3(0 To 32) As Double
    p3(0) = 43.5+txtbox3: ..... : p3(32)=0#
    Set polyline3(0) =
ThisDrawing.ModelSpace.Add3DPoly(p3)
    polyline3(0).Closed = True
    Dim region4obj As Variant
    region4obj =
ThisDrawing.ModelSpace.AddRegion(polyline3)
    Dim solid7obj As Acad3DSolid
    Dim path7(0 To 2) As Double
    Dim path8(0 To 2) As Double
    path7(0) = 0#: path7(1) = 0#: path7(2) = 0#
    path8(0) = 0#: path8(1) = -(txtbox1 - 203#):
    path8(2) = 0#
    Set path = ThisDrawing.ModelSpace.AddLine(path7,
path8)
    Set solid7obj =
ThisDrawing.ModelSpace.AddExtrudedSolidAlongPath(region4o
bj(0), path)
    .....
End Sub

```

3.4 Execute the Program

If the parametric design program is compiled with no errors, execute the program, we will obtain the whole 3D solid model of the trough type liquid distributor (figure 4)

Based on the generation of the 3D solid model, some post processing can be done to express the structure of the trough type liquid distributor. Here, two viewports are created in the AutoCAD model space drawing environment, one for the top view, another for the front view and detail drawings (figure 5).

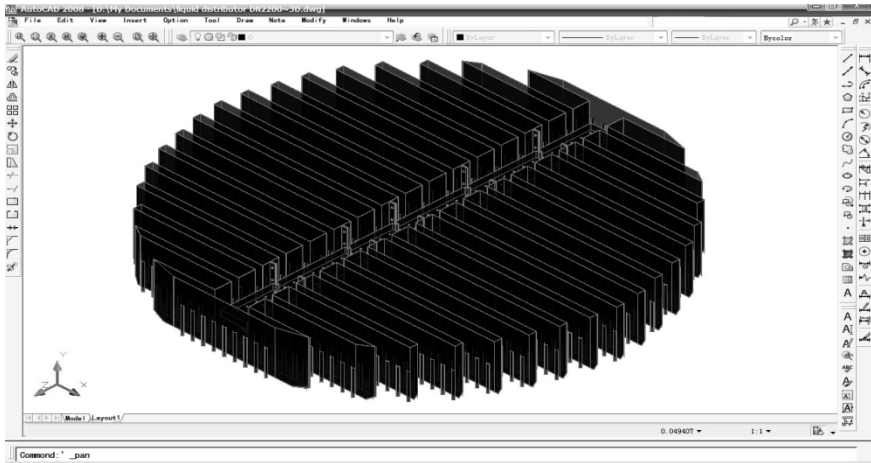


Fig. 4. The 3D model of trough type liquid distributor

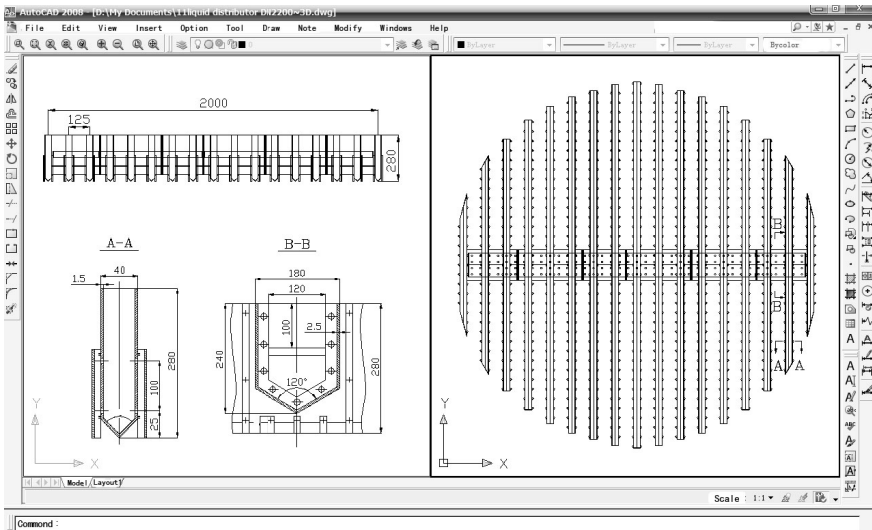


Fig. 5. The top view, front view and detail drawings of trough type liquid distributor

4 Conclusion

The secondary development of AutoCAD by taking advantage of VBA makes 3D parametric design possible and provides a design platform which is intelligent, convenient and intuitive [10]. In this way, many duplicated work in the drawing process can be avoided and at the same time the drawing quality and the efficiency can be enhanced. As a result of afore-mentioned advantages, 3D parametric design should be widely used in industry and design enterprises.

Acknowledgments. Thanks to Prof. hang for his guidance and encouragement.

References

1. Kilani, M.I.: Computer Aided Design Tools in the Development of Surface Micromachined Mechanisms. *Jordan Journal of Mechanical and Industrial Engineering* 5(2), 167–176 (2011)
2. Leng, Y.F., Liu, J., Zu, F.: Parametric Design on Belt Conveyor Drums Based on VBA. *Advanced Materials Research* 156-157, 1243–1246 (2011)
3. Inozemtsev, A.N., Troitsky, D.I., Bannatyne, M.W.M.: Parametric Modeling: Concept and Implementation. In: 4th International Conference on Information Visualization, p. 504. IEEE Press, New York (2000)
4. Wang, R.H.: Development of Parametric Drawing Program Based on AutoCAD VBA. In: International Conference on Computer Aspects of Social Networks, pp. 757–760. IEEE Press, New York (2010)
5. Wang, Y.: Develops AutoCAD 2000 Application Procedure with VBA. People's Posts and telecommunications Press, China (1999)
6. Ye, Y.N.: AutoCAD 2000 ActiveX and VBA Reference Manual. The Chinese Electric Power Press, China (2001)
7. Zhao, Y., Yu, Q.B., Wang, K.: 3D Parametric Design System for Heating Furnace. In: 2nd International Conference on Computer Engineering and Technology, pp. 173–178. IEEE Press, New York (2010)
8. Wong, K.W.W., Barford, J.P.: Teaching Excel VBA as a Problem Solving Tool for Chemical Engineering Core Courses. *Education for Chemical Engineers* 5, 72–77 (2010)
9. Chen, B., Zhou, S.P., Qiu, L.B., Peng, C.: On CAD Trough Liquid Distributor of Packing Tower. *Technology Supervision in Petroleum Industry* 22(8), 58–60 (2006)
10. Qin, S.F., Prieto, P.A., Wright, D.K.: A Novel Form Design and CAD Modeling Approach. *Computers in Industry* 59, 364–369 (2008)

A Novel Differential Evolution Algorithm with Adaptive of Population Topology

Yu Sun^{1,2,*}, Yuanxiang Li¹, Gang Liu¹, and Jun Liu²

¹ State Key Lab of Software Engineering, Computer School, Wuhan University,
Wuhan 430072, P.R. China

² School of Computer and Electronics and Information, Guangxi University, Nanning 530004,
P.R. China

sunyu1225@163.com, yxli@whu.edu.cn, lg0061408@126.com,
liujunzky@163.com

Abstract. Differential evolution is a simple and efficient algorithm. Although it is well known that the population structure has an important influence on the behavior of EAs, there are a few works studying its effect in DE algorithms. In this paper, a novel adaptive population topology differential evolution algorithm (APTDE) is proposed for the unconstrained global optimization problem. The topologies adaptation automatically updates the population topology to appropriate topology to avoid premature convergence. This method utilizes the information of the population effectively and improves search efficiency. The set of 15 benchmark functions provided by CEC2005 is employed for experimental verification. Experimental results indicate that APTDE is effective and efficient. Results show that APTDE is better than, or at least comparable to, other DE algorithms.

Keywords: Differential Evolution, Self-Adaptive, Population Topology, Global Optimization.

1 Introduction

Differential evolution (DE)[1] is a vector population based algorithm which has been successfully applied in diverse fields such as the design of digital filters [2, 3], pattern recognition [4, 5], aerodynamic design[6], power systems[7, 8], and many others. The performance of the DE algorithm is sensitive to the mutation strategy and respective control parameters such as the population size (NP), crossover rate (CR) and the scale factor (F). Besides, it is well accepted in the literature [9, 10] that the organization of individuals in the population has a major influence on the DE performance.

Recently, some researches proposed to define some structures into the population to improve DE approaches. In the traditional DE, the population model is a set of individuals with no structure. Zaharie et al.[11] designed the decentralized population, which presents a parallel distributed self-adaptive DE algorithm. Kozlov et al. [12] proposed in a simple new migration scheme in which the oldest solution in the

* Corresponding author.

population is replaced by the received one. Shing et al.[13] presented a detailed study on the parameterization of a parallel distributed DE for solving the learning scheme in asymmetric subset hood product fuzzy neural inference system.

In all the above mentioned versions of DE, a single population topology is used. However, it is possible that a particular population topology may not be suitable for all types of problems [14]. To increase diversity and avoid the problems of premature convergence, this paper presents a new adaptive DE named APTDE. In the proposed APTDE algorithm, five DE's population topology strategies are selected as candidates, if the solution of the algorithm is not improved, other candidate topology will be used instead to preserve the diversity of population and to speed up the convergence rate.

2 Differential Evolution Algorithm (DE)

The original DE algorithm follows the general procedure of an evolution algorithm. As in other evolution algorithms, four main operations drive the evolution in DE: initialize, mutation, crossover and selection. The main operations are as follows:

(1) **Initial:** At first, DE randomly generated the Initial population $\{X_{i,0} = (X_{1,i,0}, X_{2,i,0}, \dots, X_{D,i,0}) | i = 1, 2, \dots, NP\}$ according to a distribution $X_j^{min} \leq X_{j,i,0} \leq X_j^{max}$, for $j = 1, 2, \dots, D$, where NP is the population size and D is the dimension of the problem. Then, DE enters a loop of operations: mutation, crossover, and selection.

(2) **Mutation Operation:** Based on target vector $X_{i,G}$ at generation G, a mutation vector $V_{i,G+1}$ is created:

$$V_{i,G+1} = X_{r3,G} + F * (X_{r1,G} - X_{r2,G}), \quad r_1 \neq r_2 \neq r_3 \neq i \tag{1}$$

Where indexes r_1, r_2, r_3 are randomly chosen in the range $[1, NP]$. F is a factor in $[0, 2]$ for scaling differential vectors .

(3) **Crossover Operation:** To increase diversity, after mutation, a mutant vector $V_{i,G+1}$ and the target vector $X_{i,G}$ are mixed to generate the trial vector:

$$U_{i,G+1} = (u_{1,i,G+1}, u_{2,i,G+1}, \dots, u_{D,i,G+1}) \tag{2}$$

Where

$$u_{j,i,G+1} = \begin{cases} v_{j,i,G+1}, & \text{if } (rand_j[0,1] \leq CR) \text{ or } (j = j_{rand}) \\ x_{j,i,G}, & \text{otherwise} \end{cases}, j = 1, 2, \dots, D \tag{3}$$

$rand_j$ is random chosen in the range $[0, 1]$. CR is a user-specified crossover constant random chosen in the range $[0, 1)$, and j_{rand} is a randomly chosen index in the range $[1, D]$, which ensures that the trial vector $U_{i,G+1}$ gets at least one parameter from the target vector $V_{i,G+1}$.

(4) **Selection Operation:** To decide whether the trial vector $U_{i,G+1}$ become a member of generation G+1, the $U_{i,G+1}$ is compared to the target vector $V_{i,G+1}$ using a greedy scheme. The scheme is:

$$X_{i,G+1} = \begin{cases} U_{i,G+1}, & \text{if } (f(U_{i,G}) \leq f(X_{i,G})) \\ X_{i,G}, & \text{otherwise} \end{cases} \quad (4)$$

$f(U_{i,G})$ is the fitness value of the trial vector, and $f(X_{i,G})$ is the fitness value of the target vector. Compared $f(U_{i,G})$ with $f(X_{i,G})$, if $f(U_{i,G})$ is smaller than $f(X_{i,G})$, the trial vector will replace the target vector, and enter the population of the next generation. Otherwise, the target vector will remain for the next generation.

3 Adaptive Population Topology DE (APTDE)

3.1 DE with Population Topologies

Various types of population topologies in DE and their influences have been investigated in the literature [14]. In the evolution process, individuals can only interact with its neighboring individuals, which are defined on the given topology. Therefore, different topologies guide to different parent individuals selection scheme. Some common population topology structures are discussed in this section:

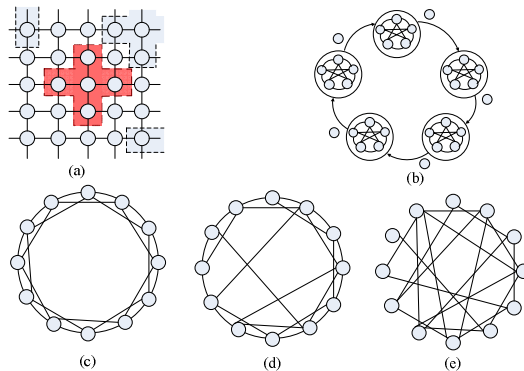


Fig. 1. (a) cellular DE (b) distributed DE (c) regular DE (d) small-world DE(e)random DE

(1) **Cellular DE:** In a cellular DE (cDE), the population is divided into a large number of very small subpopulations typically composed of only one individual, which are arranged in toroidal mesh. In the cellular topology an individual may only interact with its nearby neighbor subpopulations in the breeding loop [5].

(2) **Distributed DE:** In the case of distributed DE (dDE), the population is partitioned in a set of subpopulations arranged in a unidirectional ring topology. Every subpopulation can send information to only one other subpopulation, and in the same way it can only receive information from one single subpopulation.

(3) Regular lattice DE, Small-world DE and Random graph DE: In the regular lattice topology DE(regularDE), which is a regular non-random ring lattice, each individual is connected to its K immediate neighbors only. The random graph topology, in randDE is a theoretical construct which the links between individuals that are chosen completely at random with equal probability. The small-world topology, which generated in swDE starts with a one-dimensional regular lattice topology, which has connection between each individual and its k nearest neighbors, then rewires each connection at random with probability p , where $0 < p < 1$.

3.2 Adaptive Population Topologies in Differential Evolution Algorithm

DE suffers from the problem of premature convergence, where the population converges to some local optima of a multimodal objective function, losing its diversity. The population topology plays an essential role in the original DE algorithm. Some population topologies were brought into DE algorithms to implement in the paper[14]. To improve the performance on a specific problem by using the DE algorithm, we need to try all available topologies to fine-tune the corresponding one. And it may require a huge amount of computation time. Therefore, we attempt to develop a new DE algorithm that can automatically adapt the topologies during evolution.

The idea behind our proposed adaptation strategy is to probabilistically select suitable one out of several topologies and apply to the current population. In this Adaptive Population Topologies in Differential Evolution Algorithm (APTDE), the population is generated according to the population size NP and the dimension of the problem D . Besides, five population topologies described in Section 3.1 were selected as the candidates: cellular topology, distributed topology, regular lattice topology, small-world topology and random graph topology. And the five candidate topologies are assumed that is of the same probability to be selected. The operations of APTDE are described:

(1) Initialization

At the first, the population is structured in the specific topology, selected randomly from the five candidate topologies, and the neighbor list of each individual is also defined on the specific topology. Then, APTDE enters a loop of four operations: mutation operation, crossover operation, selection operation and updating operation.

(2) Mutation operation

The individual can only interact with its neighboring individuals. Therefore, the parents $(X_{r2,G}, X_{r3,G})$ are in this case chosen among the neighbors to generate the mutant vector based on the equation (1).

(3) Crossover operation

After mutation, a mutant vector generated by APTDE and the target vector $X_{r1,G}$ are mixed to generate the trial vector based on the equation (2).

(4) Selection operation

Then, both the mutant vector and the trial vector (the new solution) are generated, the fitness value for the latter is computed, and it is inserted instead of the current solution in the population following a given replacement policy (4).

(5) Updating operation

At the end of each generation, the best individual is recorded. Comparing the best individual among last three generations, if the best individual is improved, the previous topology would be remained; otherwise another candidate topology from the candidates would be applied in the next generation. If the new topology was applied, a new neighbor list would be generated. And then parent vectors would be selected from the new neighbor list to generate the new mutation vector.

The above operations, including mutation operation, crossover operation, selection operation and updating operation, are repeated generation after generation until some specific stopping criteria are satisfied.

4 Experiment Settings and Results

In order to evaluate the performance of the proposed APTDE algorithm discussed in Section 3.2 is compared to the DE algorithms in different topologies, namely distributed DE, cellular DE, random DE, small-world DE and regular DE, described in Section 3.1. The set of benchmark functions is proposed in IEEE CEC 2005. For all algorithms, the initial population size NP is set to 100. Besides, we have considered dimensions $D = 50$ for all the problems, and $F=0.5$ and $CR=0.9$ are used for the all DEs and no other parameters is adjusted. In addition, the number of evolution generations of all algorithms is set to 1500, and, each algorithm is run 50 times.

In the evolution process, the parents are selected randomly from their neighbor lists defined by specific topology. The topology for the distributed DE (dDE) is partitioned into four subpopulations which evolved by 25 independent individuals. The communication between the subpopulations is done every generation for dDE. In the case of the cellular DE(cDE), population is arranged in a mesh of 10×10 . Besides, individuals in the regularDE connected with their eight neighbors in a ring topology. And at least 3 connection of an individual are randomly chosen, in random DE. Moreover the population in small-world DE starts as a regular topology, and rewire 2 connections at random. In the case of adaptive population topology DE (APTDE), the population is randomly arranged on a topology mentioned in the section 3.1, and adaptive change to another if it is necessary. For all algorithms mentioned above, the connected individuals are stored in the neighbor lists of each individual.

4.1 Results and Discussions

We compare the proposed APTDE algorithms versus the other DEs with different topologies in this section. The results of 50 independent runs, including the mean and the standard deviation values, are summarized in Table 1. When the optimal solution was found by the algorithms, the percentage of runs in which it happened is indicated

between parenthesis. In Table 1, figures in bold font represent the best results we found. The background color of the cells is related to the two-tailed t-test performed to compare APTDE with the other DEs. Those figures in Tables 2 with light gray background are significantly worse than APTDE according to statistical test (with 95% confidence), while the dark gray background stands for better results. We also compare the algorithms in terms of efficiency, i.e., the number of evaluations they require to find the optimal solution. Fig. 2 shows some examples of the convergence of the best solution in the population for some test functions.

4.1.1 Comparison of the Computational Results

As the results are shown in Table 1, the test detects a different behavior between APTDE and other DEs for most of the problems. We proceed now to analyze the behavior of the same algorithms for the big problem instances. The results are provided in Table 2. We can see that the best algorithms for these problems are APTDE, better than all the other ones according to our statistical test for 6 problems. Comparing to dDE, APTDE is significantly better than dDE on 8 out of 15 functions. In the case of cDE, APTDE is significantly better than cDE on 12 out of 15 functions. Comparing to regularDE, APTDE is significantly better on 6 out of 15 functions. In the case of swDE, APTDE is significantly better on 11 out of 15 functions. At last, APTDE is significantly better than randDE on 12 out of 15 functions.

Table 1. Comparison of the Results of APTDE with other DEs (50 Variables Problems)

	APTDE	dDE	cDE	regularDE	swDE	randDE
F1	1.46E-13	1.02E-16	1.05E-09	4.20E-08	9.36E-08	2.10E-06
F2	3.48E+02	3.51E+02	1.40E+03	1.66E+03	4.07E+03	5.73E+03
F3	2.80E+06	2.92E+06	5.43E+06	4.96E+06	1.75E+07	2.25E+07
F4	2.03E+03	1.06E+03	5.95E+03	3.39E+03	1.08E+04	1.44E+04
F5	3.94E+03	6.28E+03	4.97E+03	6.48E+03	4.13E+03	5.02E+03
F6	1.08E+02	2.06E+02	2.61E+02	6.14E+02	6.87E+01	6.72E+01
F7	6.20E+03	6.20E+03	6.21E+03	6.20E+03	6.28E+03	6.37E+03
F8	2.12E+01	2.13E+01	2.12E+01	2.12E+01	2.12E+01	2.12E+01
F9	7.18E+01	1.17E+02	2.68E+02	4.84E+01	3.56E+02	3.92E+02
F10	1.37E+02	1.89E+02	3.94E+02	1.41E+02	3.90E+02	4.02E+02
F11	5.35E+01	6.21E+01	7.41E+01	6.27E+01	7.44E+01	7.45E+01
F12	2.71E+04	2.95E+04	3.80E+04	6.14E+04	2.62E+04	1.64E+05
F13	8.02E+00	1.33E+01	3.09E+01	6.40E+00	3.19E+01	3.32E+01
F14	2.29E+01	2.32E+01	2.32E+01	2.30E+01	2.33E+01	2.33E+01
F15	3.84E+02	3.85E+02	3.77E+02	3.80E+02	3.80E+02	3.92E+02

4.1.2 Comparison of the Convergence of the Best Solution

We also compare the algorithms in terms of efficiency, i.e., in terms of the number of evaluations they require to find the optimal solution. Fig. 2 shows some examples of the convergence of the best solution in the population for the different algorithms for problems F1, F2, F10, and F15. The values plotted for every generation are averaged over 1500 independent runs.

Comparing the behavior of the six algorithms, we can see that APTDE in most cases obtains the best performing, while in some problems the dDE algorithm converges faster to better solution. The swDE algorithm converges slightly slower than the dDE algorithm and obtains a better performance than the cDE algorithm. The randDE algorithm and regularDE algorithm are worst in the four problems.

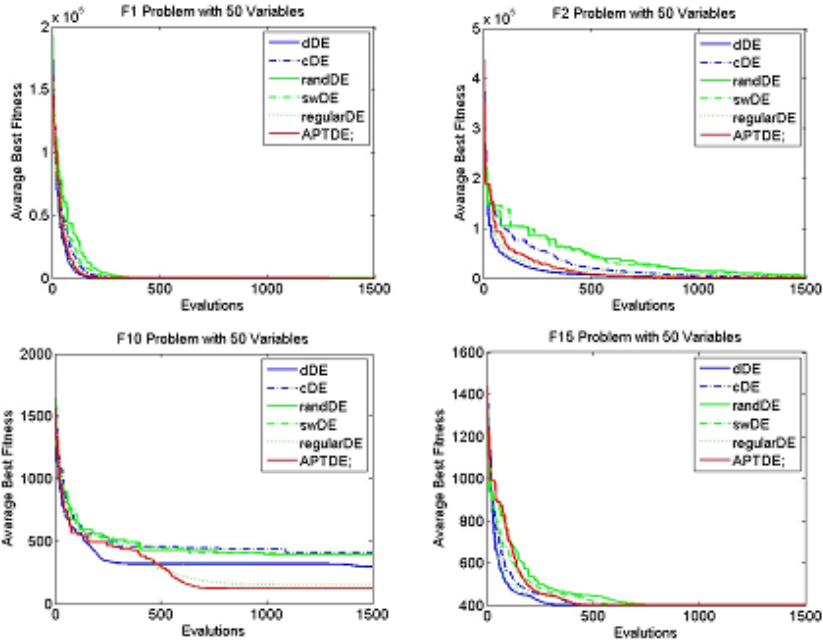


Fig. 2. Evolution of the best solution for the algorithms

5 Conclusions and Future Works

Even though DE is a good and fast algorithm, it has shown some weaknesses, especially long computational times because of its stochastic nature. This drawback sometimes limits its application to optimization problems. In this paper, an adaptive population topology DE algorithm for optimizing non-differential functions over continuous spaces is proposed. Topology adaptation is beneficial to improve the optimization performance of an evolutionary algorithm by automatically updating topologies to appropriate values during the evolutionary search. 25 benchmark functions on real-parameter optimization are used to evaluate the performance of our presented approach. Comparisons APTDE with other DEs, it demonstrates that the proposed algorithm is more effective and efficient in terms of the solution quality and the number of function evaluations in most cases. In future work, the effect of the population topology and the self-adaptation strategy will be studied in more detail. Future work also includes APTDE will be used to solve the constrained optimization problems and design support vector regression.

Acknowledgement. This work is supported by the National Nature Science Foundation of China (Grant No.61070009).

References

1. Storn, R., Price, K.: Differential Evolution – A Simple and Efficient Heuristic for global Optimization over Continuous Spaces. *Journal of Global Optimization* 11(4), 341–359 (1997)
2. Liu, G., et al.: Design of two-dimensional IIR digital filters by using a clustering-based differential evolution with chaotic sequences. *International Journal of Digital Content Technology and its Applications* 5(9), 153–163 (2011)
3. Karaboga, N., Cetinkaya, B.: Design of Digital FIR Filters Using Differential Evolution Algorithm. *Circuits, Systems, and Signal Processing* 25(5), 649–660 (2006)
4. Das, S., Konar, A.: Automatic image pixel clustering with an improved differential evolution. *Applied Soft Computing* 9(1), 226–236 (2009)
5. Maulik, U., Saha, I.: Modified differential evolution based fuzzy clustering for pixel classification in remote sensing imagery. *Pattern Recognition* 42(9), 2135–2149 (2009)
6. Rogalsky, T., Derksen, R.: Hybridization of differential evolution for aerodynamic design (2000)
7. Zhang, X., et al.: Dynamic multi-group self-adaptive differential evolution algorithm for reactive power optimization. *International Journal of Electrical Power Energy Systems* 32(5), 351–357 (2010)
8. Varadarajan, M., Swarup, K.S.: Differential evolution approach for optimal reactive power dispatch. *Applied Soft Computing* 8(4), 1549–1561 (2008)
9. Liu, G., et al.: A novel clustering-based differential evolution with 2 multi-parent crossovers for global optimization. *Applied Soft Computing* 12(2), 663–681 (2012)
10. Liu, G., et al.: Improving clustering-based differential evolution with chaotic sequences and new mutation operator. *International Journal of Advancements in Computing Technology* 3(6), 276–286 (2011)
11. Zaharie, D., Petcu, D.: Parallel implementation of multi-population differential evolution. In: *Concurrent Information Processing and Computing*, pp. 223–232 (2003)
12. Kozlov, K., Samsonov, A.: New migration scheme for parallel differential evolution. In: *Proc. 5th Int. Conf. Bioinformatics Genome Regulation Structure* (2006)
13. Singh, L., Kumar, S.: Parallel Evolutionary Asymmetric Subsethood Product Fuzzy-Neural Inference System: An Island Model Approach (2007)
14. Dorransoro, B., Bouvry, P.: Improving Classical and Decentralized Differential Evolution with New Mutation Operator and Population Topologies. *IEEE Transactions on Evolutionary Computation* 15(1), 67–98 (2011)

Empirical Study on the Relationship between Money Supply and Stock Market in Europe

Yijun Li

School of Economics, Wuhan University of Technology, Wuhan, China
Liyijun0410@163.com

Abstract. To investigate whether the European Central Bank's frequent funding action can strengthen the stock market or not, the actual relationship between money supply and stock market in Europe deserves research, especially during the debt crisis. This article first presented theoretical analysis on the channels through which money supply influences stock market. The study was then verified by empirical analysis based on ADF unit root test and Johansen cointegration test. By creating cointegration model and Vector Error Correction Model and carrying on Granger causality test, further examinations revealed the interaction between money supply and stock market capitalization in the short and long term. The results suggest that stock market capitalization is conversely related to money supply in the long run, whereas money supply has positive impact on stock market capitalization in the short-term, but it's not the Granger reason of stock market capitalization. The economic significance were summarized in the last part, then the specific proposals were put forward.

Keywords: European Debt Crisis, Money Supply, Stock Market Capitalization, Vector Error Correction Model.

1 Introduction

As the global financial crisis from the financial sector spread to the real economy, European countries implemented an active fiscal policy and moderately loose monetary policy in 2009, which eases the recession in the short term. But at the same time, a catastrophe for the wider economy was on the horizon. In December 2009, the European debt crisis broke out, and then European stock markets collapsed, even global stock market. To rescue Greece, Ireland, Italy and other debtor countries, the ECB issued loans to these countries and purchased sovereign bonds from them, causing the fund flowing into asset markets in a big way [1-4].

According to traditional theory, increasing money supply can improve safe assets ratio. In order to accumulate wealth, the public are pleasant to own venture assets with a higher yield, then stock prices rise. That would make the market value of enterprises higher than their replacement costs, which attracts more investment into the stock market and upgrades the flow of the market value of the whole stock market. So it would be interesting to determine whether European stock markets fully capture relevant data on money supply during the European debt crisis [5].

The purpose of the present paper is to contribute further to the literature on the relationship between money supply and stock market and, specifically, for the countries suffering from debt crisis. At the particular background of European debt crisis, this study aims to provide some clue to settle practical problems [6-8].

2 Methodology

The vector error correction approach is the main line of this paper. A vector error correction (VEC) model is a restricted VAR, which is commonly used for forecasting systems of interrelated time series and for analyzing the dynamic impact of random disturbances on the system of variables [9]. Because the VEC models are designed for use with nonstationary series that are known to be cointegrated, the sample data are analyzed in the following order: test for stationary of each variable and cointegration relationship between them, then establish VEC models, inspect causality in the end.

3 Empirical Study

3.1 Variable Settings and Data Sources

Monthly data for money supply and stock-market capitalization are obtained for the whole of the euro area. Choose these variables for several reasons: First, the broad money (M2) is an important implement of monetary authorities to adjust monetary policy, which have a closer relationship with the stock market than money supply figures at other levels. Therefore, this paper selected the broad money (M2) as the independent variable. Second, given various stock price indexes of the euro area and their different calculation methods, it is difficult to select a certain index which can reflect the performance of the stock market all around the euro area. Consequently, we choose the stock-market capitalization of the whole euro area (SC) as dependent variable, which is an objective reflection of the stock price [10-11].

The data comes from the European Union's statistics office website (epp.eurostat.ec.europa.eu). The period investigated for this study is January 2009 to October 2011 for euro area. In order to eliminate the interference of heteroscedasticity, we took the logarithm of original data getting LNM2 and LNSC which were more precise and of the lower order of magnitude [12].

3.2 Unit Root Tests

Time series are often nonstationary, since they include a clear time trend. Therefore, it is important to check whether a series is stationary or not before using it in a regression. Among various testing strategies, this paper firstly tests for "stationarity" of each variable by employing ADF test [13]. The ADF test with long lags is superior to the others. We chose to estimate the ADF test that employs automatic lag length selection using a Schwarz Information Criterion (SIC).

Table 1. Testing the Variables for Unit Roots

series	ADF test	Critical value		
		1%	5%	10%
LNM2	-3.098676	-4.284580	-3.562882	-3.215267
LNSC	-1.997581	-3.646342	-2.954021	-2.615817
D(LNM2)	-5.446190	-4.323979	-3.580623	-3.225334
D(LNSC)	-5.134751	-3.653730	-2.957110	-2.617434

D () refers to first order differential.

It can be seen from Table 1 that, for money supply (LNM2) series and stock-market capitalization (LNSC) series, the null hypothesis of a unit root is not rejected at the 1% significance level. To verify that the order of integration is I (1), the presence of a unit root in the first difference of the stock price indices was also tested but no unit roots in first differenced series was found. However, as the differential data cannot be used to analyze the long-term effects of the variables, we will have to rely on cointegration analysis to make up for this limitation.

3.3 Johansen Cointegration Test

Engle and Granger (1987) pointed out that a linear combination of two or more non-stationary series may be stationary. If such a stationary linear combination exists, the non-stationary time series are said to be cointegrated. The stationary linear combination is called the cointegrating equation and may be interpreted as a long-run equilibrium relationship among the variables. We choose the Johansen methodology based on the likelihood ratio with non-standard asymptotic distributions involving integrals of Brownian motions, because it is found to be the best method to proceed with cointegration estimation. According to the AIC, SC criteria, LR value, and other statistics the optimal lag orders of this test is determined.

Table 2. Tests for the Number of Cointegrating Vectors

Hypothesized No. of CE(s)	Eigenvalue	Trace Statistic	0.05 Critical Value	Prob.
None *	0.421463	28.15256	25.87211	0.0256
At most 1	0.282881	10.64045	12.51798	0.1009

means the original hypothesis is rejected at the significance level of 5%.

In Table 2, the results of cointegration tests on money supply (LNM2) and stock-market capitalization (LNSC) indicate that at the significance level of 5% the null hypothesis that the variables are not cointegrated ($r = 0$) is rejected, that is, there is a cointegration vector among these variables. Using the Johansen Cointegration Vector Estimation Approach, we can obtain the long-term equation as follows:

$$LNSZ = -40.77835LNM2 + 0.085732 @ TREND \quad (1)$$

capitalization is conversely related to the money supply, as the money supply reduces by 40.77835 units when the European stock market capitalization increases by one percentage point. The coefficient of @TREND indicates the time trend has a minor positive effect on the European stock market. While the presence of a cointegrating relation forms the basis of the VEC specification, it suggests the likelihood that money supply influences stock-market capitalization in the long run, and The interrelation between them should be verified further.

3.4 Vector Error Correction Model

A vector error correction (VEC) model is a restricted VAR designed for use with nonstationary series that are known to be cointegrated. The VEC has cointegration relations built into the specification so that it restricts the long-run behavior of the endogenous variables to converge to their cointegrating relationships while allowing for short-run adjustment dynamics. The cointegration term is known as the error correction term since the deviation from long-run equilibrium is corrected gradually through a series of partial short-run adjustments. According to the results of cointegration test, AIC, SC criteria and t, R2 value, and other statistics to determine the optimal lag orders. Thus, this paper established a group of two variable VEC model. By analyzing the output of the Eviews6.0, the only significant model is as the following equation:

$$\begin{aligned}
 D(LNM\ 2) = & -0.023165Co\ int\ Eq1 + 0.022666D(LNSZ(-1)) \\
 & [-4.33242^{**}][1.94778^*] \\
 & +0.143610D(LNM\ 2(-1)) + 0.001323C \\
 & [0.83607][2.06902^*]
 \end{aligned}
 \tag{2}$$

Where the numbers in [] are t statistic values.*and** refer to significance at the level of 5% and of 1% respectively.

As can be seen from the above equation, most of the factors through t test. From the AIC and SC statistics perspective, the whole effect of the model is pretty good, since the output of Eviews 6.0 shows that the values of AIC (-11.33231) and SC (-10.82846) are small. Judging from the coefficient, the coefficient of D(LNSZ(-1)) reaches 0.022666, which shows a slightly positive correlation between money supply and stock market capitalization in the short term which is absolutely opposed to the long-term harmonious relationship . The coefficient of error correction shows the cointegration relationship between money supply and stock-market capitalization has a negative effect on the growth of money supply.

3.5 Granger Causality Test

The Granger (1969) approach to the question of whether x causes y is to see how much of the current y can be explained by past values of y and then to see whether adding lagged values of x can improve the explanation. y is said to be Granger caused

by x if x helps in the prediction of y , or equivalently if the coefficients on the lagged x 's are statistically significant. Note that two-way causation is frequently the case; x Granger causes y and y Granger causes x . The analysis above provides the evidence that the first difference data are not integrated, non-stationary processes and further integration and cointegration analysis thus ceases at this point, we employed Granger's concept of causality to test the relationship between money supply and stock-market capitalization.

Table 3. Results of Granger-Causality Test Statistics from the VECM

Original hypothesis	P-value
D(LNM2) does not Granger cause D(LNSZ)	0.8525
D(LNSC) does not Granger cause D(LNM2)	0.0514

Table 3 shows stock market capitalization is a Granger cause of money supply at the 10% significance level, while money supply is not a Granger cause of stock market capitalization. Therefore, money supply does not have obvious degree of explaining stock market capitalization, but it does not mean that money supply has no explanation for stock market trends. The results can only reveal that money supply in such short-term interdynamic relation is liable to get into a passive position, and a more aggressive monetary policy may weaken the adverse impact.

4 Conclusion

In this thesis, money supply and stock market capitalization of euro area during the European debt crisis are used as sample data, financial measures such as ADF unit root test, Johansen cointegration test, Vector Error Correction Model (VECM) and Granger causality test are applied to study the interrelationship between money supply and stock market capitalization. From the results of this study, we get the following conclusions:

First, a standard Johansen Cointegration test for the existence of a long-run relationship between money supply and stock market capitalization for the whole euro area. Unlike traditional theory, when money supply rises, stock market capitalization drops. That is to say, over the period of January 2009 to October 2011, their negative relationship has been established.

Second, the vector error correction (VEC) model reveals that there is short-term positive equilibrium relationship between European money supply and stock market capitalization, but this influence is slight.

Besides, the results of Granger causality tests after the error correction show that there is not enough evidence to justify the predictive value of European money supply to stock market capitalization in the short term. With this just the opposite is, stock market capitalization is a Granger cause of money supply.

Above all, pure support increasing money supply does little to boost the struggling stock market and wretched economy of euro area in the long run, while the recovery and stability of the stock market show signs of a short-term strengthening economy,

which justify the funding of activities of ECB. This implies that the long-term development strategy policy should be taken to help put the economy on a firm recovery path. On the other hand, effective measures to stabilize the stock market run will enhance the confidence of investors in a short period.

References

1. Dickey, D.A., Fuller, W.A.: Likelihood ratio statistics for autoregressive time series with a unit root. *Econometrica* 49, 1057–1072 (1981)
2. Engle, R.F., Granger, C.W.J.: Cointegration and error correction: representation, estimation and testing. *Econometrica* 55, 251–276 (1987)
3. Granger, C.W.J.: Development in the study of cointegrated economic variables. *Oxford Bulletin of Economics and Statistics* 48, 213–228 (1986)
4. Johansen, S.: Statistical analysis of cointegrating factors. *Journal of Economic Dynamics and Control* 12, 231–254 (1988)
5. Johansen, S., Juselius, K.: Maximum likelihood estimation and inference on cointegration with applications to the demand for money. *Oxford Bulletin of Economics and Statistics* 52, 169–210 (1990)
6. Zhu, D.: An empirical analysis on the effect of the stock market on money demand in China. *Transaction of Shanghai Finance and Economics University* 6, 3–10 (2002) (in Chinese)
7. Sun, H., Ma, Y.: The relationship between monetary policies and the stock market in China. *Economic Research* 7, 22–27 (2003) (in Chinese)
8. Mak, W.M., Vonk, W., Schriefers, H.: Animacy in Processing Relative Clauses. *Journal of Memory and Language* 54, 466–490 (2006)
9. Ehrmann, M., Fratzscher, M.: Taking Stock: Monetary Policy Transmission to Equity Markets. *Journal of Money, Credit and Banking* 36(4), 719–737 (2004)
10. Ioannidis, C., Kontonikas, A.: Monetary Policy and the Stock Market: Some International evidence. Department of Economics, University of Glasgow (2006)
11. Qiao, G., Liu, Z.: The Empirical Analysis of the relation between money supply and the stock market. *Economic Management* 24, 12–14 (2007) (in Chinese)
12. Chen, S.: Does Monetary Policy Have Asymmetric Effects on Stock Returns? *Journal of Money Credit and Banking* 39, 667–688 (2007)
13. Ai, F.: Chinese monetary number and name of the national income relationship of empirical study. *Journal of Collective Economy, Academic Study* 8, 32–34 (2008) (in Chinese)

Evaluation of Agricultural Information Service System

Yanxia Wang and Xingjie Hui

Northeastern University at Qin Huangdao, Hebei, China 066004
wyx62256@sina.com

Abstract. Recent years, constructions of Agricultural Information Service System has completed its first stage, fundamental equipments of the system are significantly promoted and the social organization required by information service begins to take shape. But, studies about these specific fields are deficiency especially in the evaluation criteria of Agricultural Information Service System. According to the theory of system engineering and researches on structure, elements and mechanism, an index system involved the AHP method to calculate weights of index is introduced. This system will contribute to help government define key points of management and check out obstructed stages of different areas, which will compel the constructions of Agricultural Information Service System during the “12th Five Year Plan”.

Keywords: Agricultural Information, Agricultural Information Service System, System Evaluation, Evaluation Index System.

1 Introduction

As the basis of agricultural information transmission, constructions of agricultural information service system around China are developing quickly and have got some improvements. However, shortages appear while the system processing. Imbalance of the developments between different areas, deficiency of agricultural information system, comes with the unsatisfied quality and obstructed transmission of the products of agricultural information. According the mechanism of interactions among structures and elements in one system, an index system about the evaluation of Agricultural Information Service System, which is not only involved the AHP method to calculate weights of index, but also built on the analysis of classical samples of China and the basic principles of information system, is introduced to accelerate the construction of Agricultural Information Service System and provide scientific criteria for the programs and managements of government. The index system could make the main task obvious and promote the efficiency of the service system.

1.1 Concept of Agricultural Information Service System

Agricultural Information Service System is the high efficient system including transmission and orderly flow of information between the produce areas and information collection, which uses the Internet as the main carrier, along with modern

remote communication methods, diverse media and terminals collecting, transmitting, processing and developing valuable material.

1.2 Structure of Agricultural Information Service System

According to the principle of the processing of the Agricultural Information Service System (Source-Currency-Use), it could be considered as three distinctive subsystems: the production of information, transmission and application.

Production: With countless information, selection and collection of valuable ones form the potentiometer between resource and destination. Therefore, the result of collection could increase the quality and the amount of product of effective information, but also clearly show the direction of information flows.

Propagation, which means the process of the selected information's transmission and spreading from agents to users, is the critical ring of the link connecting the products and information. The potential value of information will not convert into practical effect when the selected messages are only in the agent's database. Until they are pushed to the users and combined with the process of producing, messages are useless as nothing but words.

Application: After users gain information products through diverse carriers, elements and factors, which will raise up the level of agricultural economy and farmers income, do occur according to what the comprehension of the message will be and how to combine them with other factors of produce. The effect of Agricultural information equals to the procedure of mixing production and management, agricultural activities and other produce conditions together after the arrival of information.

1.3 Goals of Agricultural Information Service System

Among plenty of benefits building such a system, three of them are on the top: First, an Agricultural Information Service Organization could be founded at the same moment, which can provide abundant human resource for construction and management of the system. Second, the structure of the system should be completed and well functioned. The flow of the Agricultural Information Product should follow the 'S-C-U' (Source-Currency-Use). Third, the system should function well and effective. Under that situation, the isolated Agricultural Information Fields could be combined with the produce of Agricultural products that will form an orderly information flows among producers, customers and markets around the world. Government could use the system to optimize the resource of agriculture and provide technical support to the development of Agricultural Economy.

2 Principles and Definition of Weights of the Indexes

2.1 Principles of Constructing an Evaluation Index System

Selecting and constructing an evaluation index system is the base of the comprehensive evaluation of system. To make the index system reflect main characters and hierarchy of

agricultural information service system, we need ascertain principia and aims of designing index system. Hereby, we should make correlative factors of system more methodic and make the evaluation index system scientific, sound and applied.

Constructing principia: systemic, purposefulness, applicability, independence, maneuverability.

Evaluation purpose: First, offering impersonal gist for system regulation and supervision; in the second, making system aims legible and concrete by evaluating; finally, forming competition mechanism by evaluating.

2.2 Modeling Process of Working Out Index Weights by AHP

AHP which is the same with multi-objective decision-making analysis is an analyzing method with both qualitative indexes and quantitative indexes. The main characteristic of the method, quantifying these experiences of decision-makers, is the same with the case with a complex structure of evaluation target and incomplete index data.

2.2.1 Main Steps of Modeling

Index system of the paper includes four layers: target layer (layer O)--subtarget layer (layer A)--rule layer(layer B)--index layer(layer X).Tiptop layer is evaluation target including only one element; middle layers are subtarget layers which include some layers; the layer above the nethermost layer is rule layer, the nethermost layer is index layer or scheme layer.

2.2.2 Constructing Judgment Matrix

To reflect proportions of subtarget layer in target measurement, and some factors of rule layer and index layer with respect to belonging targets. Pairwise comparison matrices which is judgment matrices are constructed by comparing factor.

Comparing effect of n genes $X=\{x_1, x_2, \dots, x_n\}$ on factor Z, x_i with x_j are compared every time. All comparative results of $a_{ij}=x_i/x_j$ can be shown matrix $A=(a_{ij})_{n \times n}$, so A is a relative matrix among Z-X which is judgment matrix.

Numbers 1-9 and their reciprocals are used to mark to expediently estimate importance difference between two factors.

2.2.3 Testing Consistency

Judgment matrix is comparative mark of the importance ratio of factors in matrix with respect to the same target, and index:

$$\text{Viz. } a_{ij} = \frac{a_{ik}}{b_{jk}} \quad (i, j, k=1, 2, \dots, n) \quad (1)$$

However, when we pairwise compare every factor of these complex things, we can not ensure consistency between every two factors of matrix, so we must test consistency.

If $C.R=C.I/R.I<0.1$ (R.I: average stochastic coincidence indicator, C.I: coincidence indicator), judgment matrix has accredited consistency. Otherwise, we need evaluate and amend judgment matrix again.

Expressions of C.I:

$$C.I = \frac{\lambda_{\max} - n}{n - 1} \tag{2}$$

To relax consistency requirement of multidimensional matrix, R.I is introduced to revise consistency index. In the paper, counting R.I, matrices rank($n=1\sim 11$)and samples(100~500) are needed.

Table 1. Average and stochastic coincidence indicator R.I value

matrix dimension	1	2	3	4	5	6	7	8	9
R.I	0.00	0.00	0.58	0.96	1.12	1.24	1.32	1.41	1.45

2.2.4 Figuring Out Weight Vector

Latent roots of matrix which satisfies consistency condition are worked out, most latent roots meets: $\lambda_{\max} = n$, its only one nonnegative eigenvector is simplified, the weight vector $W=(w_1, w_2, \dots, w_n)^T$ of every index is obtained. hierarchy single sequence is carried out.

2.2.5 Calculating Combination Weight-Hierarchy General Sequence and Consistency Test

Hierarchy general sequence is the relative importance quotient of every element of every layer with respect to some factor of superior layer.

The membership relation structure of evaluation index system for agricultural information service system is unattached and belongs to completely unattached hierarchy model. Establishing hierarchy general sequence, combination weight in lower layer is only the weight combination of the only super stratum factor.

The combination weight sum of all index layer factors is 1 after establishing hierarchy general sequence. According to the meaning of hierarchy general sequence, combination weights of target layer factors are sorted, the importance of some factor with respect to general target is obtained and those main factors which influence importantly general target are also obtained.

3 The Frame of Evaluation Index System

By seriously screening and carefully distinguishing, the index system frame of the comprehensive evaluation for agricultural information service system is finally established. Evaluation targets include 5 stair indexes, 17 second-level indexes and 59 third-level indexes. We can refer to these indexes' relations in combination weight and index taxis chart of comprehensive evaluation index for agricultural information service system (Table 2).

Table 2. Weight and sort of agricultural information service system index

Target layer	Subtarget layer		Rule layer		Index layer	
	Index name	Weight and sort	Index name	Weight and sort	Index name	Weight and sort
Comprehensive evaluation of agricultural information service system	Strategic status	0.0562(4)	Attention of government	0.0562(6)	Project programming of information	0.0040(36)
					Leader's capability of coordination	0.0157(17)
					Grossinvestment of system construction	0.0365(8)
					Information source of institutions	0.0009(55)
					Information source of people	0.0003(58)
					Information source of entity	0.0007(56)
					Information source of sci-technology	0.0055(30)
					Information source of documentation	0.0028(42)
					Designing index	0.0156(18)
					Coverage of collection points	0.0073(26)
	Production subsystem	0.2394(2)	Information collecting	0.0277(9)	The advance of collector	0.0033(39)
					Collection cycle	0.0015(48)
					The advance of processing technology	0.0254(12)
					The standardization level of processing	0.0284(11)
					Depth and quality of processing	0.0984(2)
					Database resource per ten thousand	0.0029(41)
					Database resource per ten thousand	0.0145(20)
					Website number per ten thousand	0.0030(40)
					Database accessing information members above county	0.0058(29)
					per ten thousand	0.0035(37-38)
Propagation subsystem	0.5521(1)	Propagation organization system	0.0483(7)	Information members of basic unit	0.0135(21)	
				per ten thousand		
				The ratio of basic-level service stations	0.0314(10)	
				The ratio of TV agricultural program	0.1462(1)	
				The ratio of Information service hall in county or town	0.0888(3)	
		resource		The ratio of radio agricultural	0.0501(5)	

Table 2. (continued)

			program	
			The ratio of rural library and their utilization situation	0.0345(9)
			Propagating service facilities and utilization	0.0155(19)
			The scale and number of science-technology demonstration zone	0.0109(23)
			Tele-voice service	0.0068(28)
			The development of interpersonal propagation	0.0052(32)
			The development of organizational propagation	0.0212(13)
			Number of innovation in mode of propagation	0.0022(45)
			Average computer number of service center above county	0.0049(33)
			Average computer number of rural service center	0.0170(14)
			The abundance of website contents and updating speed	0.0614(4)
			The ratio of network security	0.0390(7)
			Television number of per hundred households	0.0166(15)
			The usage rate of CATV	0.0010(51-52)
			Telephone number of per hundred households	0.0072(27)
			Computer number of per hundred households	0.0471(6)
			Farmers enthusiasm of training and effect	0.0091(24)
			Famers enthusiasm of information counseling	0.0024(43)
			Average schooling years	0.0010(51-52)
			Assistant decision-making load	0.0014(49-50)
			Assistant decision-making level	0.0124(22)
			Electronic commerce of agriculture	0.0053(31)
			Economic benefit	0.0035(37-38)
			Ecological benefit	0.0006(57)
			Social benefit	0.0014(49-50)
App ly- effe ct subs yste m	0.1 121 (3)	Innovati on in mode of propagati on	0.02 86 (8)	
		Network propagati on	0.12 33 (3)	
		Informati on- receiving facility	0.07 54 (4)	
		Informati on quality of farmers	0.01 25 (12)	
		Decision - supportin g	0.01 91 (11)	
		Yield energy	0.00 55 (15- 16)	

Table 2. Weight and sort of agricultural information service system index

Target layer	Subtarget layer		Rule layer		Index layer		
	Index name	Weight and sort	Index name	Weight and sort	Index name	Weight and sort	
System environment	System environment	0.0402(5)	The basic communication environment	0.0220(10)	Communications infrastructure	0.0041(35)	
			Per capita broadband quantity		0.0018(46)		
			Radio cover ratio		0.0161(16)		
			Fiscal subsidy and reward		0.0093(14)	Fiscal subsidy and reward	0.0077(25)
			Preferential loan of financial system			0.0016(47)	
			Legal environment		0.0055(15-16)	Data standardization	0.0009(53-55)
			Information security law			0.0046(34)	
			Agricultural management		0.0034(17)	Scale management	0.0002(59)
			The development of agricultural cooperation			0.0022(44-45)	
			Industrial management		0.0009(53-55)		

4 Result Analysis and Management Advices

Considering generally of the Agricultural Information Service System, the transmission is the most critical part which important index is 55.21%, much more than other factors. Therefore, we should concentrate on the construction of transmission systems as it directly relate with success of the whole system. According to the principle layer, first 3 of 17 are the integration of Propagation resource, Information processing and Network propagation (1 of production subsystem, 2 of propagation subsystem), the sum of these is over 60%.

Among of top ten of 59 indexes in index layer, there are 7 indexes in propagation subsystem. The integration of propagation resource contains four of them, and network propagation has another two. Therefore integrating transmission resources and transmission through network should be breaches where finance and management should focus on.

Agricultural Information Service System Management should carry out from following aspects: First, increase the total investment of constructions of system by system innovation; Second, considering the diversity of distinctive areas' circumstance and the demanding of farmers, promote the depth and accuracy of process of collection

of agricultural information; Third, integrating the resources of transmission, make different transmission methods interchange coverage and information affordable and convenient to be gained by farmers; Forth, government should compel computer networks construction in countries, reinforce management and monitor of the agricultural information website to promote the reliability and authority of the service system, so that farmers can use information correctly and get benefits with it.

References

1. Wang, Y.: Study on Construction and Evaluation of Chinese Agricultural Information Service System. China Agriculture Press, Beijing (2008)
2. Wang, Y., Zheng, H., Wang, J.: Elements and Mechanism of in the Agricultural Information Apply-effect System. Chinese Rural Economy 3, 29–32 (2006)
3. Li, X.: Study on Construction of Anhui Agricultural Information Service System. Journal of Anhui Agri. Sci. 39(3), 1852–1853 (2011)
4. Zhang, X., Wang, S., Yang, C., Wang, L.: Evaluation Research of Agricultural Informatization Development Level in Hebei Province. Economy and Management 1, 88–91 (2012)
5. Lu, Y., et al.: Overall Design of Agricultural Information Service Platform in New Rural. Journal of Anhui Agri. Sci. 39(19), 11914–11917 (2011)
6. Zhao, R.-Q., et al.: Analysis on the Supply and Demand of Agriculture Information Dissemination Based on Mass Communication Theory. Journal of Anhui Agri. Sci. 39(15), 9398–9400 (2011)
7. Wang, S., Tong, Z.: Measure of Agriculture Informatization Level and Development Trend in China. Research of Agricultural Modernization 3, 216–218 (2008)
8. Wang, Y.: Discussion on Causes of “Last Kilometer” and Countermeasures. Journal of Northeastern University (Social Science) 3, 180–182 (2005)
9. Lu, T., Li, X., Gao, X., Zhao, Y.: An Analysis for Constructing Urban Modern Agriculture Information Service System. Journal of Agriculture, <http://www.caaaj.org>
10. Du, H.: Study on Construction of Agricultural Information Service System in Jiansu Province. Journal of Beijing Agricultural Vocation College 1, 15–18 (2011)

Six-Mode Truncation and Chaotic Characteristics of Atmospheric Convection System

Li Zhen

Department of Basic Teaching, Tangshan College, Tangshan, 063000, China
heblizhen@163.com

Abstract. To provide a mathematical description of the chaotic behavior in atmospheric convection system, by expanding the variables in these equations in double Fourier series and selecting six basic modes in truncation, a six-modes truncation Lorenz-like equations are obtained. The basic dynamical behaviors and chaotic behaviors of the equations are simulated numerically according to control parameter changes in a certain range and the system shows a route to low-dimensional chaos through a sequence of bifurcations. This is helpful to understand turbulent flow. The data shown in numerical simulation provides a theoretical reference for the further study of turbulent flow and help chaotic characteristics of turbulent flow toward the desired direction.

Keywords: Atmospheric convection system, Truncation, Bifurcation, Chaos.

1 Introduction

Chaos theory in physics, mathematics is more complete system and has penetrated to the other disciplines, especially been widely applied in meteorology, seismology and other fields. People have achieved some research outcomes in the intersecting research area of chaos and turbulent flow [1,2]. So it is reasonable for people to use the chaos concept to explain turbulent flow. Chaos dynamics is used to study turbulent flow and this can reveal effectively the inherent law and provide a new vision for the research of turbulent flow [3]. Navier-Stokes equations of solving turbulent flow are very complex partial differential equations. It is impossible to derive the chaotic characteristics of turbulent flow in pure mathematical way. At present, truncation method is used to reduce infinite dimensional partial differential equation to finite dimensional ordinary differential equations and most main features of original system are retained [4]. Literatures [5-7] relate to this method. The appropriate mode is selected in Specific model to truncate the original equation and the problem is simplified. In the study of atmospheric turbulence, Navier-Stokes equation and heat conduction equation are expanded by Fourier series and Lorenz selected three modes in truncation. This led to the three-dimensional Lorenz system [8]. Six basic modes are selected in this paper. And a six-mode Lorenz-like system which is different from Lorenz system is obtained by introducing non-dimensional time. The new system has similar properties with Lorenz system. The basic dynamical behaviors and chaotic behaviors of the equations are simulated numerically according

to control parameter changes in a certain range. This is helpful to understand turbulent flow.

2 Mathematical Model

Atmospheric convection in two dimensions is described [8, 9] by the following system of partial differential equations:

$$\frac{\partial}{\partial t} \nabla^2 \psi = -\frac{\partial(\psi, \nabla^2 \psi)}{\partial(x, z)} + \nu \nabla^4 \psi + \alpha g \frac{\partial \theta}{\partial x} \tag{1}$$

$$\frac{\partial \theta}{\partial t} = -\frac{\partial(\psi, \theta)}{\partial(x, z)} + \frac{\Delta T}{H} \frac{\partial \psi}{\partial x} + k \nabla^2 \theta \tag{2}$$

Boundary conditions: $\psi = 0, \nabla^2 \psi = 0, \theta = 0, z = 0, H$

In this system ψ is the stream function and θ is the potential temperature. The constants α, ν, k, g denote respectively the coefficient of thermal expansion, the kinematic viscosity, the thermal conductivity and the acceleration of gravity. H is the fluid layer thickness and ΔT is the temperature difference between the upper and lower surface of the fluid. ∇^2 is the Laplacian operator, $\nabla^2 \psi = \frac{\partial^2 \psi}{\partial x^2} + \frac{\partial^2 \psi}{\partial z^2}$,

$$\nabla^2 \theta = \frac{\partial^2 \theta}{\partial x^2} + \frac{\partial^2 \theta}{\partial z^2}, \nabla^4 \psi = \frac{\partial^2 \psi}{\partial x^2} + \frac{\partial^2 \psi}{\partial z^2} + 2 \frac{\partial^4 \psi}{\partial x^2 \partial z^2}, \frac{\partial(f, g)}{\partial(x, z)} = \begin{vmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial z} \\ \frac{\partial g}{\partial x} & \frac{\partial g}{\partial z} \end{vmatrix}.$$

In order to understand the profile of the nonlinear phenomenon of this system, we find a set of six Fourier modes which helps to truncate (1) and (2) into a six-dimensional chaotic ordinary differential equations system.

3 Mathematical Analysis

The present study focuses on a Fourier mode truncation involving a suitable small number of modes in order to show the profile of the chaos arising from (1) and (2) in a simple way. Variables ψ and θ in (1) and (2) each is developed into the two-dimensional form of Fourier series by Saltzman, l is the wavelength of direction x , $2h$ is the wavelength of direction z [9].

$$\psi^*(x^*, z^*, t^*) = \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} \psi_{\min}(m, n, t^*) \exp[2\pi i(m \frac{h}{l} x^* + \frac{n}{2} z^*)] \tag{3}$$

$$\theta^*(x^*, z^*, t^*) = \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} \theta_{\min}(m, n, t^*) \exp[2\pi i(m \frac{h}{l} x^* + \frac{n}{2} z^*)] \tag{4}$$

$\psi^*, \theta^*, x^*, z^*, t^*$ are respectively the dimensionless of ψ, θ, x, z, t .

Lorenz approximated the solution of equation (1) and (2) by the three Fourier modes[8]. On this basis, six basic modes are selected in this paper:

$$\begin{aligned} \psi = \frac{k(1+a^2)\sqrt{2}}{a} \{ [X_1(t)\sin(\frac{\pi ax}{H}) + X_2(t)\sin(\frac{2\pi ax}{H})] \sin(\frac{\pi z}{H}) \\ + X_3(t)\sin(\frac{\pi ax}{H})\sin(\frac{2\pi z}{H}) \}. \end{aligned} \quad (5)$$

$$\begin{aligned} \theta = \frac{R_c \Delta T}{\pi R_a} \{ [\sqrt{2}Y_1(t)\cos(\frac{\pi ax}{H}) + \sqrt{2}Y_2(t)\cos(\frac{2\pi ax}{H})] \sin(\frac{\pi z}{H}) \\ - Z(t)\sin(\frac{2\pi z}{H}) \}. \end{aligned} \quad (6)$$

where a is a parameter and

$$R_a = \frac{g\alpha H^3 \Delta T}{kv}, \quad R_c = \frac{\pi^4(1+a^2)^3}{a^2}.$$

$X_1(t), X_2(t), X_3(t), Y_1(t), Y_2(t), Z(t)$ are functions of time and they are unrelated with x, z and the coordinate system(5) and (6) are substituted into (1), (2).According to the principle of equal of the corresponding items, the following equations are obtained:

$$\begin{aligned} X_1'(t) &= -\nu \frac{\pi^2(a^2+1)}{H^2} X_1(t) + \nu \frac{\pi^2(a^2+1)}{H^2} Y_1(t) + \frac{9\sqrt{2}}{4} \frac{k\pi^2(a^2-1)}{H^2} X_2(t)X_3(t) \\ X_2'(t) &= -\nu X_2(t) \frac{\pi^2(1+4a^2)}{H^2} + \frac{2\nu\pi^2(1+a^2)^2}{H^2(1+4a^2)} Y_2(t) + \frac{9\sqrt{2}k\pi^2(1+a^2)}{4(1+4a^2)H^2} X_1(t)X_3(t) \\ X_3'(t) &= -\nu X_3(t) \frac{\pi^2(4+a^2)}{H^2} - \frac{9\sqrt{2}k\pi^2 a^2(1+a^2)}{4H^2(4+a^2)} X_1(t)X_2(t) \\ Y_1'(t) &= \frac{R_a}{R_c} \frac{k\pi^2(1+a^2)}{H^2} X_1(t) - \frac{k\pi^2(a^2+1)}{H^2} Y_1(t) - \frac{k\pi^2(1+a^2)}{H^2} X_1(t)Z(t) \\ &\quad + \frac{3\sqrt{2}}{4} \frac{k\pi^2(1+a^2)}{H^2} X_3(t)Y_2(t) \\ Y_2'(t) &= \frac{2k\pi^2(1+a^2)}{H^2} \frac{R_a}{R_c} X_2(t) - \frac{k\pi^2(1+4a^2)}{H^2} Y_2(t) - \frac{3\sqrt{2}}{4} \frac{k\pi^2(1+a^2)}{H^2} X_3(t)Y_1(t) \\ &\quad - \frac{2k\pi^2(1+a^2)}{H^2} X_2(t)Z(t) \end{aligned}$$

$$Z'(t) = -\frac{4k\pi^2}{H^2}Z(t) + \frac{k\pi^2(1+a^2)}{H^2}X_1(t)Y_1(t) + \frac{2k\pi^2(1+a^2)}{H^2}X_2(t)Y_2(t) \tag{7}$$

On this basis, introduce a dimensionless time $\tau = \frac{k\pi^2(1+a^2)}{H^2}t$. $r = \frac{R_a}{R_c}$ is the main control parameters of system and express the ratio of the driving factors and inhibiting convection factor. $\sigma = \frac{V}{k}$ is the Prandtl number, $b = \frac{4}{1+a^2}$ representatives the shape ratio related with the convection aspect ratio, σ, b are dimensionless constants. Equations (7) transform into:

$$\left. \begin{aligned} \dot{X}_1 &= \sigma Y_1 - \sigma X_1 + \frac{9\sqrt{2}}{4} \frac{k\pi^2(a^2-1)}{H^2} X_2 X_3 \\ \dot{X}_2 &= \frac{2\sigma(1+a^2)}{(1+4a^2)} Y_2 - \frac{\sigma(1+4a^2)}{1+a^2} X_2 + \frac{9\sqrt{2}}{4(1+4a^2)} X_1 X_3 \\ \dot{X}_3 &= -\frac{\sigma(4+a^2)}{1+a^2} X_3 - \frac{9\sqrt{2}a^2}{4(4+a^2)} X_1 X_2 \\ \dot{Y}_1 &= -Y_1 + rX_1 - X_1 Z + \frac{3\sqrt{2}}{4} X_3 Y_2 \\ \dot{Y}_2 &= 2rX_2 - \frac{(1+4a^2)}{1+a^2} Y_2 - \frac{3\sqrt{2}}{4} X_3 Y_1 - 2X_2 Z \\ \dot{Z} &= -bZ + X_1 Y_1 + 2X_2 Y_2 \end{aligned} \right\} \tag{8}$$

The analysis of the mathematical model described in (1), (2) is based on six-dimensional ordinary differential equations (8) truncated by (3),(4). Truncated equation, not only hope to show the existence of complex dynamical behavior, it is more important that the model can reflect the qualitative and quantitative characteristics in experiment. That is, the characteristics of finite dimensional dynamical system must reflect the characteristics of infinite dimensional dynamical system.

4 Numerical Simulation

This truncation scheme retains a series of bifurcation of ψ in the direction of the unstable mode. With the control parameter r increasing, the stability of Lorenz-like equations (8) will change, system (8) will appear the period-doubling bifurcation and chaos in nonlinear phenomenon. Here is a detailed numerical simulation the whole process from bifurcation to chaos of system (8) when $\sigma = 10, a = 3$ and the initial values are (1.0,5.0,5.0,5.0,5.0,5.0). Fig. 1 (a)-(j) display the phase portraits of solution trajectory in the different parameters. The dynamic characteristics of system (8) are shown in the parameters of different ranges by bifurcation diagram, Fig. 2. Lyapunov exponent spectrum, Fig. 3. When the control parameter r changes

approximately in area $53 \sim 64.4$, system loses stability from the stationary state no convection and appears period doubling bifurcation and forms the unstable manifold. As is shown in Fig. 1 (a)–(d). With r increasing, the track numbers gradually increase; the unstable manifold tends to the ipsilateral equilibrium point and forms the chaotic attractor. As is shown in Fig. 1 (e)–(i). When r arrives about 112.8 and continues to increase, the system enters the new bifurcation from chaos. As is shown in Fig. 1 (j).

As can be seen in Fig. 2, some blank tapes appear in chaotic region. These ‘blank tapes’ are called the period windows in the chaotic region, such as near area $r = 67.212, 72.48, 75.31, 90.515, 98.652$. These periodic trajectories are obtained by the bifurcation of the chaotic trajectory.

The positive values of the first Lyapunov exponent λ_1 imply sensitive dependence on initial condition for the system (8) in Figs. 3 and chaotic characteristics are shown at these values in Figs.2.

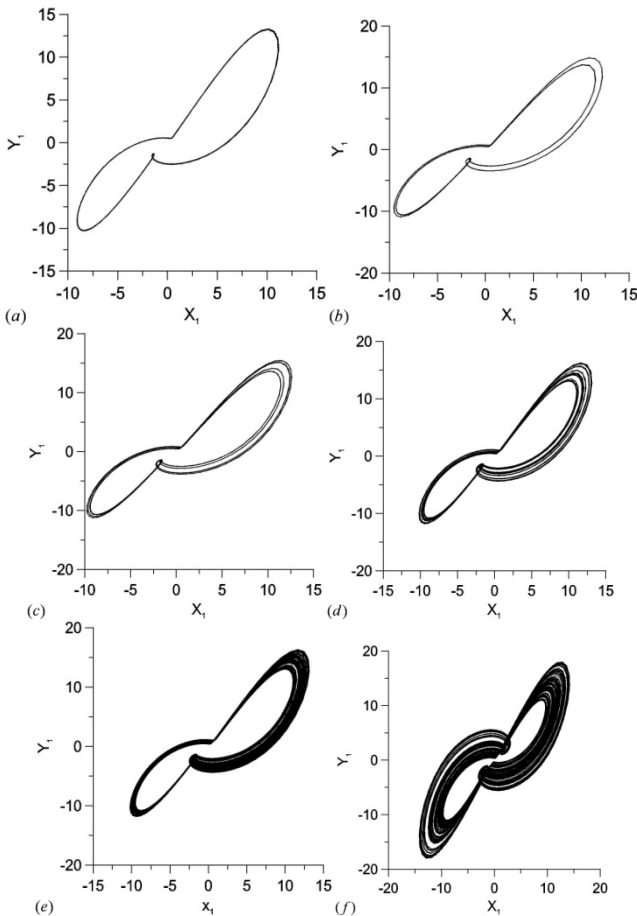


Fig. 1. Phase portraits at the specific parameter of the system (8)

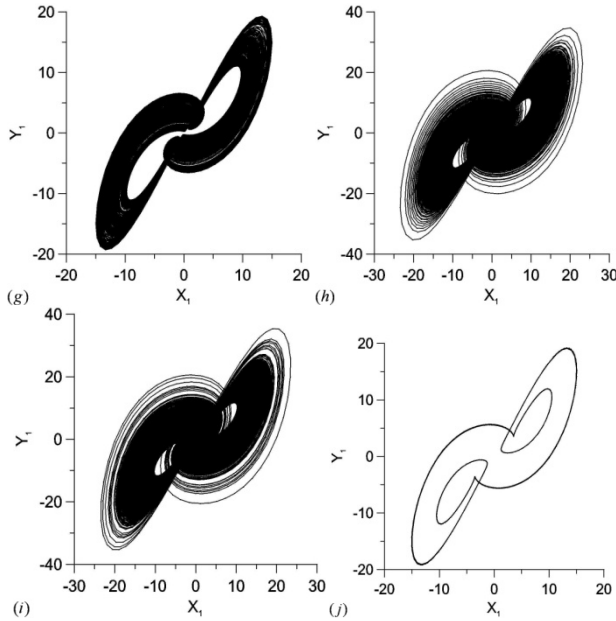


Fig. 1. (continued)

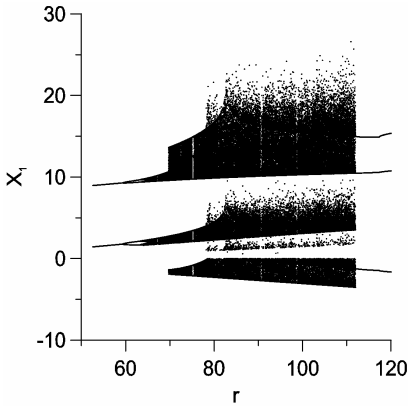


Fig. 2. Bifurcation diagram of the system (8)

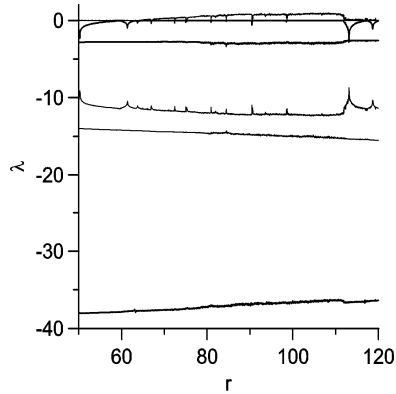


Fig. 3. Lyapunov exponents of the system (8)

5 Discussion

In the basis of Lorenz model, equations (1), (2) are again truncated to the six-dimensional ordinary differential equations by Fourier mode (3), (4). The basic dynamical behaviors and chaotic trajectories of equations are simulated numerically according to control parameter changes. The result shows the path from period doubling bifurcation to low-dimensional chaos of the system. The existence of low-dimensional chaotic trajectories helps understanding of turbulence in plane fluid

domains discussed by Lanford [10] and that the mathematical object which accounts for turbulence is attributed to one or a few low-dimensional chaotic attractors is verified. The data shown in numerical simulation provides a theoretical reference for the further study of turbulent flow and help chaotic characteristics of turbulent flow toward the desired direction.

References

1. Dong, F.: Turbulent Flow and Vortex: Discussion on the Architecture Language Expression of Chaos. Theory through the Two Non-linear Morphology. *Hua Zhong Architecture* 4, 11–14 (2010)
2. Ran, Z.: Multiscales and Cascading in Isotropic Turbulence. *Chinese Science Bulletin* 56(27), 2889–2892 (2011)
3. Cai, F., Zang, F., Liang, Y.: Nonlinear dynamic behaviors of a cracked hinged-hinged pipe conveying pulsating fluid. *Journal of Vibration and Shock* 31(4), 162–167 (2012)
4. Zhang, H., Liang, S., Song, S., Wang, H.: Truncation Error Calculation Based on Richardson Extrapolation for Variable-step Collaborative Simulation. *Science China Information Sciences* 54(6), 1238–1250 (2011)
5. Chen, Z., Price, W.G.: Transition to Chaos in a Fluid Motion System. *Chaos, Solitons and Fractals* 26, 1195–1202 (2005)
6. Xu, M., Wang, H.: Bifurcation Problems of the Model System Similar to the Lorenz Equations of the Flow between Two Concentric Rotating Spheres. *Journal of Mathematics* 27, 111–118 (2007)
7. Wang, H., Cui, Y.: Nine-modes Truncation of the Plane Incompressible Navier-Stokes Equations. *Communications in Mathematical Research* 27(4), 297–306 (2011)
8. Edward Lorenz, N.: Deterministic Nonperiodic Flow. *Journal of the Atmospheric Sciences* 20, 130–141 (1963)
9. Saltzman, B.: Finite amplitude free convection as an initial value problem. *Atmos. Sci.* 19, 329–341 (1962)
10. Landau, L.: On the Problem of Turbulence. *CR Acad. Sci. URSS* 44, 311–315 (1944)

Security Uniform Office Format Specification and API Design Based on the Java Platform

Ying Cai, Ning Li, and Chengxia Liu

Dept. of Computer Science and Technology, Beijing Information Science & Technology
University Beijing, 100101, P.R. China
{ycai, lining, liucx}@bistu.edu.cn

Abstract. “Uniform Office Format” (UOF) is the Chinese national standard for office document format. It aims to solve the problem that the office document format is not unified. It is helpful for information sharing and document exchange. Moreover, it improves the Chinese office software’s compatibility, and establishes the foundation for the document information exchange. Through the study of Specifications for the Syntax and Processing of UOF and XML Encryption/Signature, we proposed a set of APIs which ensure the confidentiality, integrity, and authentication of the UOF documents. The API is implemented by Java Cryptography Architecture, Java Cryptography Extension, Dom4j, and XPath technologies on Java platform and is confirmed to XML Digital Signature APIs (JSR105) and XML Digital Encryption APIs (JSR106) Specifications. Every interface and class is designed according to the rules that should keep high cohesion inside module and low coupling between modules. It also implements a special set of APIs for security UOF documents. Furthermore, a set of test cases together with a GUI program are developed to verify the availability and correctness of the system. It is believed that this project will contribute a lot to the research and adoption of the Specification for the Syntax and Processing of UOF Encryption/Signature.

Keywords: Uniform Office Format, Specification, Application Programming Interface.

1 Introduction

With the rapid development of network technology, office documents on the network for transmission of information and sharing has become a very common trend, which is bound to consider the document in a shared process of security and non-repudiation and involves to the encryption and signature problems. Encryption can provide confidentiality for data services. Signature provides data services for data integrity, message authentication and authentication. In order to document security, we must add the corresponding encryption and data signature in the UOF documents.

“Chinese office software application programming interface specification” [1] has been issued. We should add the security schema in the UOF documents and design corresponding API that satisfies the Chinese documents software API specification.

And we must follow the Chinese office software application programming interface specification including platform-independent, language-independent, product-independent, product portability and software reuse and application systems together. In this paper we design and implement a set of security UOF API according to the Java Cryptography Architecture (JCA), Java Cryptography Extension (JCE), Dom4j, XPath technology [2][3][4] based on the Java platform. There are also significant differences when security UOF API design follows the JSR105, JSR106 specification [5] [6]. JSR105, JSR106 is an API specification [7] [8] [9] based on XML markup language encryption and decryption / data signature. And security UOF API is designed for the higher level XML documents package. In this paper we implement the application programming interface for the security UOF document and it could encrypt and decrypt or signature the document as the whole or partial contents.

2 Outline Design

2.1 Security UOF Document Architecture

There are some main contents as following in the security UOF documents

- (1) Word processing: applications for word processing document format description;
- (2) Presentation: the presentation describes the application of the document format;
- (3) Spreadsheet: the spreadsheet application document format description;
- (4) Encryption zone: storing encrypted node for the use of decryption;
- (5) Signature zone: storage signature nodes for the use of authentication.

The structure of security UOF document is shown in Fig. 1.

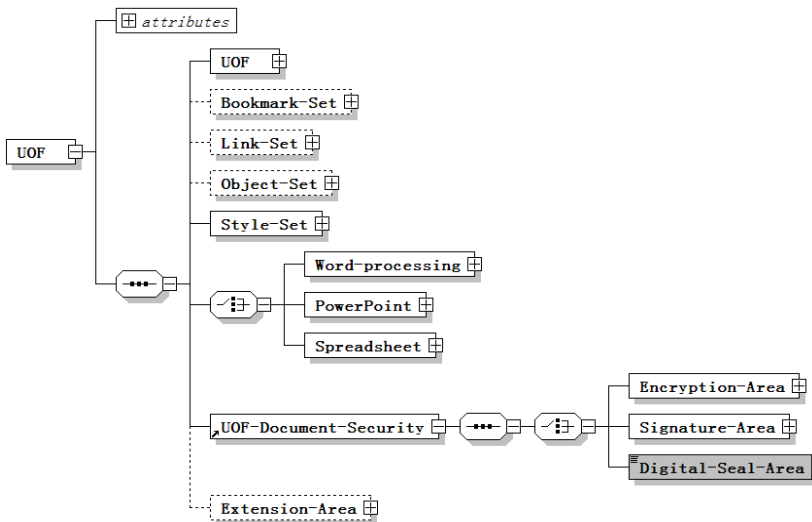


Fig. 1. Security UOF document architecture

3 Detailed System Design

The detailed system design of each module followed the principles including a single software engineering duties, opening and closing, Richter substitution, reverse dependency, interface isolation.

3.1 Key Management Module

The structure of Key elements

- (1) KeyName element: the reference to the key store alias.
- (2) KeyValue elements: storage the public key value encoded with Base64 including DSAKeyValue elements (public key information in the DSA algorithm) and RSAKeyValue sub-elements (public key information in the RSA algorithm).
- (3) RetrievalMethod elements: a reference to the key storage method.
- (4) X509Data (PGPData) elements: store the X.509 certificate information corresponding public key.
- (5) EncryptedKey elements: encrypted information node in an encrypted session key is generated and a session key when decryption the same node information.

3.2 Key Management Module Interface Design

We design it according to the definition of the KeyInfo element in the XML Signature Syntax and Processing [10][11]. A type definition is corresponding to an interface in schema that each interface has two main functions (information will be Encapsulated as XML elements and resolve the XML elements to obtain relevant information). The relationship between the interface modules is shown in Fig2.

Core interfaces are as follows.

- (1) KeyInfoType interface: Encapsulation keyinformation in the encryption or signature and key in the decryption or authentication.
- (2) KeyInfo interface: derived from KeyInfoType interface.

3.3 Key Management Module Functions in the Encryption and Signature

Using key information of encryption and signature is encapsulated into a KeyInfo element in the application environment of encryption and signature for this module.

The details are as follows.

- (1) The KeyName element is generated by encapsulation the Key alias if reference to the key through a key alias in the encryption or signature.
- (2) The public key will be encapsulated as KeyValue element in the encryption or signature.
- (3) If we use symmetric encryption system, we encrypt this session key to generate EncryptedKey element by using public key encryption algorithm.

(4) X.509 (PGP) certificate information is encapsulated into X509Data (PGPData) elements if there are X.509 (PGP) certificate as public key in the use of encryption or signature.

(5) If the key is referenced by other method, the key reference method will be encapsulated as RetrievalMethod elements. And the generated elements encapsulated in the KeyInfo element.

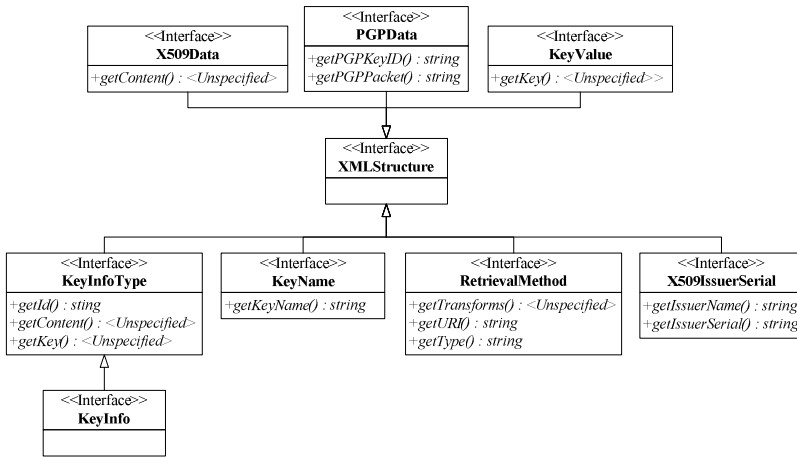


Fig. 2. Key management module interface relationship

3.4 Key Management Module Functions in the Decrypting and Authentication

We analysis the KeyInfo element to get the decryption key or authentication public key in the application environment of decryption and authentication for this module. The details are as follows.

- (1) We resolve KeyName element to get the key alias and then get the key if the KeyInfo element contains KeyName elements in the decryption or authentication.
- (2) We resolve X509Data (PGPData) element to get the private key in the decryption or public key of authentication in the decryption or authentication.
- (3) We get the decryption key through resolving the encrypted key element if the KeyInfo element contains EncryptedKey element when in the decryption.
- (4) We get the key through resolving the RetrievalMethod element if the KeyInfo element contains the RetrievalMethod element.
- (5) We could get the public key of authentication through resolving the KeyValue element if the KeyInfo element contains the KeyValue element for authentication.

4 Encryption and Decryption Module

4.1 Encryption Node Structure

EncryptedData element is a node to store encrypted information which use it in the decryption situation. The details are as follows.

- (1) EncryptionMethod elements: encryption algorithms including theKeySize elements and other information.
- (2) CipherData elements: including CipherValue elements and CipherReference elements.
- (3) CipherReference elements: including a series of optional elements (Transforms). It is an ordered list of Transform elements which describe how to get to the data object to be processed. The output of each Transform must be as the input of next Transform. The input of the first Transform is the results of resolving URI attribute of Reference element and the results of last Transform is the input of encryption algorithm to deal with.

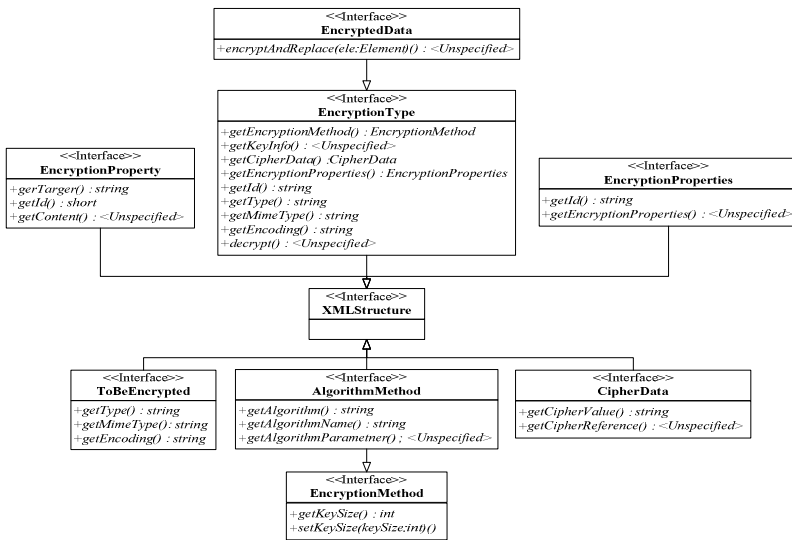


Fig. 3. Encryption module interface relationship

4.2 Encryption Module Interface Design

We design it according to the definition of the EncryptedData element in the XML Signature Syntax and Processing. A type definition is corresponding to an interface in schema that each interface has two main functions (information will be Encapsulated as XML elements and resolve the XML elements to obtain relevant information). The relationship between the interface modules is shown in Fig.3.

Core interfaces are as follows.

- (1) EncryptedType interface: Encapsulation the encryption node for the encryption algorithm, key information, cipher data and attributes in the application.
- (2) EncryptedData interface: derived from EncryptedType interface and encrypted the XML element.
- (3) EncryptedKey interface: derived from EncryptedType interface and encrypted the session key.

4.3 Encryption Process

Encryption process is as follows.

- (1) Encrypted data, type, algorithm and key.
 - (2) Encryption data according above the information.
 - (3) Construct KeyInfo element.
 - (4) Constuct EncryptedData element .
 - (5) Add EncryptedData element to the encryption area.
- The detailed process is shown in Fig. 4.

4.4 Decryption Process

The detailed decryption process is shown in Fig5.

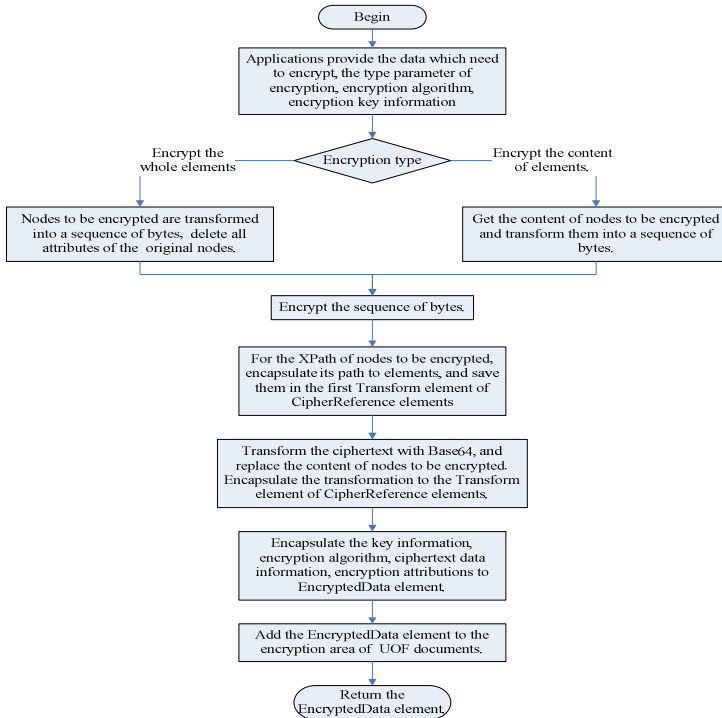


Fig. 4. Encryption process flow

5 Signature Module

5.1 Signature Node Structure

Signature node element is the node storing corresponding signature information. We do authentication operation based on this node. The details are as follows.

- (1) Signature elements: including SignedInfo element, SignatureValue elements and KeyInfo element.
- (2) SignedInfo element: including the CanonicalizationMethod algorithm element, SignatureMethod elements and one or more reference.
- (3) Reference element: including the digest algorithm (DigestMethod) elements of calculation the reference element of the signature, values (DigestValue) elements and a transform list (Transforms) elements.
- (4) Transforms elements: including an ordered list of Transform elements which describe how the signer obtained the data object. The output of each Transform is to be as the input of next Transform. The input of the first Transform is results of resolving the URI attribute of Reference element. The final result is the input of DigestMethod Transform algorithm. Each Transform element consists of a Transform algorithm attributes and parameters of the algorithm.

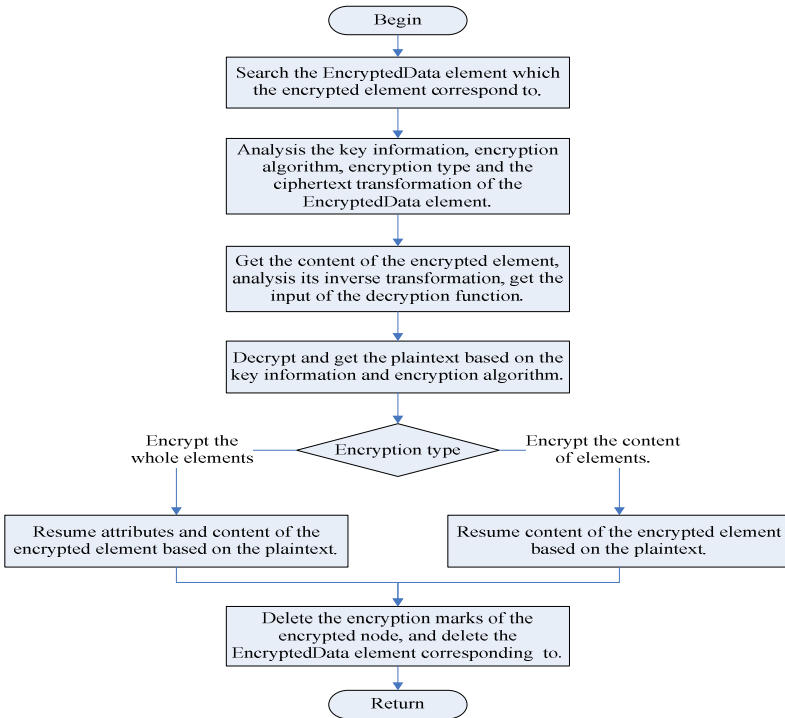


Fig. 5. Encryption process flow

5.2 The Signature Module Interface Design

The relationship between the interface module is shown in Fig. 6. The details are as follows.

(1) XMLSignature interface: we signature it according the application that provide signature algorithm, digest algorithm, regularization algorithms, key information and the signature reference data elements. We encapsulate the relative information, sign value, and digest value into the node in the signature. We resolve the same node to get the corresponding information and do authentication.

(2) SignedInfo interface: we encapsulated the application provides the signature algorithm, digest algorithm, regularization algorithms, and signature information into the SignedInfo element in the signature. We resolve the same SignedInfo element to get the corresponding information.

(3) Transform Interface: we encapsulated a series of transformations of the signature reference element in the signature. And we resolve the Transform elements to get the corresponding transform algorithm in the authentication.

(4) Reference interface: we encapsulated the signature reference element information into the Reference element in the signature. We resolve the Reference element to obtain the signature reference element and a series of transformation elements.

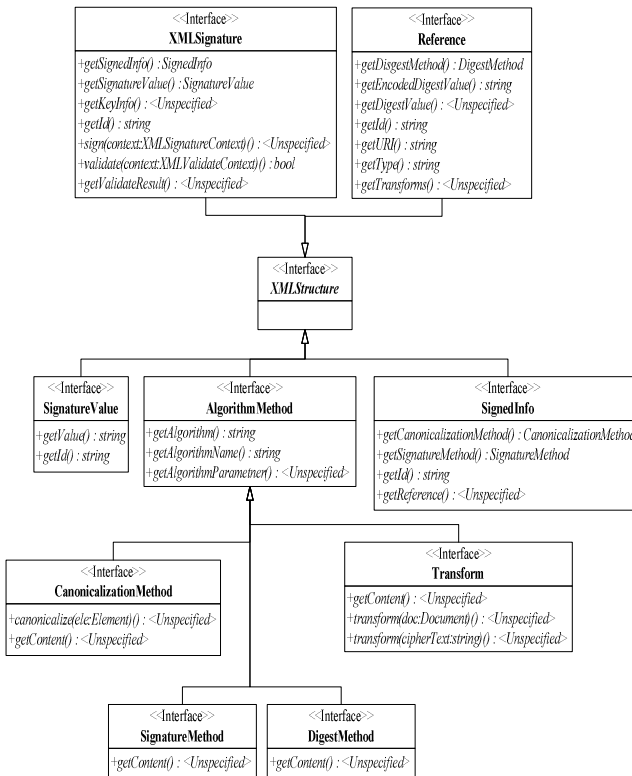


Fig. 6. Signature module interface relationship

6 Conclusion and

In this paper we give the security UOF specification and design the corresponding API. We design a test system in Java platform and test the various interfaces available to verify the usefulness and reliability of the system.

Acknowledgment. The research is supported by Funding Project for Academic Human Resources Development in Institutions of Higher Learning under the Jurisdiction of Beijing Municipality (PHR201007131) and the General program of science and technology development project of Beijing Municipal Education Commission under Grant No.KM201110772013.

References

1. Chinese office software document format [DB/OL], <http://www.uof.org.cn/>
2. JavaTM Cryptography Architecture (JCA) Reference Guide [DB/OL], <http://java.sun.com/javase/6/docs/technotes/guides/security/crypto/CryptoSpec.html>
3. O'Reilly.Java Security, 2nd edn. [DB/OL], <http://www.ibook8.com/software/Catalog113/2261.html>
4. Apache xml-security API Documentation [DB/OL], <http://santuario.apache.org/Java/api/index.html>
5. SR-000106 XML Digital Encryption APIs [DB/OL], <http://jcp.org/aboutJava/communityprocess/pr/jsr106/index.html>
6. JSR-000105 XML Digital Signature APIs [DB/OL], <http://jcp.org/aboutJava/communityprocess/final/jsr105/>
7. XML Encryption Syntax and Processing [DB/OL], <http://www.w3.org/TR/xmlenc-core/>
8. XML Signature Syntax and Processing [DB/OL], <http://www.w3.org/TR/xmlsig-core/>
9. Canonical XML [DB/OL], <http://www.w3.org/TR/xmlenc-core/>
10. XML Security Library Reference Manual [DB/OL], <http://www.aleksey.com/xmlsec/api/index.html>
11. XML Security C [DB/OL], <http://santuario.apache.org/c/apiDocs/hierarchy.html>

Bioinformatics Analysis of the Complete Nucleotide Sequence of Duck Plague Virus UL22 Gene

Li-Sha Yang^{1,2,3}, An-Chun Cheng^{1,2,3}, Ming-Shu Wang^{1,2,3}, De-Kang Zhu^{2,3},
Shun Chen^{1,2,3}, Ren-Yong Jia^{1,2,3}, and Xiao-Yue Chen^{2,3}

¹ Institute of Preventive Veterinary Medicine, Sichuan Agricultural University,
Wenjiang, Chengdu City, Sichuan, 611130, P.R. China

² Avian Disease Research Center, College of Veterinary Medicine of Sichuan Agricultural
University, 46 Xinkang Road, Ya'an, Sichuan 625014, P.R. China

³ Key Laboratory of Animal Disease and Human Health of Sichuan Province,
Sichuan Agricultural University, Wenjiang, Chengdu City, Sichuan, 611130, P.R. China
chenganchun@vip.163.com, mshuwang@163.com

Abstract. UL22 gene, also named gH gene (GenBank accession No. EU195089) from duck plague virus (DPV) CHv strain which was isolated in our laboratory, is a 2505bp segment. The analysis about the characteristics of this gene may improve our understanding for the evolution and classification of DPV. In this study, the structure of phylogenetic tree and the prediction of gH functional sites depended upon MegAlign procedure and some online analysis softwares. The results revealed that this gene sequence was quite conservative in DPV. The phylogenetic tree of the gH protein and its homologs of other 20 reference herpesviruses showed glycoprotein H had a close evolutionary relationship with certain avian herpesviruses, such as MeHV-1, GaHV-2 and GaHV-3 which just were classified to Mardivirus genus. Besides, 8 N-glycosylation sites and 52 potential phosphorylation sites were found in gH and they maybe play an important role in regulating the function of DPV glycoprotein H.

Keywords: DPV, UL22 gene, bioinformatics analysis.

1 Introduction

Duck plague virus (DPV), also known as duck enteritis virus (DEV) or Anatid herpesvirus-1 (AnHV-1), is an important pathogen causing the serious contagious disease duck viral enteritis (DVE) or duck plague (DP) in bird of anseriformes (such as ducks, geese and swans) [1]. It has been grouped belonging to the subfamily alphaherpesvirinae based on the Eighth International Committee on Taxonomy of Viruses (ICTV), however, it has not been classified into any genus [2]. The same as the genome of other herpesvirus [3,4], DPV genome consists of two extended regions of unique sequence, unique long region (UL) and unique short region (US), and each of which is flanked by terminal redundancy (TR) and internal redundancy (IR) [5].

UL22 gene, which is located at the unique long region and also called gH gene pervades all members of herpesviridae [6]. DPV UL22 gene (GenBank accession No.

EU195089) was identified and sequenced in our laboratory, is highly conservative among the α -, β -, γ -herpesvirus.

Currently, along with the development of bioinformatics and molecular biology, properties about gE[7], UL53[8], US2[9,10], UL14[11], UL13[12] genes and so on from Duck Plague Virus have been reported, but there is little information on the molecular characteristics of UL22 gene. In this article, we report the results of bioinformatics analysis characteristics about the complete nucleotide sequence of UL22 gene. We hope these results will provide some useful materials for further research about the prevention and treatment method of DPV.

2 Virus Materials and Gene Sequences

The DEV CHv strain, which is a highly virulent field strain of DPV, was obtained from Key Laboratory of Animal Disease and Human Health of Sichuan Province, Sichuan Agricultural University[13]. We got the UL22 gene (complete nucleotide sequence) from NCBI submitted by our laboratory and the accession number is EU195089. The nucleotide sequences of the UL22-like genes of 20 reference herpesviruses were obtained from NCBI GenBank nucleotide database (show in table 1).

Table 1. The information about gH amino acid sequences of DPV and 20 reference herpesviruses

Virus name	GenBank accession No.	Abbreviation	Genus designations	Length (aa)
Duck plague virus	EU195089	DPV	Genus undesigned	834 aa
Gallid herpesvirus 2	ACF94933	GaHV-2	Mardivirus	813 aa
Gallid herpesvirus 3	NP_066853	GaHV-3	Mardivirus	812 aa
Meleagrid herpesvirus 1	NP_073315	MeHV-1	Mardivirus	808 aa
Bovine herpesvirus 2	AAK55404	BoHV-2	Simplexvirus	867 aa
Human herpesvirus 1	ADD59995	HHV-1	Simplexvirus	838 aa
Macacine herpesvirus 1	AAP41440	MaHV-1	Simplexvirus	847 aa
Papiine herpesvirus 2	YP_443868	PapHV-2	Simplexvirus	857 aa
Saimiriine herpesvirus 1	YP_003933818	SaHV-1	Simplexvirus	871 aa
Gallid herpesvirus 1	CAA74690	GaHV-1	Iltovirus	779 aa
Psittacid herpesvirus 1	NP_944394	PsHV-1	Iltovirus	806 aa
Bovine herpesvirus 1	NP_045337	BoHV-1	Varicellovirus	842 aa
Bovine herpesvirus 5	AAD40580	BoHV-5	Varicellovirus	849 aa
Canid herpesvirus 1	AAK51057	CaHV-1	Varicellovirus	796 aa
Cercopithecine herpesvirus 9	NP_077452	CeHV-9	Varicellovirus	852 aa
Equid herpesvirus 4	NP_045257	EHV-4	Varicellovirus	855 aa
Equid herpesvirus 9	YP_002333521	EHV-9	Varicellovirus	852 aa
Equid herpesvirus 1	AAS45924	EHV-1	Varicellovirus	848 aa
Felid herpesvirus 1	YP_003331559	FHV-1	Varicellovirus	821 aa
Human herpesvirus 3	ABE03056	HHV-3	Varicellovirus	841 aa
Suid herpesvirus 1	ABJ97135	SHV-1	Varicellovirus	685 aa

3 Methods and Sequence Properties

3.1 The Characteristics of DPV-UL22 Gene

Using the NCBI (<http://www.ncbi.nlm.nih.gov/>) BLASTN 2.2.26+ tool to search the similarity of the nucleotide sequence of DPV-UL22 gene[14,15], we find that the nucleotide sequence of DPV-CHv UL22 gene have 100% similarities with 2085 strain, VAC strain, clone-03 strain, AV1221 strain of DPV-UL22 gene(show in figure 1). Then the open reading frame(ORF) was predicted by NCBI ORF Finder[15,16], the result shows that DPV-UL22 gene is a 2505bp fragment, that is to say, the gene is an ORF itself (show in figure 2). Using NCBI CDD (http://www.ncbi.nlm.nih.gov/Structure/cdd/docs/cdd_search.html) to analyze the 2505ORF, there is one conserved domain (PHA03294 Superfamily) in the ORF (show in figure 3).

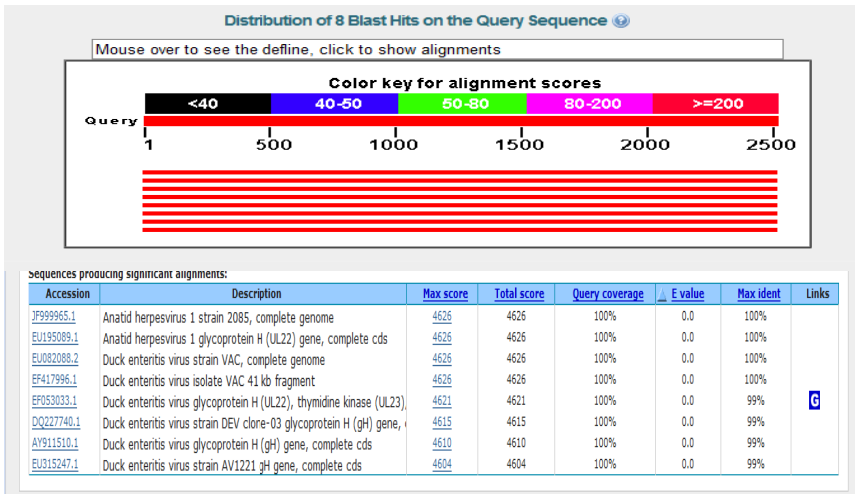


Fig. 1. The results of the similarity of DPV-UL22 gene with other herpesviruses

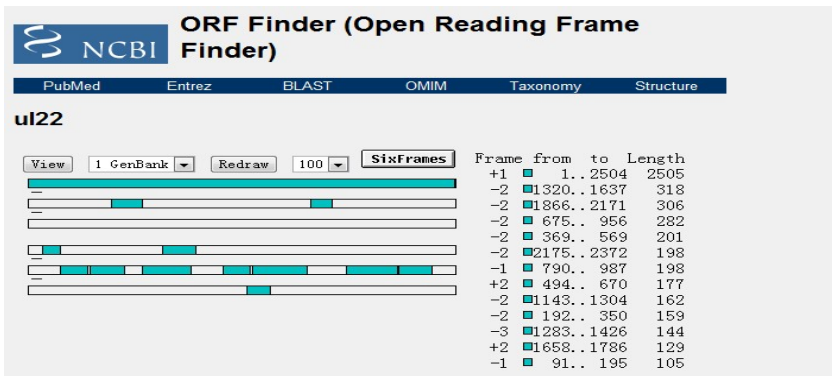


Fig. 2. The ORF of DPV-UL22 gene was discovered using NCBI ORF Finder

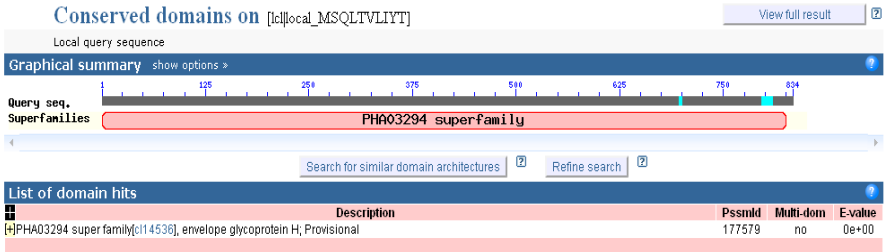


Fig. 3. The result of the conserved domain analysis of DPV-UL22 gene

3.2 Compare the Amino Acid Sequence of Glycoprotein H with Other 20 Reference Herpesviruses and Construct Phylogenetic Tree

Comparison with other herpesvirus revealed identities of 27%, 23%, 28%, 28%, 27%, 29%, 30%, 30%, 30%, 28%, 28%, 26%, 26%, 20% and 35% with the gH counterparts of the BoHV-1, BoHV-2, GaHV-2, MeHV-1, FeHV-1, GaHV-3, EHV-4, EHV-1, EHV-9, HHV-3, BoHV-5, HHV-1, HHV-2, EHV-2, SHV-1, respectively through BLASTP 2.2.26+ tool on NCBI. Meanwhile, we chose 20 reference herpesviruses from NCBI Genbank. Using the MegAlign procedure in DNASTAR 7.1 package to build the phylogenetic tree which based on DPV-CHv strain gH protein sequence and 20 reference herpesviruses gH protein sequence is shown in figure 4. The result demonstrates that DPV-CHv gH is similar to certain avian herpesviruses, such as MeHV-1, GaHV-2 and GaHV-3 which just are classified to Mardivirus genus. However, it is noteworthy that DPV is clustered within a monophyletic clade.

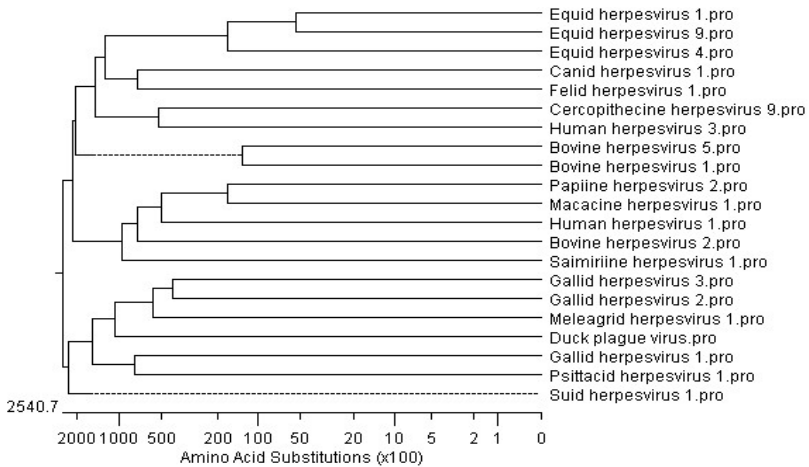


Fig. 4. Phylogenetic tree of glycoprotein H encoded by DPV-UL22 Gene based on 20 referenced gH sequences (show in Table 1) by the MegAlign program of DNASTAR 7.1

3.3 Presume the Functional Sites of Glycoprotein H

Through PROSITE motif search in PredictProtein prediction server (<https://www.predictprotein.org/submit>), we revealed that glycoprotein H had multiple functional sites. They are protein kinase C phosphorylation site, casein protein kinase II phosphorylation site, N-myristoylation site, amidation site, N-linked glycosylation site.

So as to futher analysis, we conjectured the N-linked glycosylation site (Asn-Xaa-Ser/Thr) and phosphorylation site by virtue of online analysis software NetNGlyc1.0 (<http://www.cbs.dtu.dk/services/NetNGlyc/>) and CBS website (<http://www.cbs.dtu.dk/services/NetPhos/>)[17]. The prediction results display 8 N-glycosylation sites and they lay on aa residues 32(NFTH), 77(NSTD), 147(NLSE), 493(NVTA), 659(NGTF), 643(NLSP), 752(NSST), 774(NGTI) (show in figure 5). There are 52 potential phosphorylation sites in all, including 31 serine phosphorylation sites, 14 threonine phosphorylation sites and 7 tyrosine phosphorylation sites(show in figure 6).

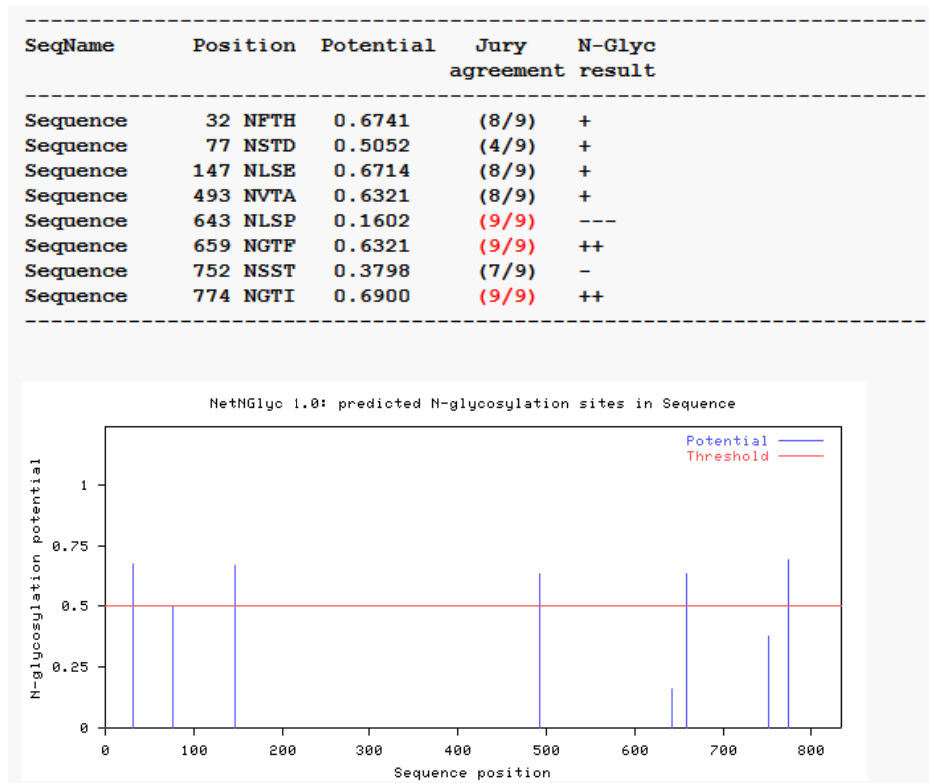


Fig. 5. The results of the N-glycosylation sites analysis of glycoprotein H encoded by DPV-UL22 gene

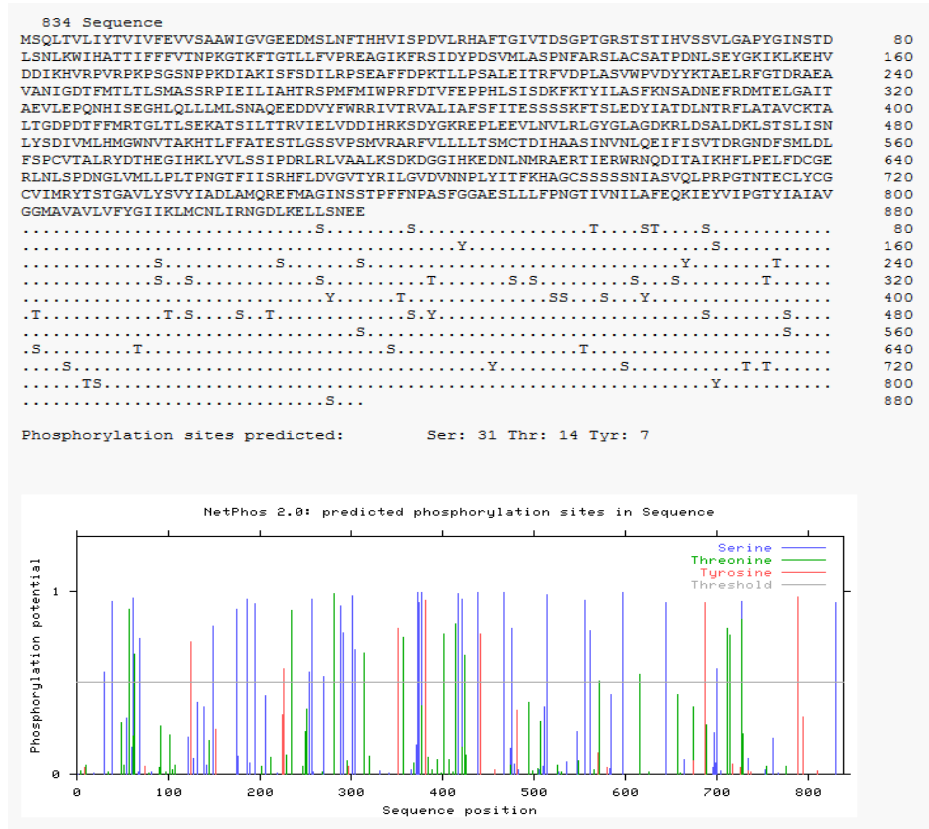


Fig. 6. The results of the phosphorylation sites analysis of glycoprotein H encoded by DPV-UL22 gene

4 Discussion

According to the similarity search of the nucleotide sequence of the DPV-UL22 gene, it suggested that the gene sequence was terrific conservative in DPV. If the nucleotide sequence in a DNA molecule may be encoded polypeptide or protein, it would able to be translated into protein by the ribosome. The 5'-end of the nucleotide sequence contains with start codons and the end of the translation contains with stop codons, so the location between start codons and stop codons is called the open reading frame(ORF). From the analysis consequences by NCBI ORF Finder, the ORF of UL22 gene is 2505bp length. Moreover, there are still some very short ORFs in this sequence, but generally speaking, they are difficut to encode proteins.

The construction of the phylogenetic tree on the base of DPV UL22 protein sequence and 20 reference herpesviruses in the research is shown in figure 4. The phylogenetic branch indicates that DPV gH protein has close evolutionary relationship with the Mardivirus genus in which MeHV-1, GaHV-2 and GaHV-3 are

included, but DPV is clustered within a monophyletic clade. All these results is extremely helpful for us to further understand the genetic characteristics of DPV UL22 gene and genetic relationships with other herpesviruses so that these information may provide a reliable reference for DPV classification studies.

Glycosylation is one of the most crucial post-translation modifications in eukaryotes, it can effect the antigenic determinants, charge properties and enzymatic properties, especially, the thermostability of protein. N-glycosylation is a co-translational mechanism involving the transfer of a precursor oligosaccharide to Asn residues in a small sequence Asn-Xaa-Ser/Thr, where X is any amino acid except Proline[18,19]. The phosphorylation also widely presents in the process of protein synthesis as a kind of chemical modification which is an important step in regulating the activity of protein[20,21]. Collectively, these N-glycosylation sites and potential phosphorylation sites may act a significant role in regulating the biological function of DPV glycoprotein H.

Acknowledgements. The research was supported by China 973 program (2011CB111606), China Agricultural Research System(CARS-43-8), the Changjiang Scholars and Innovative Research Team of Sichuan Agricultural University(PCSIRT0848),and National science and technology support program for agriculture(2011BAD34B03) An-chun Cheng and Ming-shu Wang are the corresponding authors at: Institute of Preventive Veterinary Medicine, Sichuan Agricultural University, &Key Laboratory of Animal Disease and Human Health of Sichuan Province, Sichuan Agricultural University, Wenjiang, Chengdu city, Sichuan, 611130, P.R.China & Avian Disease Research Center, College of Veterinary Medicine of Sichuan Agricultural University, 46# Xinkang Road, Yucheng district, Yaan 625014, China. Tel.: +86 835 2885774;fax: +86 835 2885774. E-mail address: chenganchun@vip.163.com.(A. Cheng); mshwang@163.com (M. Wang).

References

1. Sandhu, T.S., Metwally, S.: Duck Virus Enteritis (Duck Plague). Blackwell, New York (2008)
2. Fauquet, C.M., Mayo, M.A., Maniloff, J., Desselberger, U., Ball, L.A.: Virus Taxonomy: Classification and Nomenclature of Viruses: Eighth Report of the International Committee on the Taxonomy of Viruses. Elsevier Academic, London (2005)
3. Ushijima, Y., Luo, C.H., Goshima, F., Yamauchi, Y., Kimura, H., Nishiyama, Y.: Determination and Analysis of the DNA Sequence of Highly Attenuated Herpes Simplex Virus Type 1 Mutant HF10, a Potential Oncolytic Virus. *Microbes Infect.* 9, 142–149 (2007)
4. Boehmer, P.E., Lehman, I.R.: Herpes Simplex Virus DNA Replication. *Annu. Rev. Biochem.* 66, 347–384 (1997)
5. Wu, Y., Cheng, A.C., Wang, M.S., Yang, Q., Zhu, D.K., Jia, R.Y., Chen, S., Zhou, Y., Wang, X.Y., Chen, X.Y.: Complete Genomic Sequence of Chinese Virulent Duck Enteritis Virus CHv Strain. *J. Virol.* (accepted, March 2012)
6. Nicolson, L., Cullinane, A.A., Onions, D.E.: The Nucleotide Sequence of an Equine Herpesvirus 4 Gene Homologue of the Herpes Simplex Virus 1 Glycoprotein H Gene. *J. Gen. Virol.* 71, 1793–1800 (1990)

Evidence Conflict Analysis Approach to Obtain an Optimal Feature Set for Bayesian Tutoring Systems

Choo-Yee Ting, Kok-Chin Khor, and Yok-Cheng Sam

Multimedia University, Jalan Multimedia, 63100 Cyberjaya, Malaysia
{cyting, kckhor, ycsam}@mmu.edu.my

Abstract. Identifying the appropriate features for constructing a Bayesian student model is crucial to ensure that the model is always optimal. Feature sets can be identified via two types of feature selection algorithms: (i) algorithms that return a discrete set of features, and (ii) algorithms that rank features from the highest to the lowest importance with respect to a class label. To determine the optimal feature set from the second type of feature selection algorithm has always been a challenge, mainly because indifference in overall predictive accuracies between feature sets often occurs. In this light, this paper proposes evidence conflict analysis approach to tackle the challenges. This approach analyzes the conflicts in evidence when a Bayesian Network is employed as a student model. To demonstrate the proposed method, the experiments in this study had utilized two datasets that were transformed from 244 students' log data. The empirical findings suggested that evidence conflict analysis can differentiate the performance of feature sets having the same overall predictive accuracy.

Keywords: Evidence Conflict Analysis, Feature Selection, Student Modeling, Bayesian Networks.

1 Introduction

Research work in student modeling has recently shown that Bayesian Networks (BNs) can be employed to develop student models in Intelligent Tutoring Systems (ITSs) [1, 2]. The acceptance of BNs as student modeling technique is largely because of its capability in handling uncertainty, encoding expert knowledge, and performing automatic probability update in light of new evidence. Examples of ITSs that employed BN as their student models are PrimeClimb [1], and ANDES [2]. These ITSs share a set of common characteristics: (i) small in BN size, (ii) containing manageable number of nodes (features), and (iii) tractable Conditional Probability Tables (CPTs). Such characteristics relax the process of feature selection when constructing a student model.

Although BNs have widely been recognized in ITSs, limited work, however, has reported on employing BNs in scientific inquiry learning environments. This type of tutoring environments emphasizes on acquisition of scientific inquiry skills such as formulating hypotheses, identifying and manipulating variables, conducting simulated

experiments, comparing results, and deriving conclusions. By acquiring such skills, a student's engagement and experience in learning science can further be enhanced [6]. In addition, researchers have reported that scientific inquiry skills are able to improve problem solving skills [3]. Often, students are granted the freedom to interact within the learning environment. Knowledge is constructed and idea can be refined as they interact with such learning environments. In general, the main characteristics of scientific inquiry learning environment are that it (i) employs exploratory learning approach, and (ii) enriches learning experience via rich user interface design (i.e., textbox, drag-and-drop, sliding, and etc). An example of scientific inquiry learning environment are BGuILE [3].

To our best knowledge, little work has been done to employ BN as a student model for scientific inquiry learning environments. This is because inferring the acquisition level of scientific inquiry skills from student interactions with ITSs can be a major challenge. The rich interaction between students and ITSs has resorted to huge amount of evidential nodes (features) in BN and therefore, identifying the important nodes is not trivial. To resolve this challenge, researchers in ITSs often employ feature selection techniques to reduce the number of features. Feature selection is important. Suppose there is a dataset with medium size of features, say N . It will generate a total of 2^N feature sets. If the features and the class involves are both binary type, then the hypothesis space for the dataset is 2^{2^N} [4]. However, three common questions that require attention are "What is the difference in overall predictive accuracy by having x , $x+1$, or $x-1$ features?", "How to indifference the performance of two feature sets if they elicit the same overall predictive accuracy?", or "How to detect conflicts in evidence?" In this light, this study attempted the evidence driven conflict analysis [15], in addition to feature selection and classifier approaches, to tackle challenges rose by the questions.

This paper shall proceed with an overview of INQPRO, a scientific inquiry learning environment in Section 2 while Section 3 provides a high-level discussion of evidence driven conflict analysis. The dataset, procedures, and experiment design to employ evidence conflict analysis shall be discussed in Section 4. The empirical findings are discussed in Section 5. Lastly, Section 6 highlights the major contributions of this work and provide directions this work can be further improved.

2 An Overview of INQPRO Learning Environment

This study used the scientific inquiry tutoring environment named INQPRO, developed by the author in previous work [5]. It has six different interfaces with different learning objectives and activities. Students are required to formulate hypothesis statement and identify the appropriate variables (i.e., manipulated, responding, and constant). Students can request assistance from the animated pedagogical agent (hereafter Peedy). Upon requests, Peedy will only provide indirect help to students.

In this study, all student interactions with INQPRO were logged and a log file was created for each student. The navigated interfaces and actions performed by each student were logged with timestamp. Examples of logged information are clicking of a button, number of times clicking Peedy, answering questions prompted by Peedy, selecting a value for manipulated variable, and drag-and-drop different mass. Features were then extracted from logged files before they can be used to infer a student's acquisition level of scientific inquiry skills.

3 Evidence Conflict

Most research work in Bayesian student model reported iterative designs of their student models before optimal model can be acquired [1, 2]. The ACE learning environment for instance employed two different student models, with the second revised model containing self-regulation nodes [2]. The ANDES Physics Tutor took almost several years to refine its student model [2]. Similar to ACE and ANDES, the INQPRO learning environment has its static and dynamic student model evolved from a preliminary version to a more stable version [5]. Various models were proposed mainly because it is difficult and almost impossible to construct the best student model, and therefore, the constructed model is often an approximation of the best model. In many previous ITSs works, overall predictive accuracy has been employed as an indicator to determine the fitness of a particular student model. In this study, not only that the overall predictive accuracy was used to measure the performance of a feature set, evidence conflict analysis was also employed to ensure that the feature set was optimal. Evidence conflict analysis can be used to detect possible conflicts among the evidence (features), but also between the evidence and the constructed model. If conflict occurs between the evidence and the model, it implies that the model might not be optimal and therefore, the model should be reconstructed by using other feature sets. The following is a brief discussion about evidence conflict analysis.

Let ε_i and ε_j be two pieces of evidence. Conflict is a measure that compares the joint probability of the evidence, $P(\varepsilon_i, \varepsilon_j)$, with the product of the probabilities of the individual pieces of evidence $P(\varepsilon_i)P(\varepsilon_j)$ [6]. Under normal conditions, $P(\varepsilon_i\varepsilon_j) > P(\varepsilon_i)$.

Conflict happens between the two pieces of evidence ε_i and ε_j , if $\log \frac{P(\varepsilon_i)P(\varepsilon_j)}{P(\varepsilon_i, \varepsilon_j)} > 0$.

Finally, the conflict measure can be defined as

$$\text{conf}(\varepsilon) = \log \frac{P(\varepsilon_i)P(\varepsilon_j)}{P(\varepsilon)}, \text{ where } \varepsilon = \{ \varepsilon_i, \varepsilon_j \} \quad (1)$$

In general, there is a conflict between the evidence ε if $\text{conf}(\varepsilon) > 0$. Leveraging on this idea, it was hypothesized that any given two or more feature sets that elicit the same overall predictive accuracy could be differentiated by using the conflict measure. That is, a feature set with the lowest conflict measure is considered optimum.

4 Methodology

Fig. 2. depicts the high-level presentation of the overall experiment design to determine the optimal feature set. The process begins with dataset preparation, which involves data cleaning and data transformation. The dataset was divided into the training and testing set using 80-20 split. Features were ranked based on their importance to the class before the model construction phase. The constructed models were tested with the testing set to elicit their overall predictive accuracies and conflict measures. Detailed discussion about the dataset preparation, model development, and accuracy measurement is given in the following subsections.

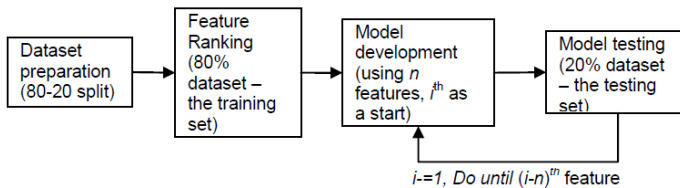


Fig. 2. High-level presentation of experiment design

4.1 Dataset Preparation

A total of 244 records were used in this study. They were generated using students' log data that was gathered over a period of one academic year from the Faculty of Information Technology, Multimedia University, Malaysia. The log data was preprocessed and transformed into two datasets. Both datasets consist of 59 observable features. Observable features are features that can directly be extracted from student interactions with INQPRO. Examples of these features are cnt_Scenario, t_scenario, drag_progressbar, and ans_agent_maniVar. Although both datasets share the same observable features, however, they can be distinguished by their class labels. Let sk represent scientific inquiry skills acquisition dataset and have class label "AcquisitionLevel". The class label takes one of the values {mastery, partial-mastery, non-mastery}. Let Mcc represent the conceptual change dataset and has "Changed" as class label. The class has two possible values {yes, no}.

In this study, each dataset was given an 80-20 split. Let D_{sk_80} denotes the 80% of D_{sk} while D_{sk_20} represents 20% of D_{sk} . Similarly for Mcc, the 80% and 20% of D_{cc} are denoted by D_{cc_80} and D_{cc_20} , respectively. The 80% of each dataset was meant for feature ranking and model development phase while 20% of the dataset was used for model evaluation and conflict analysis.

4.2 Model Construction

In this study, the feature selection technique used was the GainRatioAttributeEval evaluator with Ranker search method in WEKA [19]. Rather than identifying a distinct set of features, the algorithm ranked features from the most to the least

influence with respect to the class labels, which are AcquisitionLevel and Changed for M_{sk} and M_{cc} , respectively.

To determine the performance of a feature set, a classifier must be firstly developed. The classifier needs to be trained using the training dataset before tested using the testing dataset. In this study, for each dataset, the model development process began by using 59 features to design the first model, and followed by 58 features for the second model. Such procedure was repeated by removing one feature at a time based on the resulted feature ranking until 59 models were constructed. All the models constructed in this work took the form of a Naïve Bayesian Network (NBN). The main reason to use NBN is because of the simple construction method and easy in computation [7]. Let M_{sk_n} denotes the model built using portion of D_{sk_80} while M_{cc_n} was built using portion of M_{cc_80} , with n denotes the number of features used for building that particular model. For instance, M_{sk_30} represents that a classifier had been built using top 30 features using D_{sk} . Similarly, M_{cc_5} denotes that the top five features from D_{cc} were used to build the model.

4.3 Performance Evaluation

The evaluation methods employed in this study were predictive accuracy measurement and conflict measurement. To calculate the overall predictive accuracy of a NBN model developed for a particular dataset, each record of the dataset is fed into the model to elicit the predictive accuracy, with the given formula below:

$$\text{Overall Predictive Accuracy} = \frac{n_{\text{correct}}}{N} \times 100\%$$

The numerator n_{correct} represents the total number of correctly classified P in the testing dataset. The denominator N represents the total number of P in the testing dataset.

Evidence conflict analysis was subsequently performed using E.q. (1). To perform evidence conflict analysis, the testing dataset was fed into a particular model to elicit the conflict measure. The number of records with conflict > 0, average conflict measure, standard deviation for conflict measure, the minimum and maximum of conflict measure were calculated.

5 Results and Discussion

5.1 Optimal Feature Set for Dataset D_{sk}

Table 3 shows the overview results of experiment conducted using D_{sk_20} . In summary, the overall predictive accuracy was 57.1454% while the average number of records showed conflict with a particular model (i.e., $\text{conf}(\epsilon) > 0$) was 7.32. The average conflict measure was smaller than 0, indicating that the dataset generally fits most of the classifiers. The maximum conflict measure was 0.4359, which again indicated that D_{sk_20} has low evidence conflict given the models.

Table 3. Result summary for testing dataset \mathcal{D}_{sk_20}

Total features except Class label	Overall Predictive Accuracy (%)	Average number of records with conflict	Average conflict measure	Standard deviation for conflict measure	Max Conflict measure	Min Conflict measure
59	57.1454	7.32	-1.0095	1.0405	0.4359	-4.1350

Table 4. Top five feature sets ranked based on matching accuracy using \mathcal{D}_{sk_20}

M_{sk_n} (n represents the number of features used)	Overall Predictive Accuracy (%)	Total number of records with conflict > 0	Average Conflict	Standard Deviation	Max conflict	Min conflict
M_{sk_26}	62.3711	7	-1.0442	1.1959	0.4961	-4.0245
M_{sk_25}	62.3711	8	-0.8704	1.1112	0.5229	-3.6675
M_{sk_24}	63.4021	8	-0.7888	0.9997	0.4790	-3.5834
M_{sk_23}	63.4021	9	-0.7535	0.9853	0.4881	-3.6314
M_{sk_22}	62.8866	9	-0.7202	0.9582	0.4784	-3.7747

Table 4 shows the first five models having the highest overall predictive accuracies. The models are those having feature set of 26 features, 25 features, 24 features, 23 features, and 22 features. From the results, M_{sk_26} , M_{sk_25} , and M_{sk_22} can be ignored because they depicted lower overall predictive accuracies as compared to models with 23 and 24 features, which depicted the highest overall predictive accuracy (63.4021%). Although the M_{sk_24} and M_{sk_23} having the same accuracy, they are different in the total number of records with conflict and average conflict measure. By taking into consideration of overall predictive accuracy, total number of conflicted records, and average conflict measure when comparing the models with 23 and 24 features, it can be concluded that the latter performed optimally. Due to space limitation, the features are not listed out here.

5.2 Optimal Feature Set for Dataset \mathcal{D}_{cc}

Table 5 shows the overall results when the dataset \mathcal{D}_{cc_20} was used. The main difference between the tables is that the evidence conflict is higher for \mathcal{D}_{cc_20} . This can be observed from the average number of records that have $\text{conf}(\epsilon) > 0$, and average conflict measure.

Table 5. Result summary for testing dataset \mathcal{D}_{cc_20}

Total features except Class label	Overall Predictive Accuracy (%)	Average number of records with conflict	Average conflict measure	Standard deviation for conflict measure	Max Conflict measure	Min Conflict measure
59	72.0939	16.81	-0.2001	0.5284	1.2581	-1.9010

Table 6. Top five feature sets ranked based on overall predictive accuracy using dataset \mathcal{D}_{cc_20}

M_{cc_n} (n represents the number of features used)	Overall Predictive Accuracy (%)	Total number of records with conflict > 0	Average Conflict	Standard Deviation	Max conflict	Min conflict
M_{cc_11}	74.7423	30	-0.0300	0.3330	1.4831	-1.2972
M_{cc_10}	74.7423	13	-0.0336	0.3341	1.4826	-1.3052
M_{cc_9}	74.7423	13	-0.0320	0.3352	1.4822	-1.3119
M_{cc_8}	74.2268	14	-0.0185	0.2831	1.4006	-0.8509
M_{cc_7}	74.2268	19	-0.0244	0.2828	1.3978	-0.8386

Table 6 presents the extracted first five feature sets with highest overall predictive accuracies using dataset \mathcal{D}_{cc_20} . The feature sets that comprised of 9, 10, and 11 features demonstrated same overall predictive accuracy (74.7423%). M_{sk_7} and M_{sk_8} depicted lower matching accuracies (74.2268%) as compared to M_{sk_11} , M_{sk_10} , and M_{sk_9} , and therefore can be relaxed for further processing. From the first three models M_{sk_11} , M_{sk_10} , and M_{sk_9} , the model with 11 features depicted highest number of conflicted records and were therefore discarded. From the two remaining feature sets (represented by models M_{sk_11} , M_{sk_10}), the feature set with 10 features (M_{sk_10}) demonstrated lowest average conflict measure (-0.0336), and therefore, can be considered optimal. The ten features are `t_SimExp_ani`, `Path_After_Hypo`, `Path_Before_Variable`, `Path_After_Variable`, `t_CompareData`, `t_Scenario`, `t_Drag_Drop_Mass_SimExp`, `t_total_KeyinData_Graph`, and `cnt_Verify_SelfRegulate`.

6 Conclusion and Future Directions

There are many features that can be extracted and identify in a scientific inquiry learning environment due to the freedom granted to student in exploration. Therefore, merely ranking features from the highest importance to the lowest importance with respect to the class label is not meaningful. Predictive accuracy measurement is not sufficient to determine the optimal feature set because there are feature sets with the

same predictive accuracy (refer to Table 2 and 4). In this light, this research work proposed the evidence conflict analysis approaches to tackle the problem. Findings from empirical studies reported in this paper have concluded that the evidence conflict analysis has successfully differentiated feature sets that exhibit the same matching accuracy.

The datasets used in this study were relatively small because they were transformed from 244 student log files. Therefore, further investigation about the evidence conflict analysis could be performed with more log files. Apart from the limitation in dataset size, the findings reported in this study were confined to Naïve Bayesian Network. Future investigation could be performed to study the performance of evidence conflict analysis using different types of Bayesian Networks, which include the machine-learned Bayesian Networks, and expert elicited network.

References

1. Conati, C., Maclaren, H.: Empirically building and evaluating a probabilistic model of user affect. *Journal of User Modeling and User-Adapted Interaction* 19(3), 267–303 (2009)
2. Muldner, K., Burleson, W., van de Sande, B., VanLehn, K.: An analysis of students' gaming behaviors in an intelligent tutoring system: predictors and impacts. *Journal of User Modeling and User-Adapted Interaction* 21(1-2), 99–135 (2011)
3. Hulshof, C.D., Wilhelm, P., Beishuizen, J.J., Van Rijn, H.: FILE: A Tool for the Study of Inquiry Learning. *Computers in Human Behavior* 21, 945–956 (2005)
4. Liu, H., Motoda, H.: *Computational Methods of Feature Selection*. Chapman & Hall/CRC, Boca Raton, FL (2008)
5. Ting, C.Y., Phon-Amnuaisuk, S.: Factors influencing the performance of Dynamic Decision Network for INQPRO. *Computers & Education* 52(4), 762–780 (2009)
6. Kjaerulff, U., Madsen, A.: *Bayesian Networks and Influence Diagrams: a Guide to Construction and Analysis*. Springer, New York (2008)
7. Han, J.W., Kamber, M.: *Data Mining: Concepts and Techniques*. Morgan Kaufmann (2006)

Binary Vote Assignment Grid Quorum for Managing Fragmented Database

A. Noraziah, Ainul Azila Che Fauzi, Noriyani Mohd Zin, and Tutut Herawan

Faculty of Computer Systems and Software Engineering,
Universiti Malaysia Pahang Lebuhraya Tun Razak, 26300 Kuantan Pahang, Malaysia
{noraziah, tutut}@ump.edu.my, ainulazila@yahoo.com,
noriyanimz@gmail.com

Abstract. Replication usually referred as mechanism to increase availability and performance in distributed databases. Handling fragmented database replication becomes challenging issue since the distributed database is scattered into split fragments. Fragmentation in distributed database is very useful in terms of usage, efficiency, parallelism and also for security. In this paper, we manage fragmented database replication and transaction management using a new proposed algorithm called Binary Vote Assignment on Grid Quorum (BVAGQ). This technique combines replication and fragmentation. This strategy partition the database into disjoint fragments. The result shows that managing replication and transaction through proposed BVAGQ able to preserve data consistency. It also increases the degrees of parallelism. This is because by using fragmentation, replication and transaction can be divided into several sub-queries that operate on the fragments.

Keywords: Data replication, Replication algorithm, Fragmentation, BVAGQ, Quorum.

1 Introduction

Replication is the process of copying and maintaining database objects in multiple databases that make up a distributed database system [1]. It is broadly installed in disaster tolerance systems to replicate data from the primary system to the remote backup system dynamically and on-line [2]. Data replication may occur if the same data is stored in multiple storage devices. Replication is also frequently referred as mechanism to increase availability and performance in distributed databases [3]. Handling fragmented database replication becomes challenging issue to administrator since the distributed database is scattered into split replica partitions or fragments. Each partition or fragment of a distributed database may be replicated into several different sites in distributed environment [4]. Changes applied at one site are captured and stored locally before being forwarded and applied at each of the remote locations. Fragmentation in distributed database is very useful in terms of usage, efficiency, parallelism and also for security. This strategy will partition the database into disjoint fragments. If data items are located at the site where they used most frequently,

locality of reference is high. In fragmentations, similarly, reliability and availability are low [5]. But by combining fragmentation with replication, performance should be good [5]. Even if one site becomes unavailable, users can continue to query or even update the remaining fragments.

There are several protocols on replication. Read-One-Write-All Monitoring Synchronization Transaction System (ROWA-MSTS) Protocol proposed by Noraziah, A. et al. [6] replicates consistencies which is guaranteed by the consistency of execution on one replica, but the client replicas are only updated and cannot provide accurate responses to queries. The other protocol is Binary Vote Assignment on Data Grid (BVADG) [7,8] where a data will replicate to the neighbouring sites from its primary site. Four sites on the corners of the grid have only two adjacent sites, and other sites on the boundaries have only three neighbours.

In this paper, we manage fragmented database replication and transaction management using a new proposed algorithm called Binary Vote Assignment on Grid Quorum (BVAGQ). This technique combines replication and fragmentation. This strategy partitions the database into disjoint fragments.

The rest of the paper is organized as follows. Section 2 describes the related works. Section 3 describes system model of the proposed BVAGQ algorithm and experiment result. Finally, the conclusion of this work is presented in Section 4.

2 Related Works

2.1 Data Replication

The process of creating and maintaining multiple copies of some computing resource is called replication. This technique has been widely used for achieving qualities like scalability, high performance, high availability and fault tolerance in computer systems. Consistency of the system state is another quality that may also vary to different degrees in a replicated system because of concurrent operations in the replicas. There are three categories of fragmented replication scheme data replication. They are all-data-to-all-sites, some-data-to-all-sites and some-data-to-some-sites. Read-One-Write-All (ROWA) [9,10] and Hierarchical Replication Scheme (HRS) [10] are the examples of all-data-to-all-sites protocols. ROWA has been proposed preserving replicated data in network environment [9,6]. Meanwhile, replication in HRS starts when a transaction initiates at site 1. All the data will be replicated into other sites. All sites will have all the same data. For some-data-to-all-sites category, The Majority Quorum protocol and Weighted Voting protocol employ voting to decide the quorum techniques [11]. A tree structure has been assigned to the set of replicas in this technique. The replicas are positioned only in the leaves, whereas the non-leaf nodes of the tree are regarded as “logical replicas”, which in a way summarize the state of their descendants [12]. Besides Voting Protocol, Tree Quorum (TQ) [11] can also be categorized in some-data-to-all-sites. These replication protocols make use of a logical tree structure. The cost and availability vary according to the failure condition, whereas they are constant for other replication protocols [13]. One more protocol in this category is Branch replication scheme [2]. Its goals are to increase the scalability, performance, and fault tolerance. Replicas are created as close as possible to the clients that request the data files. Using this technique, the growing of the replica

tree is driven by client needs. Binary Vote Assignment on Data Grid (BVADG) [7] is one of the protocols in some-data-to-some-sites protocol. A data will replicate to the neighbouring sites from its primary site. Four sites on the corners of the grid have only two adjacent sites, and other sites on the boundaries have only three neighbours. Thus, the number of neighbours of each sites is less than or equal to four.

2.2 Database Fragmentation

Fragmentation in distributed database is very useful in terms of usage because usually, applications work with only some of relations rather than entire of it. In data distribution, it is better to work with subsets of relations as the unit of distribution. The other benefit from fragmentation is the efficiency. Data is stored close to where it is most frequently used and for data that is not needed, it is not stored. By using fragmentation, a transaction can be divided into several subqueries that operate on fragments. So, it will increase the degrees of parallelism. Besides, it also good for security as data not required for local applications is not stored. So, it will not available to unauthorized users. There are two main types of fragmentation which are horizontal and vertical [4]. Fragmentation in a single database needs to be divided into two or more pieces such that the combination of the pieces yields the original database without any loss of information. Each resulting piece is known as a database fragment [8]. Fragmentation in distributed database is very useful in terms of usage, reliability and efficiency. Fragmentation phase is the process of clustering the information accessed simultaneously by applications in fragments, while the process of distributing the generated fragments over the database system sites is called allocation phase. To fragment a class, it is possible to use two basic methods which are vertical fragmentation and horizontal fragmentation. Other than that two methods, it is also possible to execute mixed or hybrid fragmentation on a class by combining both techniques. In the object model, vertical fragmentation breaks the class logical structure (its attributes and methods) and distributes them across the fragments, which will logically contain the same objects, but with different structures. On the other hand, horizontal fragmentation distributes class instances across the fragments, which will have exactly the same structure but different contents. Thus, a horizontal fragment of a class contains a subset of the whole class extension. Each partition/fragment of a distributed database may be replicated [8]. Changes applied at one site are captured and stored locally before being forwarded and applied at each of the remote locations. Synchronous replication can be categorized into several schemes, i.e., all data to all sites (full replication), all data to some sites and some data to all sites. Expensive synchronization mechanisms are needed in order to maintain the consistency and integrity of the replicated data in distributed environment.

2.3 A Heuristic Approach to Vertical Fragmentation Incorporating Query Information

The aim of this research is a heuristic approach to vertical fragmentation, which uses a cost model and is targeted at globally minimizing these costs. Firstly, take the Attribute Usage Matrix and Attribute Access Matrix in [14] to construct an Attribute Usage Frequency Matrix (AUFM) grouped by site that issues the queries. Secondly, we compute the *request* for each attribute at each site. Thirdly, assuming that we have

been given the values of transportation cost factors, and then we calculate the *pay* of each attribute at each site using the values of the *request* and values of cost factors. Finally, for each attribute, compare all the *pay* at all sites to find the minimal one. Then allocate attribute *ai* to the site that of the minimal *pay*. This approach leads to a better design because the change of input query information is reflected in the decision of fragmentation to reduce the total query costs. It also improves the deficiencies of all the other vertical fragmentation approaches, which make decision according to the affinities between each pair of attributes, by introducing a simplified cost model at the stage of vertical fragmentation [14].

2.4 Read-One-Write-All Monitoring Synchronization Transaction System (ROWA-MSTS)

In ROWA-MSTS techniques, replicas consistencies is guaranteed by the consistency of execution on one replica, but the client replicas are only updated and cannot provide accurate responses to queries. Synchronous replication methods guarantee that all replicas are maintained consistent at all times by executing each transaction locally only after all replicas have agreed on the execution order. Through this, a very strict level of consistency is maintained. However, because of the strict synchronization between replicas that is required for each transaction, synchronous replication methods have been deemed impractical and often times a centralized or client-server approach is preferred for systems that critically require strict consistency [6].

3 System Model and Experiment Results

Binary Vote Assignment Grid Quorum (BVAGQ) technique will be used to approach the research. Each site has a premier data file. In the remainder of this paper, we assume that replica copies are data files. A data will replicate to the neighbouring sites from its primary site. Four sites on the corners of the grid have only two adjacent sites, and other sites on the boundaries have only three neighbours. Thus, the number of neighbours of each sites is less than or equal to 4.

3.1 BVAGQ Fragmentation Database Model

In this section, we proposed the new BVAGQ algorithm by considering the distributed database fragmentation.

The following notations are defined.

V is a transaction.

S is relation in database.

S_i is vertical fragmented relation derived from S , where $i = 1, 2, \dots, n$.

P_k is a primary key.

x is an instant in T which will be modified by element of V .

T is a tuple in fragmented S .

S_{iP_k} is a horizontal fragmentation relation derived from S_i .

P_i is an attribute in S , where $i = 1, 2, \dots, n$.

$M_{i,j}$ is an instant in relation S , where $i, j = 1, 2, \dots, n$.

i represent a row in S .

j represent a column in S .

η and ψ are groups for the transaction V .

$\gamma = \alpha$ or β where it represents different group for the transaction V (before and until get quorum).

V_η is a set of transactions that comes before V_ψ , while V_ψ is a set of transactions that comes after V_η .

D is the union of all data objects managed by all transactions V of BVAG.

Target set = $\{-1,0,1\}$ is the result of transaction V , where -1 represents unknown status, 0 represents no failure and 1 represents accessing failure.

BVAG transaction elements $V_\eta = \{V_{\eta x,qr} \mid r = 1, 2, \dots, k\}$, where $V_{\eta x,qr}$ is a queued element of V_η transaction.

BVAG transaction elements $V_\psi = \{V_{\psi x,qr} \mid r = 1, 2, \dots, k\}$, where $V_{\psi x,qr}$ is a queued element of V_ψ transaction.

BVAG transaction elements $V_\eta = \{V_{\lambda x,qr} \mid r = 1, 2, \dots, k\}$, where $V_{\lambda x,qr}$ is a queued element either in different set of transactions V_η or V_ψ .

$\hat{V}_{\lambda x,ql}$ is a transaction that is transformed from $V_{\lambda x,qr}$.

$V_{\mu x,q1}$ represents the transaction feedback from a neighbour site. $V_{\mu x,q1}$ exists if either $V_{\lambda x,qr}$ or $\hat{V}_{\lambda x,ql}$ exists.

Successful transaction at primary site $V_{\lambda x,qr} = 0$ where $V_{\lambda x,qr} \in D$ (i.e., the transaction locked an instant x at primary). Meanwhile, successful transaction at neighbour site $V_{\mu x,q1} = 0$, where $V_{\mu x,q1} \in D$ (i.e., the transaction locked a data x at neighbour).

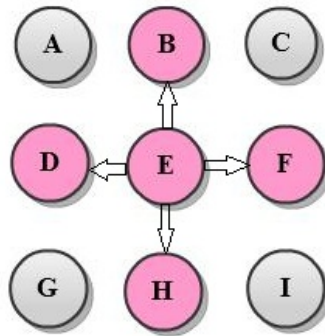


Fig. 1. Five replication servers connected to each

3.2 Experiment Result

To demonstrate BVAGQ transaction, 5 replication servers are deployed as in Figure 1. Each server or node is connected to one another through a fast Ethernet switch hub. There have been done two experiments. In the first experiment, there is only a transaction

request to update instant at one server meanwhile in second experiment there are two concurrent transactions request to update same instant at two different servers. Theoretically, each of the neighbours replication servers and the primary replication server should have to be connected each other logically.

Table 1. Experiment Results

REPLI CA	B	D	E	F	H
TIME					
t1	unlock(x)	unlock(x)	unlock(x)	unlock(x)	unlock(x)
t2	begin_tran saction	begin_tran saction	begin_transaction	begin_trans action	begin_tran saction
t3			$V_{\eta_{x,q_1}}$ write lock(x), counter_w(x)=1		
t4			$V_{\eta_{x,q_1}}$ propagate lock:D		
t5		$V_{\eta_{x,q_1}}$ lock(x) from E			
t6			$V_{\eta_{x,q_1}}$ get lock:D, counter_w(x)=2		
t7			$V_{\eta_{x,q_1}}$ propagate lock:B		
t8	$V_{\eta_{x,q_1}}$ lock(x) from E				
t9			$V_{\eta_{x,q_1}}$ get lock:B, counter_w(x)=3		
t10			$V_{\eta_{x,q_1}}$ propagate lock:F		
t11				$V_{\eta_{x,q_1}}$ lock(x) from E	
t12			$V_{\eta_{x,q_1}}$ get lock:F, counter_w(x)=4		
t13			$V_{\eta_{x,q_1}}$ propagate lock:F		
t14					$V_{\eta_{x,q_1}}$ lock(x) from E
t15			$V_{\eta_{x,q_1}}$ get lock:H, counter_w(x)=5		
t16			$V_{\eta_{x,q_1}}$ obtain quorum		
t17			$V_{\eta_{x,q_1}}$ update x		
t18			S is fragmented into S_1 and S_2		
t19			S_1 is fragmented into $S_{1(p,k,x)}$ and $S_{1(p,k,y)}$		
t20	commit $\hat{V}_{\lambda_{x,q_1}} \in V_{\eta}$	commit $\hat{V}_{\lambda_{x,q_1}} \in V_{\eta}$	commit $\hat{V}_{\lambda_{x,q_1}} \in V_{\eta}$	commit $\hat{V}_{\lambda_{x,q_1}} \in V_{\eta}$	commit $\hat{V}_{\lambda_{x,q_1}} \in V_{\eta}$
t21	unlock(x)	unlock(x)	unlock(x)	unlock(x)	unlock(x)

Using BVAGQ rules, each primary replica will copy other database to its neighbor replicas. Client can access other database at any server that has its replica.

From the result at Table 1, at time equal to 1 (t_1), instant x at all servers are unlock. At t_2 , the transaction begin. At t_3 , we can see that there is a transaction, $V_{\eta x, q1}$ requests to update instant x at server E . The transaction initiates lock. Hence, write counter for server E now is equal to 1. At t_4 , $V_{\eta x, q1}$ propagates lock at its neighbour replica D . At server D , $V_{\eta x, q1}$ lock(x) from E . Thus at t_6 , the transaction achieves in getting lock from the neighbor then write quorum for is equal to 2. Next, $V_{\eta x, q1}$ propagates lock at server B at t_7 and at t_8 , $V_{\eta x, q1}$ lock(x) from E . After that, $V_{\eta x, q1}$ propagates to server F and H , further successfully lock them.

At t_{16} , $V_{\eta x, q1}$ obtain all quorums and then instant x is updated at t_{17} . At t_{18} , the relation S is fragmented into S_1 and S_2 using vertical fragmentation. At t_{19} , the relation S_1 fragmented using horizontal fragmentation into $S_{1(Pk,x)}$ and $S_{1(Pk,y)}$. Finally, at t_{20} , $\hat{V}_{\lambda x, q1} \in V_{\eta}$ is commit and at t_{21} , instant x at all replica servers will unlock and ready for next transaction to take place.

4 Conclusion

It is important to manage a transaction in order to preserve data consistency and reliability of the systems [15,16]. Due to this, we have designed a new model called Binary Vote Assignment on Grid Quorum (BVAGQ). BVAGQ's goals are to manage fragmented database replication and transaction. From the experiment result, it is shown that the system achieve to manage replication and transaction through proposed BVAGQ. It is also shown that the system preserves the data consistency through a synchronization approach for all replicated sites.

Acknowledgement. This research is supported by FRGS from Ministry of Higher Education of Malaysia, Vote. No RDU 100109.

References

1. Kun, R., Zhanhuai, L., Chao, W.: LBD RP: A Low-bandwidth Data Replication Protocol on Journal-based Application. In: The Proceeding of 2nd International Conference on Computer Engineering and Technology (ICCET), vol. 80(12), pp. 1489–1498 (2010)
2. Deris, M.M., Abawajy, J.H., Taniar, D., Mamat, A.: Managing data using neighbour replication on a triangular-grid structure. International Journal of High Performance Computing and Networking 6(1), 56–65 (2009)
3. Kemme, B., Alonso, G.: A suite of database replication protocols based on group communication primitives. In: The Proceeding of 18th International Conference on Distributed Computing Systems, pp. 156–163 (1998)

4. Fauzi, A.A.C., Noraziah, A., Zain, N.M., Beg, A.H., Khan, N., Elrasheed, I.S.: Handling Fragmented Database Replication through Binary Vote Assignment Grid Quorum. *Journal of Computer Science* 7(9), 1338–1342 (2011)
5. Distributed Database, <http://www.scribd.com/doc/6142386/-distributed-database> (referred on March 2, 2011)
6. Noraziah, A., Ahmed, N.A., Sidek, R.M.: Data Replication Using Read-One-Write-All-Monitoring Synchronization Transaction Systems in Distributed Environment. *Journal of Computer Science* 6(10), 1033–1036 (2010)
7. Noraziah, A., Fauzi, A.A.C., Zain, N.M., Beg, A.H.: Lowest Data Replication Storage of Binary Vote Assignment Data Grid. In: Zavoral, F., Yaghob, J., Pichappan, P., El-Qawasmeh, E. (eds.) *NDT 2010, Part II. CCIS*, vol. 88, pp. 466–473. Springer, Heidelberg (2010)
8. Noraziah, A., Fauzi, A.A.C., Deris, M.M., Saman, M.Y.M., Zain, N.M., Khan, N.: Managing Educational Resource - Student Information Systems Using BVAGQ Fragmented Database Replication Model. *Procedia Social and Behavioral Sciences* 28, 127–132 (2011)
9. Noraziah, A., Sidek, R.M., Klaib, M.F.J., Jayan, T.L.: A Novel Algorithm of Managing Replication and Transaction through Read-One-Write-All Monitoring Synchronization Transaction System (ROWA-MSTS). In: *Proceeding of the 2nd International Conference on Network Applications, Protocols & Services (NETAPPS 2010)*, pp. 20–25 (2010)
10. Perez, J.M., García-Carballeira, F., Carretero, J., Calderón, A., Fernández, J.: Branch Replication Scheme: A New Model for Data Replication in Large Scale Data Grids. *Future Generation Computer Systems* 26(1), 12–20 (2010)
11. Sung, C.C., Hee, Y.Y.: Dynamic hybrid replication effectively combining tree and grid topology. *The Journal of Supercomputing* 59(3), 1289–1311 (2010)
12. Storm, C., Theel, O.: A General Approach to Analyzing Quorum-Based Heterogeneous Dynamic Data Replication Schemes. In: Garg, V., Wattenhofer, R., Kothapalli, K. (eds.) *ICDCN 2009. LNCS*, vol. 5408, pp. 349–361. Springer, Heidelberg (2008)
13. Henry, K., Swanson, C., Xie, Q., Daudjee, K.: Efficient Hierarchical Quorums in Unstructured Peer-to-Peer Networks. In: Meersman, R., Dillon, T., Herrero, P. (eds.) *OTM 2009, Part I. LNCS*, vol. 5870, pp. 183–200. Springer, Heidelberg (2009)
14. Navathe, S.B., Ra, M.: Vertical partitioning for database design: a graphical algorithm. *SIGMOD Record* 14, 440–450 (1989)

WLAR-Viz: Weighted Least Association Rules Visualization

A. Noraziah¹, Zailani Abdullah², Tutut Herawan¹, and Mustafa Mat Deris³

¹ Faculty of Computer Systems and Software Engineering, Universiti Malaysia Pahang
Lebuhraya Tun Razak, 26300 Kuantan Pahang, Malaysia

² Department of Computer Science, Universiti Malaysia Terengganu, 21030 Kuala Terengganu,
Terengganu, Malaysia

³ Faculty of Computer Science and Information Technology, Universiti Tun Hussein Onn
Malaysia, Parit Raja, Batu Pahat 86400, Johor, Malaysia
{noraziah, tutut, zailania}@ump.edu.my,
mmustafa@uthm.edu.my

Abstract. Mining weighted least association rules has been an increasing demand in data mining research. However, mining these types of rules often facing with difficulties especially in identifying which rules are really interesting. One of the alternative solutions is by applying the visualization model in those particular rules. In this paper, a model for visualizing weighted least association rules is proposed. The proposed model contains five main steps, including scanning dataset, constructing Least Pattern Tree (LP-Tree), applying Weighted Support Association Rules (WSAR*), capturing Weighted Least Association Rules (WELAR) and finally visualizing the respective rules. The results show that by using a three dimensional plots provide user friendly navigation to understand the weighted support and weighted least association rules.

Keywords: Weighted least association rules, Data mining, Visualization.

1 Introduction

In the past decades, mining association rules or patterns from transaction database has attracted many research interests. The aim of mining association rules is to uncover all interesting and useful patterns that are presented in data repositories. It was first introduced by Agrawal *et al.* [1] and still attracts many attentions from knowledge discovery community [2,3,4,5,6]. In association rules, a set of item is defined as an itemset. The itemset is said to be frequent, if it occurs more than a predefined minimum support. Besides that, confidence is another alternative measurement used in pair in association rules. The association rule is said to be strong if it meets the minimum confidence. In contradiction with the previous itemset, least itemset is a set of item that is infrequently found in the database. However, it may produce an interesting result for certain domain applications such as to detect the air pollution [7], serious diseases [8], educational decision support [9,10,11,12] and many more. Normally, the least itemset can be only captured by lowering the minimum support threshold. As a result, this approach may produce the enormous number of association

and it is enormously difficult to identify which association rules are most significant. Furthermore, the lowering the minimum support will also proportionally increase the computational performance in generating the complete set of association rules.

In our previous work, we have proposed the Weighted Least Association Rules framework (WELAR-f) in [13] to extract the significant association rules. Basically, the WELAR framework contains an enhanced version of existing prefix tree and frequent pattern growth algorithm called LP-Tree and LP-Growth algorithm, respectively. Moreover, Weighted Support Association Rules (WSAR*) measurement is also suggested in [13]. We have shown that by modifying this framework into suitable visualization model, the significant rules can be captured and visualized. In this paper, significant rules based on Breast-Cancer Wisconsin and Mushroom datasets are finely presented and 3-Dimensionally visualized.

The rest of the paper is organized as follows. Section 2 describes the related works. Section 3 explains in details the proposed methods. This is followed by experimental results in section 4. Finally, conclusion and future direction are reported in section 6.

2 Related Works

Association rules visualization is one of the exciting subset in association rules. Its main objective is to display data that can facilitate and comprehend the user interpretation. Until this recent, many authors have come forward to develop visualization techniques to support them in analyzing and comprehensively viewing the association rules.

Wong et al. [14] used 3-Dimensional method to visualize association rules for text mining. Bruzzese and Buono [15] presented a visual strategy to analyze huge rules by exploiting graph-based technique and parallel coordinates to visualize the results of association rules mining algorithms. Ceglar et al. [16] reviewed the current association visualization techniques and introduced a new technique for visualizing hierarchical association rules. Kopanakis et al. [17] proposed 3-Dimensional methods of visual data mining technique for the representation and mining of classification outcomes and association rules. Lopes et al. [18] presented a framework for visual text mining to support exploration of both general structure and relevant topics within a textual document collection. Leung et al. [19,20], developed a visualizer technique for frequent pattern mining. Later, Herawan and Deris employed soft set theory for mining maximal association rules and visualized them. Besides that, Abdullah et al. visualized the Construction of Incremental Disorder Trie Itemset Data Structure (DOSTrieIT) for Frequent Pattern Tree (FP-Tree).

Agrawal et al. was the first introduced the support-and-confidence measurement for evaluation and classification of association rules. However, this measurement is not covered weighted items. Cai et al. [25] introduced Weighted Association Rules (WAR) with MINWAL(O) and MINWAL(W) algorithms based on the support bounds approach to mine the weighted binary ARs. It can be considered as among the first attempt to allow the single item to carry the weight rather than 0 or 1. Tao et al. [26] proposed an algorithm namely Weighted Association Rule Mining (WARM) to discover significant weight of itemset. In summary, the study of weighted association rules is still very limited and worth to explore.

3 Proposed Method

Throughout this section the set $I = \{i_1, i_2, \dots, i_{|A|}\}$, for $|A| > 0$ refers to the set of literals called set of items and the set $D = \{t_1, t_2, \dots, t_{|U|}\}$, for $|U| > 0$ refers to the data set of transactions, where each transaction $t \in D$ is a list of distinct items $t = \{i_1, i_2, \dots, i_{|M|}\}$, $1 \leq |M| \leq |A|$ and each transaction can be identified by a distinct identifier TID.

3.1 Definition

Definition 1. (Least Items). An itemset X is called least item if $\alpha \leq \text{supp}(X) \leq \beta$, where α and β is the lowest and highest support, respectively. The set of least item will be denoted as Least Items and

$$\text{Least Items} = \{X \subset I \mid \alpha \leq \text{supp}(X) \leq \beta\}$$

Definition 2. (Frequent Items). An itemset X is called frequent item if $\text{supp}(X) > \beta$, where β is the highest support.

The set of frequent item will be denoted as Frequent Items and

$$\text{Frequent Items} = \{X \subset I \mid \text{supp}(X) > \beta\}$$

Definition 3. (Merge Least and Frequent Items). An itemset X is called least frequent items if $\text{supp}(X) \geq \alpha$, where α is the lowest support.

The set of merging least and frequent item will be denoted as LeastFrequent Items and

$$\text{LeastFrequent Items} = \{X \subset I \mid \text{supp}(X) \geq \alpha\}$$

LeastFrequent Items will be sorted in descending order and it is denoted as

$$\begin{aligned} \text{LeastFrequent Items}^{\text{desc}} = & \\ & \left\{ \begin{array}{l} X_i \mid \text{supp}(X_i) \geq \text{supp}(X_j), 1 \leq i, j \leq k, i \neq j, \\ k = |\text{LeastFrequent Items}|, x_i, x_j \subset \text{LeastFrequent Items} \end{array} \right\} \end{aligned}$$

Definition 4. (Ordered Items Transaction). An ordered items transaction is a transaction which the items are sorted in descending order of its support and denoted as t_i^{desc} , where $t_i^{\text{desc}} = \text{LeastFrequentItems}^{\text{desc}} \cap t_i, 1 \leq i \leq n, |t_i^{\text{least}}| > 0, |t_i^{\text{frequent}}| > 0$.

An ordered items transaction will be used in constructing the proposed model, so-called LP-Tree.

Definition 5. (Significant Least Data). Significant least data is one which its occurrence less than the standard minimum support but appears together in high proportion with the certain data.

Definition 6. (Item Weight). A weight of an item is defined as a non negative real number and it denoted as $\text{Item Weight} = \{X \subset I \mid 0 \leq \text{weight}(X) \leq 1\}$

Definition 7. (Itemset Length). A weight of an item is defined as a non negative real number and it denoted as $\text{Itemset Length} = \{X \subset I \mid 0 \leq \text{weight}(X) \leq 1\}$

Definition 8. (Weighted Support Association Rules). A Weighted Support Association Rules (WSAR*) is a weight of itemset by formulating the combination of the support and weight of item, together with the total number of support in either of them.

The value of Weighted Support Association Rules denoted as WSAR* and

$$\text{WSAR}^*(I) = \frac{((\text{supp}(A) \times \text{weight}(A)) + (\text{supp}(B) \times \text{weight}(B)))}{(\text{supp}(A) + \text{supp}(B) - \text{supp}(A \Rightarrow B))}$$

WSAR* value is determined by multiplying the summation of items weight from both antecedent and consequence, with the support of the itemset.

3.2 Model

There are five major components involved in visualizing the weighted least association rules (WELAR). All these components are closely interrelated and the process flow is moving in one-way direction. A complete overview model of visualizing the critical least association rules is shown in Figure 1.

Dataset. All datasets used in this model are in a flat file format. Each record (or transaction) is written in a line in the file and stored separately from others. The flat file takes up much less space than the structure file.

LP-Tree. The construction of Least Pattern Tree (LP-Tree) structure is based on the support descending orders of items. After that, several iterations will take place at LP-Tree to mine and generate the desired patterns or least association rules.

Weighted Support Association Rules (WSAR*). WSAR* for each association rule is computed. The items weight are assigned randomly (also can be fixed) in a range of 0.1 and 1.0. The antecedent, consequence and antecedent-consequence are utilized in calculating the WSAR*.

Weighted Least Association Rules (WELAR). Only association rules that have WSAR* value equal or more than predefined minimum WSAR* are classified as significant rules. These rules consist the combination of both least and frequent itemset.

Visualization. Weighted least association rules are presented in 3-D. Beside the value of WLAR*, the values from others measurements are also displayed for comparison and further analysis.

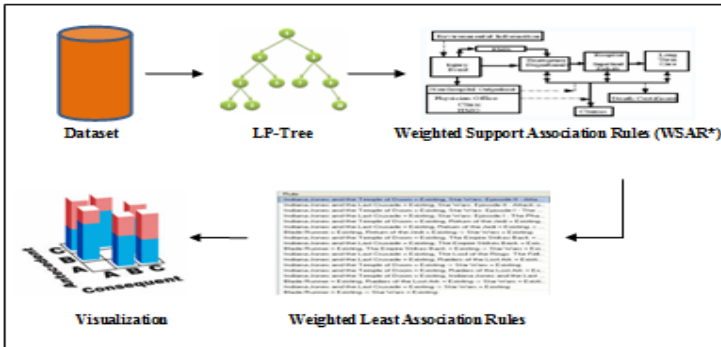


Fig. 1. WLAR visualization model

4 Experimental Results

In this section, we do comparative analysis of weighted least rules being generated using current weighted association rules, and the proposed measure, Weighted Support Association Rules (WSAR*). Abbreviations for each measurement are shown in Table 1.

Table 1. Abbreviations of different weighted measures

Measures	Description
CAR	Classical Association Rules
WAR	Weighted Association Rules
WSSAR	Weighted Support Significant Association Rules
WSAR*	Weighted Support Association Rules (proposed measure)

4.1 Breast-Cancer Wisconsin Dataset

The first experiment, we evaluate our proposed visualization model to Breast-Cancer-Wisconsin dataset from UCI Machine Learning Repository. The aim of the dataset is to diagnose the breast cancer according to Fine- Needle Aspirates (FNA) test. The dataset was obtained from a repository of a machine-learning database University of California, Irvin. It was compiled by Dr. William H. Wolberg from University of Wisconsin Hospitals, Madison, Madison, WI, United States. It has 11 attributes and 699 records (as of 15 July 1992) with 158 benign and 241 malignant classes, respectively. Figure 2 visualizes top 20 weighted least association rules in 3-Dimensional Bar Form.

4.2 Mushroom Dataset

In the second experiment, we evaluate our proposed visualization model to Mushroom dataset which is also from UCI Machine Learning Repository. It is a dense dataset and consists of 23 species of gilled mushroom in the Agaricus and Lepiota

Family. Table 12 shows the fundamental characteristics of the dataset. Figure 16 shows some portion of data taken from the Mushroom dataset. It has 23 attributes and 8,124 records. Fig. 3 visualizes top 26 weighted least association rules in 3-Dimensional Bar Form.

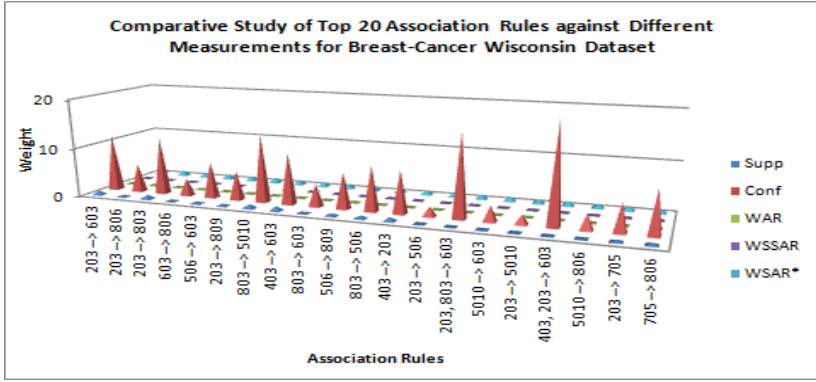


Fig. 2. A three dimensional bar form of visualizing weighted least association rules (WSAR*) of Breast-Cancer Wisconsin dataset

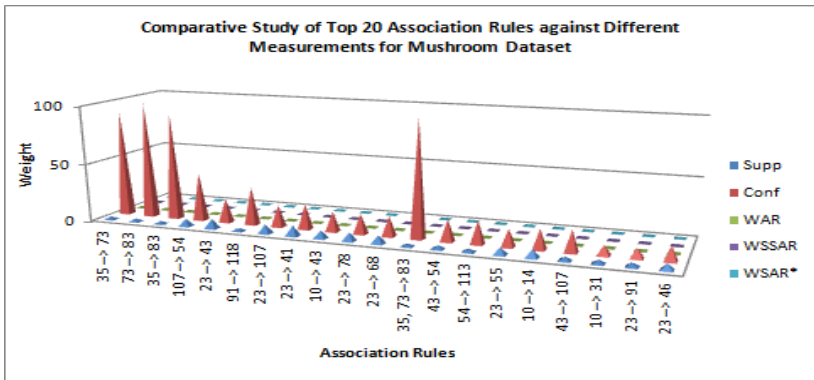


Fig. 3. A three dimensional bar form of visualizing weighted least association rules (WSAR*) of Mushroom dataset

5 Conclusion

The current approaches for visualizing association rules are still focusing on common rules. From our knowledge, no research has been carried out to visualize the weighted least association rules. Therefore in this paper, we proposed WLAR-Viz (*Weighted Least Association Rules Visualization*), an approach for visualizing the significant and weighted least association rules using Weighted Support Association Rules (WSAR*) measurement. We evaluate the proposed model through the benchmarked

Breast-Cancer Wisconsin and Mushroom dataset. The results show that using 3-Dimensional Bar form can provide useful analysis in comparing the different types of measurements in association rules. With this approach, we believe that our proposed approach can also be used in capturing weighted least association rules from other domain applications.

Acknowledgement. This research is supported by Fundamental Research Grant Scheme (FRGS) from Ministry of Higher Education of Malaysia, Vote. No RDU 100109.

References

1. Agrawal, R., Imielinski, T., Swami, A.: Database Mining: A Performance Perspective. *IEEE Transactions on Knowledge and Data Engineering* 5(6), 914–925 (1993)
2. Abdullah, Z., Herawan, T., Deris, M.M.: An Alternative Measure for Mining Weighted Least Association Rule and Its Framework. In: Zain, J.M., Wan Mohd, W.M.B., El-Qawasmeh, E. (eds.) ICSECS 2011, Part II. CCIS, vol. 180, pp. 480–494. Springer, Heidelberg (2011)
3. Herawan, T., Yanto, I.T.R., Deris, M.M.: Soft Set Approach for Maximal Association Rules Mining. In: Ślęzak, D., Kim, T.-H., Zhang, Y., Ma, J., Chung, K.-I. (eds.) DTA 2009. CCIS, vol. 64, pp. 163–170. Springer, Heidelberg (2009)
4. Abdullah, Z., Herawan, T., Deris, M.M.: Mining Significant Least Association Rules Using Fast SLP-Growth Algorithm. In: Kim, T.-H., Adeli, H. (eds.) AST/UCMA/ISA/ACN 2010. LNCS, vol. 6059, pp. 324–336. Springer, Heidelberg (2010)
5. Herawan, T., Deris, M.M.: A soft set approach for association rules mining. *Knowledge Based Systems* 24(1), 186–195 (2011)
6. Abdullah, Z., Herawan, T., Deris, M.M.: Tracing Significant Information using Critical Least Association Rules Model. To Appear in Special Issue of ICICA 2010, *International Journal of Innovative Computing and Applications* x(x), xxx–xxx (2012)
7. Mustafa, M.D., Nabila, N.F., Evans, D.J., Saman, M.Y., Mamat, A.: Association rules on significant rare data using second support. *International Journal of Computer Mathematics* 83(1), 69–80 (2006)
8. Abdullah, Z., Herawan, T., Deris, M.M.: Detecting Critical Least Association Rules in Medical Databases. *International Journal of Modern Physics: Conference Series* 9, 464–479 (2012)
9. Abdullah, Z., Herawan, T., Noraziah, A., Deris, M.M.: Extracting Highly Positive Association Rules from Students' Enrollment Data. *Procedia Social and Behavioral Sciences* 28, 107–111 (2011)
10. Abdullah, Z., Herawan, T., Noraziah, A., Deris, M.M.: Mining Significant Association Rules from Educational Data using Critical Relative Support Approach. *Procedia Social and Behavioral Sciences* 28, 97–101 (2011)
11. Herawan, T., Vitasari, P., Abdullah, Z.: Mining Interesting Association Rules of Student Suffering Mathematics Anxiety. In: Zain, J.M., Wan Mohd, W.M.B., El-Qawasmeh, E. (eds.) ICSECS 2011, Part II. CCIS, vol. 180, pp. 495–508. Springer, Heidelberg (2011)
12. Herawan, T., Vitasari, P., Abdullah, Z.: Mining Interesting Association Rules on Student Suffering Study Anxieties using SLP-Growth Algorithm. *International Journal of Knowledge and Systems Science* 3(2), 24–41 (2012)

13. Abdullah, Z., Herawan, T., Deris, M.M.: An Alternative Measure for Mining Weighted Least Association Rule and Its Framework. In: Zain, J.M., Wan Mohd, W.M.B., El-Qawasmeh, E. (eds.) ICSECS 2011, Part II. CCIS, vol. 180, pp. 480–494. Springer, Heidelberg (2011)
14. Wong, P.C., Whitney, P., Thomas, J.: Visualizing Association Rules for Text Mining. In: Proceedings 1999 IEEE Symposium on Information Visualization (Info Vis 1999), pp. 120–123 (1999)
15. Bruzzese, D., Buono, P.: Combining Visual Techniques for Association Rules Exploration. In: Proceedings of the Working Conference on Advanced Visual Interfaces (AVI 2004), pp. 381–384. ACM Press (2004)
16. Ceglar, A., Roddick, J., Calder, P., Rainsford, C.: Visualizing hierarchical associations. *Knowledge and Information Systems* 8, 257–275 (2005)
17. Kopanakis, I., Pelekis, N., Karanikas, H., Mavroudkis, T.: Visual Techniques for the Interpretation of Data Mining Outcomes. In: Bozaris, P., Houstis, E.N. (eds.) PCI 2005. LNCS, vol. 3746, pp. 25–35. Springer, Heidelberg (2005)
18. Lopes, A.A., Pinho, R., Paulovich, F.V., Minghim, R.: Visual text mining using association rules. *Computers & Graphics* 31, 316–326 (2007)
19. Leung, C.K.S., Irani, P., Carmichael, C.L.: WiFiViz: Effective Visualization of Frequent Itemsets. In: Proceeding of Eighth IEEE International Conference on Data Mining (ICDM 2008), pp. 875–880. IEEE Press (2008)
20. Leung, C.K.-S., Irani, P.P., Carmichael, C.L.: FlsViz: A Frequent Itemset Visualizer. In: Washio, T., Suzuki, E., Ting, K.M., Inokuchi, A. (eds.) PAKDD 2008. LNCS (LNAI), vol. 5012, pp. 644–652. Springer, Heidelberg (2008)

Ontology-Based Genes Similarity Calculation with TF-IDF

Yue Huang¹, Mingxin Gan^{1,*}, and Rui Jiang²

¹ School of Economics and Management, University of Science and Technology Beijing,
Beijing 100083, P.R. China

² Bioinformatics Division, TNLIST/Department of Automation, Tsinghua University, Beijing
100084, P.R. China

mylinnia@126.com, ganmx@ustb.edu.cn, ruijiang@tsinghua.edu.cn

Abstract. The Gene Ontology (GO) provides a controlled vocabulary of terms for describing genes from different data resources. In this paper, we proposed a novel method determining semantic similarity of genes based on GO. The key principle of our method relies on the introduction of Term Frequency (TF) and Inverse Document Frequency (IDF) to quantify the weights of different GO terms to the same gene. Different from previous leading methods, our method needs no parameters and computes the gene similarity directly rather than term similarity first. Experimental results of clustering genes in biological pathways from *Saccharomyces Genome Database* (SGD) have demonstrated that our method is quite competitive and outperforms leading method in certain cases.

Keywords: Gene Ontology (GO), gene semantic similarity, TF-IDF, biological pathways.

1 Introduction

The Gene Ontology aims to represent gene products from heterogeneous databases with a standardized vocabulary of terms [1]. The GO comprises three orthogonal ontologies: biological process, molecular function, and cellular component. They are represented as directed acyclic graphs (DAGs) in which nodes correspond to terms and their relationships are represented by edges [2]. Each GO term is related to one or more parent terms by relationships and the most common types of relationships are "is-a" and "part-of". Calculating the semantic similarity of genes with annotated GO terms is an important application of GO, since they can explicitly reflect the closeness in meaning between two genes [3], and there are some available tools online, such as FuSSiMeG [4], G-SESAME [5], and GOSemSim [6]. Furthermore, comparison of function-unknown genes with function-known genes can be of reference for biologists to determine the function of unknown genes.

Although investigating similarity measures based on GO is recent, several gene similarity measures have been proposed in the past decade. Semantic similarity can be

* Corresponding author.

defined for both the GO terms and genes. Nearly all gene similarity measurements calculate the semantic similarity of GO terms first, and then compute the similarity of genes based on the similarities of GO terms annotating them. Obviously, when only comparison of genes is desired, this way of computing adds complexity.

Existing methods of ontology-based semantic similarity calculation for GO terms are usually divided into three categories: information-content based methods, which determine the similarity between two concepts by the extent to which they share information in common [7]; ontology-structure based methods, in which structural factors such as shortest path, network density, node depth, type of link, link strength, locations in the ontology graph are taken into account; and a hybrid of the two. Information-content based methods are based on the assumption that the more the information two concepts share in common, the more similar they are. Typical methods in this group include Resnik [8], Jiang [9], Lin [10], Schlicker [2], and their variants. These methods prove to be useful in tasks of natural language processing, and are originally designed for WordNet [11]. However, these approaches highly depend on the annotation data being used, while ignore the information contained in the structure of the ontology [12], and are only suitable for "is-a" relationship in the ontology. Whereas ontology-structure based methods consider not only the number of common ancestor terms but also locations of these ancestor terms related to these two specific terms in ontology graph [12]. These methods include Wang [12], Zhang [13], DOPCA [14], and etc. However, it is not a trivial task to determine parameters introduced by these methods: the two different semantic contribution factors according to relationship types in Wang [12]; the two factors for tuning the contribution of the length of the shortest path and the depth into similarity function in Zhang [13]. Thus, it is significant that a similarity method is immune to parameters. There are also hybrid methods that take into consideration not only IC but also the structure of GO, such as Combine [15] and SP [16].

Based on the semantic similarity of GO terms computed by the above methods, some methods define gene similarity as the maximum or average value between all the pairs of gene annotation terms [17]. Others represent genes in vector space using vector similarity measures to calculate functional similarity, but they couldn't avoid zero result, which will be discussed later in Section 2. Still others simply use the number of GO terms two genes share as the functional similarity [18], the outcomes of which are not quite competitive.

In this paper, we proposed a method measuring genes similarity based on annotation data. To overcome the drawback of other methods mentioned above, Term Frequency (TF) and Inverse Document Frequency (IDF) which are popular employed in the field of information retrieval are employed to calculate the role of different GO terms annotating a certain gene. Our method needs no parameters and directly computes the semantic similarity of genes, which has been proved to be impressive compared with the recognized Wang method [12] by evaluating clustering genes in biological pathways from Saccharomyces Genome Database (SGD).

The rest of this paper is organized as follows: Section 2 represents our novel method for assessing semantic similarity over genes. Some experimental evaluation is discussed in Section 3. We finish the paper with a conclusion in Section 4.

2 The Proposed Measure

2.1 Background

Given a gene G , it can be represented as a n -dimensional vector which consists of GO terms $G=\{t_1, t_2, \dots, t_n\}$, where t_i ($i=1, 2, \dots, n$) denotes a GO term, and n is the total number of GO terms in GO. Traditionally, each gene vector is binary valued, with 1 representing the presence of the GO term in the gene's annotation and 0 representing its absence. The most common way to compute the similarity between genes is to calculate cosine value of two genes G_1 and G_2 as follows [17]:

$$sim_{\cos}(G_1, G_2) = \frac{G_1 \cdot G_2}{|G_1||G_2|} \quad (1)$$

A variation on the cosine measure is to replace the non-zero values in the binary vector with scaled values, i.e. weight w_i which denotes each GO term based on the frequency of its occurrence in the corpus as follows [19]:

$$w_i = \log\left(\frac{N}{n_i}\right) \quad (2)$$

where N is the total number of genes in the corpus and n_i is the number of genes in the corpus annotated with that term t .

Both the cosine measure and weighted cosine measure face the problem that they compute a zero result when two genes share no direct GO term annotations, which is apparently not ideal for semantic dissimilarity measurement.

2.2 Definitions

The true path rule every GO term must obey indicates that a GO term's all parent terms contribute to the term's semantics [12]. Furthermore, each GO term can be represented as a DAG, which starts from the specific term and ends at any of the three gene ontologies. To demonstrate the contribution of all terms in a DAG, we define the extended version (EV) of a term t as the terms set from t 's DAG. For any term t in GO, $EV(t)$ is defined as:

$$EV(t) = \{t' \mid t' \in DAG(t)\} \quad (3)$$

where $DAG(t)$ denotes the set of GO terms including all of t 's ancestor terms and itself. For instance, Fig. 1 shows the DAG for GO term cell-cell adhesion (GO: 0016337) whose ancestor terms include GO: 0008150, GO: 0009987, GO: 0022610, and GO: 0007155. Then the extended version of term t (GO: 0016337) is $EV(t)=\{GO: 0008150, GO: 0009987, GO: 0022610, GO: 0007155, GO: 0016337\}$.

For one specific gene, if all of its annotation GO terms are represented by their EV form, then one GO term perhaps occurs more than once in the gene's EV set. We

define the extended annotation (*EA*) of a gene *G* as a set of the gene’s annotation GO terms with their frequency in this set. For any gene *G*, *EA*(*G*) is defined as:

$$EA(G) = \{(t, freq) \mid t \in G_a\} \tag{4}$$

where *G_a* denotes the annotation set of terms in *G* with their *EV* forms, and *freq* denotes the frequency of term *t* in *G_a*. For instance, see Fig. 1, if we assume gene *G* is annotated by GO: 0009987 and GO: 0022610, then its extended annotation could be *EA*(*G*)={(*GO*: 0008150,2),(*GO*: 0009987,1),(*GO*: 0022610,1)}.

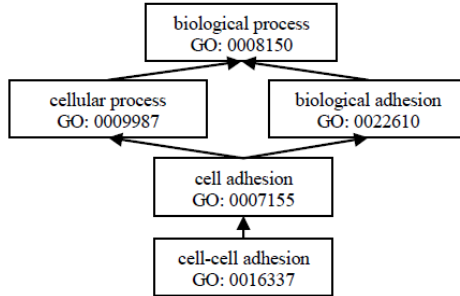


Fig. 1. DAG for GO term ‘cell-cell adhesion’ (GO: 0016337)

2.3 Gene Semantic Similarity

Let $G = \{G_1, G_2, \dots, G_N\}$ be a set of genes and $T = \{t_1, t_2, \dots, t_n\}$ be a set of distinct GO terms annotating *G*, where *N* and *n* are the total number of genes and terms, respectively. We use TF to estimate the importance of a GO term to a specific gene: the more a GO term appears in a gene, the more it is estimated to be important in this gene. In addition, we use IDF to measure the role of a GO term to the whole gene database: the frequency of a GO term in the gene database and its importance are negatively correlated. There are many variants of TF-IDF [20], and we compute the weight in the following way:

$$w_{t,G} = TF_{t,G} \cdot IDF_t = \frac{f_{t,G}}{\max_z f_{t,G}} \ln\left(\frac{N}{n_t} + 0.01\right) \tag{5}$$

where *f_{t,G}* is the frequency of GO term *t* in *EA*(*G*) obtained with equation (4), *z* denotes the number of GO terms annotating *G*, *N* is the number of genes in a specific gene database, and *n_t* is the number of genes where term *t* occurs.

Formally, a gene is represented as a scaled valued vector $G = \{w_1, w_2, \dots, w_n\}$, where *w_i* (*i*=1, 2, ..., *n*) describes the weight of GO term *t_i* to the gene *G* and 0 representing its absence. Finally, we use the cosine value to calculate the semantic similarity of genes:

$$sim(G_1, G_2) = \frac{\sum_{i=1}^n w_{i,G_1} w_{i,G_2}}{\sqrt{\sum_{i=1}^n (w_{i,G_1})^2} \sqrt{\sum_{i=1}^n (w_{i,G_2})^2}} \quad (6)$$

where n denotes the total number of terms in GO. The range of the similarity value obtained by our method is between 0 and 1. The higher the value obtained the more similar between them.

3 Evaluation

The evaluation methods for gene functional similarity measurements include three categories [13]: comparing the similarity values calculated by the measure with gene sequence similarity, comparing with gene expression profile, and using gene pathways and clusters information to validate the results. In our paper, we used the third approach and compared our method with the leading method Wang [12]. We implemented our method in C#, and used G-SESAME [5] to compute similarity as Wang proposed. We extracted 104 biological pathways which contained at least three genes from SGD. Clustering results of genes from these pathways proved that our method was competitive and outperformed Wang in certain cases. We take pathway "glyoxylate cycle" (Table 1) as an example to show the performance of our method.

Table 1. Functions of genes in glyoxylate cycle pathway

Id	E.C number of reaction	Gene name
1	EC:4.2.1.3	ACO1
	EC:4.2.1.3	ACO2
2	EC:2.3.3.1	CIT1
	EC:2.3.3.1	CIT2
	EC:2.3.3.1	CIT3
3	EC:4.1.3.1	ICL1
4	EC:1.1.1.37	MDH1
	EC:1.1.1.37	MDH2
	EC:1.1.1.37	MDH3
5	EC:2.3.3.9	DAL7
	EC:2.3.3.9	MLS1

In the field of biology, if different genes participate in the same reaction, then these genes perhaps have the same biological function. Thus if the outcomes of clustering genes based on the similarity calculated by a gene similarity method are consistent with the functional categorization suggested by biological reactions, then it proves the effectiveness of the gene similarity method indirectly. Table 1 shows that there are 5 reactions and 11 genes in glyoxylate cycle pathway: ACO1 and ACO2 should be assigned into one cluster; CIT1, CIT2, and CIT3 should be in another; and so on.

	ACO1	ACO2	CIT1	CIT2	CIT3	ICL1	MDH1	MDH2	MDH3	DAL7	MLS1
ACO1	1.000	0.763	0.510	0.476	0.298	0.436	0.250	0.197	0.258	0.238	0.234
ACO2		1.000	0.269	0.177	0.173	0.371	0.222	0.178	0.163	0.202	0.195
CIT1			1.000	0.868	0.642	0.514	0.376	0.332	0.364	0.494	0.474
CIT2				1.000	0.634	0.595	0.316	0.302	0.464	0.597	0.603
CIT3					1.000	0.406	0.270	0.241	0.343	0.440	0.422
ICL1						1.000	0.382	0.345	0.491	0.522	0.503
MDH1							1.000	0.760	0.688	0.296	0.293
MDH2								1.000	0.650	0.264	0.252
MDH3									1.000	0.457	0.550
DAL7										1.000	0.783
MLS1											1.000

Fig. 2. Similarities among genes in glyoxylate cycle pathway obtained by our method

See Fig. 2 and Fig. 3, our method puts ACO1 and ACO2 together in the third clustering step, since the similarity 0.763 means they are very similar to each other. However, Wang doesn't put them together until the fifth step, since the similarity is only 0.544. While the descriptions of ACO1 and ACO2 from SGD clearly depict that ACO2 is similar to ACO1 and they are both aconitase required for the TCA cycle, but ACO1 is also independently required for mitochondrial genome maintenance. Thus the similarity of ACO1 and ACO2 by our method is consistent with their functions, whereas their similarity obtained by Wang is not consistent with their functions.

	ACO1	ACO2	CIT1	CIT2	CIT3	ICL1	MDH1	MDH2	MDH3	DAL7	MLS1
ACO1	1.000	0.544	0.144	0.144	0.155	0.213	0.311	0.130	0.311	0.144	0.144
ACO2		1.000	0.199	0.199	0.199	0.338	0.155	0.180	0.155	0.199	0.199
CIT1			1.000	1.000	0.909	0.199	0.155	0.180	0.155	0.728	0.728
CIT2				1.000	0.909	0.199	0.155	0.180	0.155	0.728	0.728
CIT3					1.000	0.199	0.161	0.180	0.161	0.728	0.728
ICL1						1.000	0.155	0.180	0.155	0.199	0.199
MDH1							1.000	0.697	1.000	0.155	0.155
MDH2								1.000	0.697	0.180	0.180
MDH3									1.000	0.155	0.155
DAL7										1.000	1.000
MLS1											1.000

Fig. 3. Similarities among genes in glyoxylate cycle pathway obtained by Wang

Threshold	Initial	0.868	0.783	0.763	0.760	0.650	0.634	0.503	0.406	0.241	0.163
Clustering Result	ACO1	ACO1	ACO1	ACO1	ACO1	ACO1	ACO1	ACO1	ACO1	ACO1	ACO1
	ACO2	ACO2	ACO2	ACO2	ACO2	ACO2	ACO2	ACO2	ACO1	ACO2	ACO1
	CIT1										ACO1
	CIT2	CIT1	CIT1	CIT1	CIT1	CIT1	CIT1	CIT1	CIT1		ACO2
	CIT3	CIT2	CIT2	CIT2	CIT2	CIT2	CIT2	CIT2	CIT1		
		CIT3	CIT3	CIT3	CIT3	CIT3	CIT3		CIT2		
									CIT3		
		ICL1	ICL1	ICL1	ICL1	ICL1	ICL1	ICL1	ICL1		
									DAL7		
		DAL7	DAL7	DAL7	DAL7	DAL7	DAL7	DAL7	DAL7		
		MLS1	MLS1	MLS1	MLS1	MLS1	MLS1	MLS1	MLS1		
	MDH1	MDH1	MDH1	MDH1	MDH1	MDH1	MDH1	MDH1	MDH1	MDH1	
	MDH2	MDH2	MDH2	MDH2	MDH2	MDH2	MDH2	MDH2	MDH2	MDH2	
	MDH3	MDH3	MDH3	MDH3	MDH3	MDH3	MDH3	MDH3	MDH3	MDH3	

Fig. 4. Clustering results of genes in glyoxylate cycle pathway based on the similarity values obtained by our method

Threshold	Initial	1.000	0.909	0.728	0.697	0.544	0.338	0.311	0.199	
Clustering Result	ACO1	ACO1	ACO1	ACO1	ACO1	ACO1	ACO1	ACO1		
	ACO2	ACO2	ACO2	ACO2	ACO2	ACO2	ACO2	ACO1		
	ICL1	ICL1	ICL1	ICL1	ICL1	ICL1	ICL1	ACO2	ACO1	
	MDH1							ICL1	ACO2	
	MDH3	MDH1	MDH1	MDH1	MDH1	MDH1	MDH1	MDH1	ACO1	
	MDH2	MDH3	MDH3	MDH3	MDH3	MDH3	MDH3	MDH3	ACO2	
	CIT1							MDH2	ICL1	
	CIT2							MDH2	DAL7	
	CIT3							MDH2	MLS1	
	DAL7							MDH2	MDH1	
	MLS1							MDH2	MDH2	
		CIT1	CIT1	CIT1	CIT1	CIT1	CIT1	CIT1	CIT1	
		CIT2	CIT2	CIT2	CIT2	CIT2	CIT2	CIT2	CIT2	
		CIT3	CIT3	CIT3	CIT3	CIT3	CIT3	CIT3	CIT3	
		DAL7	DAL7	DAL7	DAL7	DAL7	DAL7	DAL7	DAL7	
		MLS1	MLS1	MLS1	MLS1	MLS1	MLS1	MLS1	MLS1	

Fig. 5. Clustering results of genes in glyoxylate cycle pathway based on the similarity values obtained by Wang

When the similarity threshold reaches 0.634, there are five clusters by our method and the clustering results are exactly the same as Table 1 suggests (Fig. 4). However Fig. 5 shows that Wang fails to identify the correct five clusters. This proves our method of gene similarity is successful and outperforms Wang in this pathway.

4 Conclusion

In this paper, we proposed a novel method calculating semantic similarity based on ontology between genes. First we extend each GO term of a gene annotation, which reflect the semantic contribution of a GO term’s ancestors to it. Then, we use a common variant of TF-IDF to quantify the importance of different GO terms in the same gene, which helps represent each gene in a vector space model. Finally, gene similarity is calculated as the cosine value of two gene vectors. Our method is easy to implement, needs no extra parameters, and computes the gene similarity directly. Experimental results have proved its competitiveness compared with leading method. It is worthwhile that the completeness of the ontology used is of importance in this calculation, and with the ongoing construction of GO the similarity obtained by our method can be of greater help to biologists.

Acknowledgments. This work was partly supported by the National Natural Science Foundation of China under Grants No. 71101010, 61175002 and 71172169, the Fundamental Research Funds for the Central Universities under Grant No. FRF-BR-11-019A, and TNLIST Cross Discipline Foundation.

References

1. Ashburner, M., Ball, C.A., Blake, J.A., Botstein, D., Butler, H., Cherry, J.M., et al.: Gene Ontology: Tool for the Unification of Biology. *Nat. Genet.* 25, 25–29 (2000)
2. Schlicker, A., Domingues, F.S., Rahnenführer, J., Lengauer, T.: A New Measure for Functional Similarity of Gene Products Based on Gene Ontology. *BMC Bioinformatics* 7, 302 (2006)

3. Pesquita, C., Faria, D., Falcão, A.O., Lord, P., Couto, F.M.: Semantic Similarity in Biomedical Ontologies. *PLoS Comput. Biol.* 5, e1000443 (2009)
4. Couto, F.M., Silva, M.J., Coutinho, P.: Implementation of a Functional Semantic Similarity Measure between Gene-Products. Technical report, Department of Informatics, University of Lisbon (2003)
5. Du, Z., Li, L., Chen, C.F., Yu, P.S., Wang, J.Z.: G-SESAME: Web Tools for GO-Term-Based Gene Similarity Analysis and Knowledge Discovery. *Nucleic Acids Res.* 37, W345–W349 (2009)
6. Yu, G., Li, F., Qin, Y., Bo, X., Wu, Y., Wang, S.: GOSemSim: an R Package for Measuring Semantic Similarity among GO Terms and Gene Products. *Bioinformatics* 26, 976–978 (2010)
7. Resnik, P.: Semantic Similarity in a Taxonomy: an Information-Based Measure and its Application to Problems of Ambiguity in Natural Language. *J. Artif. Intell. Res.* 11, 95–130 (1999)
8. Resnik, P.: Using Information Content to Evaluate Semantic Similarity in a Taxonomy. In: 14th International Joint Conference on Artificial Intelligence, pp. 448–453. Morgan Kaufmann Publishers, San Francisco (1995)
9. Jiang, J.J., Conrath, D.W.: Semantic Similarity Based on Corpus Statistics and Lexical Taxonomy. In: 10th International Conference Research on Computational Linguistics (ROCLING X), pp. 19–33. Scandinavian University Press, Taiwan (1997)
10. Lin, D.: An Information-Theoretic Definition of Similarity. In: 15th International Conference on Machine Learning, pp. 296–304. Morgan Kaufmann Publishers, California (1998)
11. Lord, P.W., Stevens, R.D., Brass, A., Goble, C.A.: Investigating Semantic Similarity Measures across the Gene Ontology: the Relationship between Sequence and Annotation. *Bioinformatics* 19, 1275–1283 (2003)
12. Wang, J.Z., Du, Z., Payattakool, R., Yu, P.S., Chen, C.F.: A New Method to Measure the Semantic Similarity of GO Terms. *Bioinformatics* 23, 1274–1281 (2007)
13. Zhang, S., Shang, X., Wang, M., Diao, J.: A New Measure Based on Gene Ontology for Semantic Similarity of Genes. In: International Conference on Information Engineering, pp. 85–88. IEEE Press, New York (2010)
14. Gan, M., Dou, X., Wang, D., Jiang, R.: DOPCA: A New Method for Calculating Ontology-Based Semantic Similarity. In: 10th IEEE/ACIS International Conference on Computer and Information Science, pp. 110–115. IEEE Press, New York (2011)
15. Li, R., Cao, S., Li, Y., Tan, H., Zhu, Y., Zhong, Y., et al.: A Measure of Semantic Similarity between Gene Ontology Terms Based on Semantic Pathway Covering. *Prog. Nat. Sci.* 16, 721–726 (2006)
16. Shen, Y., Zhang, S., Wong, H.S.: A New Method for Measuring the Semantic Similarity on Gene Ontology. In: IEEE International Conference on Bioinformatics and Biomedicine, pp. 533–538. IEEE Press, New York (2010)
17. Popescu, M., Keller, J.M., Mitchell, J.A.: Fuzzy Measures on the Gene Ontology for Gene Product Similarity. *IEEE/ACM Trans. on Comput. Biol. and Bioinfo.* 3, 263–274 (2006)
18. Lee, H.K., Hsu, A.K., Sajdak, J., Qin, J., Pavlidis, P.: Coexpression Analysis of Human Genes across Many Microarray Data Sets. *Genome Res.* 14, 1085–1094 (2004)
19. Mistry, M., Pavlidis, P.: Gene Ontology Term Overlap as a Measure of Gene Functional Similarity. *BMC Bioinformatics* 9, 327 (2008)
20. Soucy, P., Mineau, G.W.: Beyond TFIDF Weighting for Text Categorization in the Vector Space Model. In: 19th International Joint Conference on Artificial Intelligence, pp. 1130–1135. Morgan Kaufmann Publishers, California (2005)

Quantitative Study of Oilfield Casing Damage

Deng Rui^{1,2}, Zhang Liang³, and Guo Haimin^{1,2,*}

¹ Key Laboratory of Exploration Technologies for Oil and Gas Resources, Ministry of Education (Yangtze University), Jingzhou City, Nanhuan Road 1#, 434023 Hubei, China

² School of Geophysics and Oil Resources, Yangtze University, Jingzhou City, Nanhuan Road 1#, 434023 Hubei, China

Zhonghai Petroleum (China) Co., Ltd. Zhanjiang Branch, Zhanjiang City, Nandiao Road, 524057, China

{guoHaimin,ghm3819}@163.net

Abstract. Carrying on the analysis to the cause of casing failure based on engineering factors and geological factors, and choosing the improved analytic hierarchy process to have a quantitative study for oilfield of the casing failure, improved the influence of quantitative study to casing failure when subjective factors been existent, it's extremely effective to determine the weights of each affecting factor. Establishing the analytical model of casing failure based on the improved ahp analysis, and analyzing the casing failure situation quantitatively, then comparing to the actual damage datas, the accuracy rate is above 85%, which shows that it can be effectively evaluated to casing failure by this method.

Keywords: Casing damage, Analytical hierarchy process (AHP), Quantitative analysis.

1 Introduction

With the continuous development of oil fields, the serious chip casing damage phenomenon occurred one after another which seriously affected the normal production of oilfield. To this end, a series of causes and mechanisms research of casing damage has launched, the factors led to casing damagement are generalized seven kinds of geological factors [1, 2] five kinds of development factors [3] and corrosion factors. Overseas, the oil casing damage is divided into two types of strata cut casing and special strata (gypsum strata or rock salt) casing collapse on the mechanical mechanism [4] Casing damage related to many factors, such as the geology, development, engineering and oil field chemistry, which are interrelated constraints to form a complex system, so now the oilfield casing damage research is still on the basis of the qualitative research.

Determining the weight of oilfield casing damage is the key to quantitative research. However, there has not been a good way to determine the weights over the years. At present, the commonly used method in general is Delphi method, which is based on the knowledge, wisdom, experience, information and values of number of experts to analyze, judge, weighed and give the corresponding weight to the indicators prepared [5]. Since

* Corresponding author.

Delphi method is mainly built on the experts' subjective judgments, there inevitably are some anthropogenic factors in scoring because of the different professional background, experience and the familiarity to every evaluation indicator of expert team member. In addition, a number of evaluation indicators are difficult to grade quantitatively and there aren't appropriate test conditions, which affect the results of the evaluation. In addition to this method, there are α -method, the method using ideal point to determine the weights and entropy method [6]. The advantages are their strong theory, which overcomes the subjectivity of Delphi method. But it is required quantitative evaluation indicators, while there are many factors that are qualitative in many practical problems, which are the shortcomings of these two ways. Compared with many kinds of methods, Analytic Hierarchy Process (AHP) is thought to be an effective method to determine the weight of various impact factors in oilfield casing damage.

2 AHP to Determine the Weight in the Quantitative Study of Casing Damage

The most important feature of AHP method is that on the basis of further analysis of the nature, external and its internal links of the complex practical problems, with less quantitative information realize the mathematization of decision-making thinking process to provide simple evaluation method for the complex decision-making issues with multiple objectives, criteria or without structural characteristics. It avoids that the predicted weight is incompatible with the actual situation due to the subjectivity of people and overcomes the problem that the decision maker and the decision analyzers are difficult to communicate. Therefore, AHP method is suited to the quantitative research of oilfield where the impact factors are difficult to accurately calculate. The core idea of AHP method is to construct a matrix meeting the consistent test, then apply the eigenvalue method to get the largest eigenvalues of the paired comparison matrix, and last to obtain the normalized feature vectors corresponding to the largest eigenvalues, which is the weight vector.

2.1 The Establishment of a Hierarchical Structure Model

On the depth analysis of the formation mechanism of oil field casing damage, decompose various factors into a number of top-down hierarchy according to their different properties. Factors of the same layer depend or impact on the factors of the upper layer, at the same time, they control the factors of the under layer or are impacted by the factors of the under layer. The top layer is the target layer here for the quantitative study of casing damage. There are one or several levels in the middle, often as the standard layer, which include the geological factors, the development factors and engineering factors in this model. The undermost layer is usually the program level or the object layer. The hierarchical structure model of the quantitative study of oilfield casing damage is shown in Fig. 1.

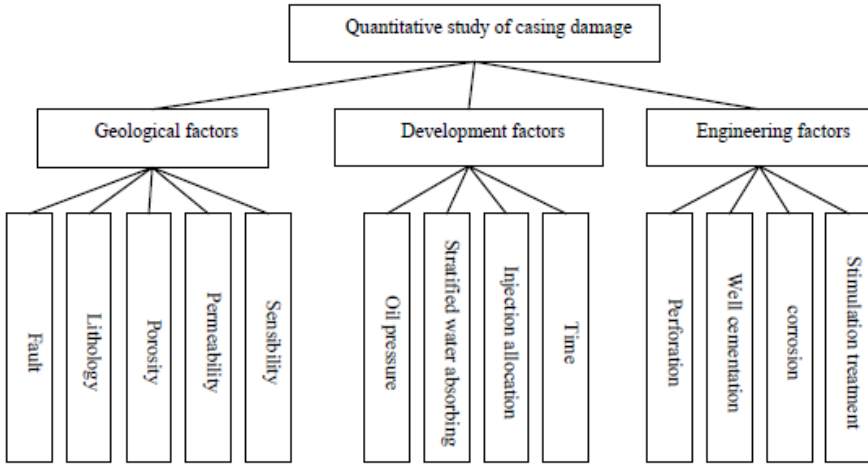


Fig. 1. The hierarchical structure model of quantitative study of oilfield casing damage

2.2 The Establishment of Comparative Matrix Pairs

After the establishment of a hierarchical structure, the subordinate relationship of the elements between the upper and lower layers has been determined. Compare paired the importance of the elements of each structural level. The traditional AHP method use 1-9 scale method to obtain judgment matrix by pair-wise comparison of each element [7].

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix}$$

Where: (1) $a_{ij} > 0$; (2) $a_{ii} = 1$; (3) $a_{ij} = 1/a_{ji}$; (4) $a_{ij} = a_{ik} \times a_{kj}$;

The matrix meeting the four characteristics is known as the "consistency" matrix and the only non-zero characteristic root of the matrix A is n [8].

2.3 Calculation of Eigenvalues and Eigenvectors

According to the special nature of judgment matrix the square root method can be used to calculate. First calculate the product M_i of elements of each line of the judgment matrix, then calculate the nth root \bar{W}_i of M_i , and finally normalization process is done to get the feature vector. Namely:

$$M_i = \prod_{j=1}^n b_{ij} \Rightarrow \bar{W}_i = \sqrt[n]{M_i} \Rightarrow W_i = \bar{W}_i / \sum_{j=1}^n \bar{W}_j \Rightarrow W = [W_1, W_2, \dots, W_n]^T$$

2.4 Consistency Test

In order to make the judge results coincide with the actual situation, consistency test is needed. The consistency test formula of judgement matrix is $CR = \frac{CI}{RI}$, where CI represents the consistency test index, $CI = \frac{\lambda_{max} - n}{n - 1}$; n represents the order of the matrix; and RI represents the average random consistency index (values shown in Table 1) (Kamal and Al-Harbi, 2001).

Table 1. Values of the average random consistency index RI

n	1	2	3	4	5	6	7	8	9	10	11
RI	0	0	0.58	0.90	1.12	1.24	1.31	1.41	1.45	1.49	1.51

When $CR < 0.001$, the matrix consistency is acceptable, otherwise, the judgment matrix is needed to adjust until it meets the consistency test.

3 The Parameters Identification Model of Casing Damage

3.1 Model Structure

Comprehensive consideration of geological factors, development factors and engineering factors and based on the causes and formation mechanism of casing damage, a parameter of quantitative identification of casing damage is obtained, which is called casing parameters (para-TS). The index values of geological factors F_{Di} are multiplied by their weights W_{Di} and accumulate, and the accumulated values are recorded as F_D , which is called the geological sub-degree of casing damage; the index values of development factors F_{Ki} are multiplied by their weights W_{Ki} and accumulate, and the accumulated values are recorded as F_K , which is called the development sub-degree of casing damage; the index values of engineering factors F_{Gi} are multiplied by their weights W_{Gi} and accumulate, and the accumulated values are recorded as F_G , which is called the engineering sub-degree of casing damage.

Name

$$F_D = \sum_{i=1}^n (F_{Di} \times W_{Di}) \tag{1}$$

$$F_K = \sum_{i=1}^n (F_{Ki} \times W_{Ki}) \tag{2}$$

$$F_G = \sum_{i=1}^n (F_{Gi} \times W_{Gi}) \tag{3}$$

The geological sub-degree, development sub-degree and engineering sub-degree and their weights are respectively multiplied and add together to obtain the casing parameters:

$$para-TS = F_D \times \overline{W_D} + F_K \times \overline{W_K} + F_G \times \overline{W_G} \quad (4)$$

According to the practical condition of the research area, the values of the model are $\overline{W_D} = 0.35$, $\overline{W_K} = 0.45$, $\overline{W_G} = 0.20$.

3.2 The Type And Standards Of Casing Damage Evaluation

Through a comprehensive analysis of the causing of damage factors, combined the data and the comprehensive analysis of the research area base on expert experience, the oilfield casing damage evaluation in different geological and development conditions is divided into 4 types, namely: no abnormality, abnormality led by geological factors, abnormality led by development factors, abnormality led by engineering factors.

Model evaluation criteria are:

First, when the casing damage parameters (para-TS) < 0.4, it is regarded as no abnormality.

Second, when the casing damage parameters (para-TS) > 0.4, it is divided into the following situations:

- (1) If $F_D > 0.7$, it is regarded as the abnormality led by geological factors;
- (2) if $F_K > 0.5$, it is regarded as the abnormality led by development factors;
- (3) If $F_G > 0.8$, it is regarded as the abnormality led by engineering factors;

4 The Analysis of the Practical Application Result on Oilfield

Combined with casing well to do the statistical analysis of the regional geological characteristics of the study area, it can be concluded that the main control factors of the casing damage in the study are: the “fault conducting water” phenomenon of the mudstone dipping waters of the fault plane; the hold pressure phenomenon led by the poor reservoir properties (especially in the permeability zone); the swelling after adsorbing water of the mudstone with fast water rising, the compression of casing damage because of the creep.

4.1 To Determine Weights by the Application of AHP

Geological factors:

The consistency test of Geological factors is RI=0.0002.

Development factors: The consistency test Development factors is RI=0.00001.

Engineering factors: The consistency test of Engineering factors is RI=0.00025.

Table 2. The computation process of geological factors

Comparison matrix						Judgment matrix					Calculate weight
B1	C1	C2	C3	C4	C5	C1	C2	C3	C4	C5	W_i
C1	2	3	4	3	4	1.0000	1.4422	6.2403	3.0000	9.0000	0.4352
C2	1	2	4	3	4	0.6934	1.0000	4.3267	2.0801	6.2403	0.3017
C3	0	0	2	1	3	0.1602	0.2311	1.0000	0.4807	1.4422	0.0698
C4	1	1	3	2	3	0.3333	0.4087	2.0803	1.0000	3.0000	0.1450
C5	0	0	1	1	2	0.1111	0.1602	0.6934	0.3333	1.0000	0.0483

Finally the weight of geological factors is (0.4352,0.3017,0.0698,0.1450,0.0483). The weight of development factors is (0.2806 0.1165 0.5425 0.0602).The weight of engineering factors is(0.4951 0.3038 0.0551 0.1460).

4.2 The Test of Judgment Model of Casing Damage Parameter

There are a total of 77 wells in the research block, here we only used the Bei 54-60 wells to show the inspection process of model instance. The inspection results of model instance of the nine casing damage wells that has been found in the research area are shown in Table 3.

4.2.1 Geological sub-degree of Bei 54-60

a. Fault: The causing point at the depth of 1297.0-1299.0 of Bei 54-60 is near faults of the No. 10 small layer, therefore the index value of the fault is 1 (the principle is that the value is taken 1 at the fault, otherwise 0 should be taken).

b. Lithological characters: The shale content is taken to evaluate the lithological characters of the research area according to the actual geological conditions and Statistical analysis of the previous. The data are normalized by formula five, then calculated the average of the shale data to obtain the index value of shale content, which is 0.2413:

$$Y = X - \min val / \max val - \min val \quad (5)$$

c. Porosity: The index value of porosity after normalization is 0.8517.

d. Permeability: Because 10 small layer of Bei 54-60 is near the fault, the permeability should be the data of the small layers near the fault, which has a little influence on the weight. The index value of permeability after normalization is 0.7024.

e. Sensitivity: Research shows that there is serious rate sensitivity in the reservoir of research block, therefore the index value of sensitivity is taken 1 (the principles to get the value of the sensitivity is that if there is sensitivity in the reservoir, take one as the sensitivity, otherwise take 0).

The geological sub-degree is

$$F_p = 0.4325 * 1 + 0.3017 * 0.2413 + 0.0698 * 0.8517 + 0.1450 * 0.7024 + 0.0483 * 1 = 0.7149.$$

4.2.2 Development sub-degree of Bei 54-60

a. Oil pressure: The index value of oil pressure after normalization is 0.6412.

b. Layered water absorption: The index value of layered water absorption after normalization is 0.4813.

c. Injection allocation: The index value of injection allocation is 0.3207 after normalization of data from April 2004 to September 2009.

d. The time factor: The index value of time factors of Bei 54-60 after normalization of the time data of nine casing damagement wells is 0.8214.

By calculation the development sub-degree

$$F_K = 0.2806 * 0.6412 + 0.1165 * 0.4813 + 0.5425 * 0.3207 + 0.0602 * 0.8214 = 0.4594$$

4.2.3 Engineering sub-degree of Bei 54-60

a. Perforation: There are perforations at set-point depth of 1297.0-1299.0, so the index value of perforation is 1 (the principles to get index value of perforations is that if there are perforations at set-point depth, it is taken as 1, or it is taken as 0.).

b. Well cementation: Combined with actual field data, cementing of Bei 54-60 is good, thus the index value of cementing is 0.4 (the principles to get the index value of cementing quality is that the excellent value is 0.2, good value is 0.4, general value is 0.6, and the poor value is 0.8).

c. Corrosion: Because the groundwater salinity of the whole area is lower by overall analysis, the index value of corrosion factors of research block is taken as 0.200.

d. Stimulation treatment: Bei 54-60 wells has no any fracturing, acidizing and other stimulation, the index value of stimulation is 0.0000.

By calculation that the engineering factor is

$$F_G = 0.4951 * 1 + 0.3038 * 0.4 + 0.551 * 0.200 + 0.1460 * 0 = 0.7268$$

$para - TS = F_D \times \overline{W_D} + F_K \times \overline{W_K} + F_G \times \overline{W_G} = 0.6023$, the result showed that the casing damage well Bei 54-60 is the abnormality led by geological factors, which is consistent with the actual judgment result of the oilfield (Table 3).

Table 3. Model test results of nine casing damage wells of 301 block

Well number	Geological sub-degree	Development sub-degree	Engineering sub-degree	Comprehensive evaluation	Interpretation of model
Bei 54-60	0.7	0.4	0.7	0.6023	Geological abnormality
Bei 58-60	0.8	0.4	0.6	0.5800	Geological abnormality
Bei 46-55	0.7	0.6	0.4	0.5950	Development abnormality
Bei 3-1	0.6	0.8	0.5	0.6700	Development abnormality
Bei 58-54	0.5	0.6	0.6	0.5650	Development abnormality
Bei 32-56	0.5	0.4	0.2	0.3950	No abnormality
Bei 44-54	0.5	0.5	0.8	0.5600	Engineering abnormality
Bei 60-62	0.5	0.7	0.6	0.6100	Development abnormality
Bei 52-54	0.5	0.5	0.9	0.5800	Engineering abnormality

5 Conclusion

There are a total of 77 wells in research area, which 9 casing wells, 68 non-casing wells in the current. Applying the method to write the computer program for batch data automatically determine that there are 14 casing damaged wells, 63 non-casing wells. Compared with the actual judgment result, the compliance rate is much more than 85%. So the main conclusions are:

(1) The results show that it is better to determine the weight of impact factors of casing damage in oilfield applying AHP method and the judgment model of casing damage parameter is reasonable.

(2) Writing computer programs about the AHP method to quantitatively identify and forecast oil casing is a simple, practical, economical and effective way.

(3) Because there are some subjective factors to get the index values of various impact factors, AHP method provides a preliminary quantitative method for the Quantitative Study of oilfield casing damage.

References

1. Altuzarra, A., Moreno-Jiménez, J.M., Salvador, M.: A Bayesian prioritization procedure for AHP-group decision making. *European Journal of Operational Research* 182, 367–382 (2007)
2. Chiotis, E., Vrellis, G.: Analysis of casing failures of deep geothermal wells in Green. *Geothermics* 24(5), 695 (1995)
3. Dagdeviren, M., Yavuz, S., Kilinç, N.: Weapon selection using the AHP and TOPSIS methods under fuzzy environment. *Expert Systems with Applications* 36, 8143–8151 (2009)
4. Dai, L., Xu, S.Y.: Integratal study of geological factors causing casing damage in oil wells. *Geological Disaster and Environmental Protection* 16(3), 331–334 (2005)
5. Dong, Y.C., Xu, Y.F., Li, H.Y., Dai, M.: A comparative study of the numerical scales and the prioritization methods in AHP. *European Journal of Operational Research* 186, 229–242 (2008)
6. Dou, Z.L., Zeng, L.F., Zhang, Z.H., et al.: Research on the diagnosis and description of wormhole. *Petroleum Exploration and Development* 28(1), 75–77 (2001)
7. Hua, J.P., Chen, X.D.: Application of engineering well logging data in casing damage. *Petroleum Geology and Oilfield Development in Daqing* 23(4), 78–80 (2004)
8. Jiang, Q.Y., Xie, J.X., Ye, J.: Mathematical model, pp. 227–228. *The Press of High Education, Beijing* (2003)

Complete SAT Solver Based on Set Theory^{*}

Wensheng Guo¹, Guowu Yang¹, Qianqi Le¹, and William N.N. Hung²

¹ University of Electronic Science and Technology Chengdu, Sichuan, China

² Synopsys Inc., Mountain View, California, USA

{gws,guowu}@uestc.edu.cn, leqqi777@163.com, william_hung@alumni.utexas.net

Abstract. SAT problem is a NP-complete. Many SAT benchmarks that come from the different real life SAT problems are proposed to verify the performance of solvers. Our research focuses on the Model RB benchmark which can be mapped by the coloring problem and others. We propose a translating method based on set for Model RB instances of CNF formulas, and a complete search algorithm. We use the weight of clauses based on the set to determine the order of the search. The results show our solver has the best runtime for the mostly instances and is comparable to the best SAT solvers.

Keywords: SAT, Set, Exclusive Set, Relative Set

1 Introduction

Boolean satisfiability(SAT) is the problem that decides an assignment for the variables in a boolean formula that makes the formula true. SAT is a NP-complete problem and is central in the theory of computation. Algorithms for the satisfiability problems have been studied since the early years of the discipline of computer science. Many real-world problems can be transformed into SAT problems and many of these problems can be effectively solved via satisfiability [1–13]. It is hard to design the general SAT algorithm to solve all SAT problems. There are, however, many theoretical and experimental results which show good average-case performance for certain class of SAT problems [14–17]. In this paper, we focus on the Model RB benchmark and we introduce a kind of SAT translating and searching method based on set. We use the set to describe the SAT clauses and the relative clauses as the condition for search and backtracking. Experimental results indicate that our SAT solver provides the better performance among used complete satisfiability solvers to establish satisfiable assignments of the instances based on Model RB benchmark.

2 Preliminaries

2.1 DPLL

Most SAT solving algorithms are based on DPLL [18, 19], which is a search strategy prioritized based on the depth of backtracking. The each iteration of

^{*} This research supported in part by NSFC Program(No.60973016) of China.

the algorithm selects a variable and assigns it with a value. It essentially reduces the problem because clauses that are already satisfied under the current set of variable assignments can be removed, and literals that are false (based on current set of variable assignments) can be removed from their clauses. The algorithm iterates until all variables have been assigned or it encounters an empty clause (conflict clause). An empty clause is a clause without literals (because all literal were removed as *FALSE* under the current set of variable assignments). The algorithm backtracks upon encountering an empty clause, and assigns a different value to one of the variables on its search tree (mostly related to the conflict). The algorithm repeats the above process until it finds a satisfiable solution or exhausts the search space (unsatisfiable).

2.2 Model RB Benchmark

The SAT benchmarks are usually divided into three categories: randomly, applications and crafted. Model RB is a type of random CSP model, which is a revision to the standard Model B. Coloring problem can be translated into the same form. Its CNF form (conjunctive normal form - conjunction of disjunctions) has the following form [20].

1. First n clauses are disjoint sets of all boolean variables, every set has the same number of variables;
2. Other clauses are the random 2-clause $\neg x \vee \neg y$, x and y are any two mutually exclusive variables.
3. Exactly one variable from every disjoint set can take value *TRUE*, and if no random clause is violated by this assignment, then the instance is satisfiable.

3 SAT Search Algorithm

3.1 Mapping the CNF to Set

We use variable sets MC_1, MC_2, \dots, MC_n to denote the first n clauses. Every clause is disjoint set of limited boolean variables, every set has the same number of variables. All other clauses are of the form $(\neg x \vee \neg y)$, the clause is satisfied when variables x, y are not simultaneously *TRUE*. Hence we can create a exclusion set MN_e for each variable e , such that if e is *TRUE*, the variables inside MN_e are all *FALSE*. We then have the following theorem:

Theorem 1. *The problem instance S is satisfiable if and only if there exists a set of variables $\{e_1, e_2, \dots, e_n | e_i \in MC_i, i = 1, 2, \dots, n\}$ such that $\{e_1, e_2, \dots, e_n\} \cap (MN_{e_1} \cup MN_{e_2} \cup \dots \cup MN_{e_n}) = \emptyset$.*

Proof. Necessary condition: If S is satisfiable, then there is an assignment to all variables such that all clauses in S are satisfied. Under such assignment, each disjunction of variables in MC_i of S is *TRUE*, which means this assignment will assign *TRUE* to variable e_i (where $e_i \in MC_i, i = 1, 2, \dots, n$). We can construct

a set from the above variables: $\{e_1, e_2, \dots, e_n\}$. This set of variables must satisfy $\{e_1, e_2, \dots, e_n\} \cap (MN_{e_1} \cup MN_{e_2} \cup \dots \cup MN_{e_n}) = \emptyset$. Suppose the above set of variables does not satisfy **(II)**, there exists e_k , where $e_k \in \{e_1, e_2, \dots, e_n\} \cap (MN_{e_1} \cup MN_{e_2} \cup \dots \cup MN_{e_n})$. Hence $e_k \in MN_{e_j}$, which means $e_k \neq e_j$, i.e., there is a clause $(\neg e_k \vee \neg e_j)$ in S . Since e_k, e_j are both *TRUE* under the assignment, we know $(\neg e_k \vee \neg e_j)$ is *FALSE*, which contradicts with the statement that $(\neg e_k \vee \neg e_j)$ is a satisfiable clause under S . Hence the condition is necessary.

Sufficient condition: Suppose there exists a set of variables $\{e_1, e_2, \dots, e_n | e_i \in MC_i, i = 1, 2, \dots, n\}$ such that $\{e_1, e_2, \dots, e_n\} \cap (MN_{e_1} \cup MN_{e_2} \cup \dots \cup MN_{e_n}) = \emptyset$. We now assign *TRUE* to variables in the set $\{e_1, e_2, \dots, e_n | e_i \in MC_i, i = 1, 2, \dots, n\}$, and *FALSE* to all other variables. Obvious, the first n clauses that are denoted by MC_i are satisfiable under such assignment. All other clauses in S are of the form $(\neg x \vee \neg y)$. Hence if S is unsatisfiable under this assignment, then there exists a clause $(\neg x \vee \neg y)$ that is unsatisfiable, which means both x, y are *TRUE*. Based on our variable assignment, we have $x, y \in \{e_1, e_2, \dots, e_n\}$. Since $y \in MN_x$, we have $y \in \{e_1, e_2, \dots, e_n\} \cap (MN_{e_1} \cup MN_{e_2} \cup \dots \cup MN_{e_n})$, which contradicts with our condition. Hence S must be satisfiable.

For any i th clause, $i = 1, 2, \dots, n$, we create the relative set

$$MCR_i = \{MC_j | \exists (\neg x \vee \neg y), \text{ while } x \in MC_i, y \in MC_j, i \neq j, 1 \leq j \leq n\} \quad (1)$$

The sets MCR describe the relations of the first n clauses, We then have the following theorem:

Theorem 2. *If $\{e_1, e_2, \dots, e_n | e_i \in MC_i, i = 1, 2, \dots, n\}$ is a satisfiable solution, then $\forall MC_k \in MCR_i, MC_k \setminus MN_{e_i} \neq \emptyset$.*

Proof. Suppose there exists $MC_k \in MCR_i$ and $MC_k \setminus MN_{e_i} = \emptyset$, that is $MC_k \subseteq MN_{e_i}$, all the variables in MC_k is *FALSE*, thus $e_k \in MC_k$ is *FALSE* which contradicts with the statement that $\{e_1, e_2, \dots, e_n | e_i \in MC_i, i = 1, 2, \dots, n\}$ is a satisfiable solution(that is $e_k = \text{TRUE}$). Hence theorem is correct.

For any clause $C_i, i = 1, 2, \dots, n$, we set the weight of clause

$$Weight(v_i) = \max_{e_k \in v_i} Max(e_k) \quad (2)$$

where the function Max is defined as follows:

$$Max(e) = \max_{MC_j \in MCR_i} |MN_e \cap MC_j|, e \in v_i \quad (3)$$

The Formula **(3)** computes the intersection of the exclusive set of a variable $e \in C_i$ and all other clauses, then it chooses the maximum one as the function value. The Formula **(2)** chooses the maximum one among $Max(e_k), e_k \in C_i$ as the weight of clause C_i , which indicates the relative degree for the clause and is used to select the search root.

3.2 Basic Ideas

Before the SAT searching, we translate the CNF to set form. To reduce the search complexity, the algorithm should begin to search from the clause that has the maximum relative degree as the root clause. After determining the root clause, the algorithm assigns the *TRUE* to a variable in the root clause. According to the Theorem 1, we can construct the new $MC_new_i = \{x | x \in MC_i, x \notin MN_e\}$, where e is the variable that has been assigned to *TRUE*, i is the clause sequence number that is not searched. Then we choose the minimum MC_new_i set as the next search clause and assign the *TRUE* to the first variable in the MC_new_i . The solver repeats the selection and assignment until MC_new_i is *NULL* that indicates the current search path is wrong. For reducing the computation complexity of empty determination, we only need compute the MC_new in MCR_i , where i is the current selected clause.

3.3 The Algorithm

Based on the relative degree and sets, we purpose a algorithm that is named by the relative set SAT solving algorithm(R-SAT). The pseudo code for the algorithm is shown in Algorithm 1. The solving process is similar to depth-first-search on a tree, with maximum tree depth $n + 1$. The root node of the tree is an empty set, with child nodes determined by MC_j where j is decided at the initialization of the algorithm. Every element of MC_j is a child node of the root node. The nodes of the tree are gradually constructed during the search process. A satisfiable solution is reached when the search found a path from the root node to a leaf node, with path length n . We describe the main functions of this algorithm in the following subsections.

Selecting an Index. The search tree is implicitly constructed during the solving process. The size of the search space and the size of the tree is directly related to the selection of index, which greatly affects the efficiency of the algorithm. Hence the main purpose of the 5th line *select_clause_index* function and the 29th line *select_search_index* function should be reducing the search space. We adopt two strategies here:

1. During initialization of the algorithm, the *select_clause_index* function computes a weight based on the Formula (2) for each set MC_1, MC_2, \dots, MC_n , and selects the set with the largest weight for the search. The weight for the set MC_i is the sum of attribute value of every element in that set. To compute the attribute value of an element variable e , we take its exclusive set MN_e and intersect with other variable sets $MC_k (k \neq i)$. The largest number of elements in the intersection results is the attribute value of this element variable e .
2. During the search, the *select_search_index* function picks an index j from the index set $MV(j \in MV)$, where j satisfies $Size(MC_j) = \min Size(MC_k) | k \in MV$. If there are more than one satisfiable indices, we randomly pick one among them.

Algorithm 1. The R-SAT Algorithm

```

1: Function R-SAT(CNF for problem instance  $S$ )
2: Initialization: Construct  $MC_i, MN_e, MCR_i$  based on CNF file
3: BOOL  $result \leftarrow false$ 
4:  $MV \leftarrow [1, 2, \dots, |MC|]$ 
5:  $root \leftarrow select\_clause\_index(MV)$ 
6:  $MV \leftarrow MV \setminus [root]$ 
7:  $Result \leftarrow recursive\_search(1, root, MC, MV)$ 
8: if  $Result$  then
9:   Print("SAT")
10: else
11:   Print("UNSAT")
12: end if
13: End Function
14: BOOL Function  $recursive\_search(depth, search\_clause, MC\_in, MV\_in)$ 
15: BOOL  $result \leftarrow false$ 
16:  $Order(MC\_in_{search\_clause})$  {order the variable in the selected clause based on the
   size of  $MN_e$ }
17: for  $element \in MC\_in_{search\_clause}$  do
18:    $MCR\_intersection \leftarrow MCR_{search\_clause} \cap MV\_in$ 
19:   for  $number \in MCR\_intersection$  do
20:      $MC\_new\_number \leftarrow MC\_in_{number} \setminus MN_{element}$ 
21:     if  $MC\_new\_number = NULL$  then
22:       break
23:     end if
24:   end for
25: end for
26: if  $depth = MC\_set\_count$  then
27:   return  $true$ 
28: end if
29:  $new\_mc\_selected \leftarrow select\_search\_index(MC\_new, MV\_in)$ 
30:  $MV\_in \leftarrow MV\_in \setminus [new\_mc\_clause]$ 
31:  $result \leftarrow recursive\_search(depth + 1, new\_mc\_clause, MC\_new, MV\_in)$ 
32: return  $result$ 
33: End Function

```

Selecting an Element. The 16th line *Order* function is used to sort the variable in the selected clause based on the size of MN_e . The ordering of variables in MC has substantial impact on the efficiency of the algorithm. If the SAT problem instance is satisfiable, then picking a variable that has a *TRUE* assignment in the satisfying solution will greatly speed up the search process. In this paper, we sort the variables of the set MC by cardinalities of each variables exclusive set MN , and pick the variables with cardinalities from small to large.

Backtracking. The purpose of the *recursive_search* function is to recursively search and restore the status to prior level at the given parameter depth level. It handles some data processing for selecting a new variable at that level, such

Table 1. Experimental comparison of SAT solvers over Model RB benchmark instances. We used the CPU time limit of 5000 seconds for all the tests.

Sample	Runtime(s)			
	MiniSAT	glucose	clasp	R-SAT
Frb30-15-1.cnf	1.94	0.056	1.18	0.01
Frb30-15-2.cnf	0.468	0.036	0.3	0.35
Frb30-15-3.cnf	1.66	0.32	1.06	0.26
Frb30-15-4.cnf	0.524	0.336	0.51	0.24
Frb30-15-5.cnf	0.536	0.484	0.18	0.01
Frb35-17-1.cnf	5.016	5.6324	1.49	2.26
Frb35-17-2.cnf	49.339	6.24	12.89	4.75
Frb35-17-3.cnf	13.0848	8.2845	1.98	0.48
Frb35-17-4.cnf	0.42	1.844	1.27	0.87
Frb35-17-5.cnf	18.0811	8.6365	5.3	2.93
Frb40-19-1.cnf	0.356	0.264	28.86	0.54
Frb40-19-2.cnf	25.5136	32.662	3.87	5.99
Frb40-19-3.cnf	652.753	234.143	317.71	12.45
Frb40-19-4.cnf	601.094	118.615	51.31	61.87
Frb40-19-5.cnf	158.494	123.136	90.23	99.62
Frb45-21-1.cnf	>5000	>5000	2579.6	251.2
Frb45-21-2.cnf	>5000	983.553	1536.76	536.07
Frb45-21-3.cnf	>5000	71.3045	2529.36	2248.28
Frb45-21-4.cnf	287.594	167.082	812.67	574.94
Frb45-21-5.cnf	>5000	1916.22	636.23	424.95

as restoring data (MC, MV) back to the $depth - 1$ level. This level is the closest level to the search node with empty set, and has some new variables that have not been selected yet. To reduce the computation complexity, the 17th-25th lines use the relative sets as basis for the empty determination based on the Theorem 2. Our algorithm keeps on going through the index selection, variable selection, and backtracking (if necessary) to complete the entire search. If the algorithm can find a set of variables without conflict, then the SAT problem instance is satisfiable, otherwise it is unsatisfiable.

4 Experimental Results

Our experiments use the instances from Model RB benchmark [21]. We chose three state-of-the-art SAT solvers for comparison with R-SAT. The SAT solvers were MiniSAT, glucose and clasp (we used the latest available versions to the time of writing this paper). All the tests were run on a machine with Intel i5 M430 CPU 2.27GHz with 2GB of memory under Ubuntu10.10. Although the testing machine has two processors, no parallel processing was used. We used the CPU time limit of 5000 seconds for all the tests. For each SAT instance, we measured the overall time necessary to decide its satisfiability. The results are shown in Table 1, where > 5000 means the runtime exceeds 5000 seconds (our

timeout threshold). The speedup obtained by using our method compared to the selected SAT solvers. It is obvious that our solver is better than the other SAT solvers. Our solver has the best runtime for almost all the instances. This indicates our approach is a very promising for this kind of problem.

5 Conclusion

In this paper, we purpose a novel translating method of the set. Based on the weight of clause, we describe the selection of the search root clause, clause search order and variable search order. Experimental results indicate that our solver is faster and uses much less memory than the other SAT solvers.

References

1. Larrabee, T.: Test pattern generation using Boolean Satisfiability. *IEEE Trans. CAD* 11(1), 4–15 (1992)
2. Wood, R.G., Rutenbar, R.A.: FPGA routing and routability estimation via Boolean satisfiability. In: *International Symposium on Field-programmable Gate Arrays*, Monterey, California, United States, pp. 119–125. ACM, New York (1997)
3. Biere, A., et al.: Symbolic Model Checking using SAT procedures instead of BDDs. In: *Proc. Design Automation Conference ACM/IEEE*, pp. 317–320 (1999)
4. Bjesse, P., Leonard, T., Mokkedem, A.: Finding Bugs in an Alpha Microprocessor Using Satisfiability Solvers. In: *Berry, G., Comon, H., Finkel, A. (eds.) CAV 2001. LNCS*, vol. 2102, pp. 454–464. Springer, Heidelberg (2001)
5. Song, X., et al.: Board-Level Multiterminal Net Assignment for the Partial Cross-Bar Architecture. *IEEE Trans. VLSI Systems* 11(3), 511–514 (2003)
6. Hung, W.N.N., Narasimhan, N.: Reference Model Based RTL Verification: An Integrated Approach. In: *Proc. IEEE International High Level Design Validation and Test Workshop (HLDVT)*, pp. 9–13. IEEE (November 2004)
7. Hung, W.N.N., et al.: Segmented channel routability via satisfiability. *ACM Transactions on Design Automation of Electronic Systems* 9(4), 517–528 (2004a)
8. Hung, W.N.N., et al.: Routability Checking for Three-Dimensional Architectures. *IEEE Trans. VLSI Systems* 12(12), 1398–1401 (2004b)
9. Hung, W.N.N., et al.: Optimal Synthesis of Multiple Output Boolean Functions using a Set of Quantum Gates by Symbolic Reachability Analysis. *IEEE Trans. CAD* 25(9), 1652–1663 (2006)
10. He, F., et al.: A satisfiability formulation for FPGA routing with pin rearrangements. *International Journal of Electronics* 94(9), 857–868 (2007)
11. Hung, W.N.N., et al.: Defect Tolerant CMOL Cell Assignment via Satisfiability. *IEEE Sensors Journal* 8(6), 823–830 (2008)
12. Han, H., Somenzi, F., Jin, H.: Making Deduction More Effective in SAT Solvers. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 1271–1284 (2010)
13. Audemard, G., Katsirelos, G., Simon, L.: A restriction of extended resolution for clause learning SAT solvers. In: *Proceedings of the 24th AAAI Conference on Artificial Intelligence, AAAI* (2010)
14. Franco, J., Paull, M.: Probabilistic analysis of the Davis Putnam procedure for solving the satisfiability problem. *Discrete Applied Mathematics* 5, 77–87 (1983)

15. Mendonca, M., Wasowski, A., Czarnecki, K.: SAT-based Analysis of Feature Models is Easy. In: Proceedings of the International Software Product Line Conference (SPLC). Software Engineering Institute, Carnegie Mellon University (2009)
16. Kemper, S.: SAT-based verification for timed component connectors. *Electronic Notes in Theoretical Computer Science (ENTCS)* 255, 103–118 (2009)
17. Zhou, M., He, F., Gu, M.: An Efficient Resolution Based Algorithm for SAT. In: 2011 Fifth International Symposium on Theoretical Aspects of Software Engineering (TASE), pp. 60–67 (2011)
18. Davis, M., Putnam, H.: A computing procedure for quantification theory. *J. ACM* 7, 201–215 (1960)
19. Davis, M., Logemann, G., Loveland, D.: A machine program for theorem proving. *Comms. ACM* 5, 394–397 (1962)
20. Xu, K., Li, W.: Exact Phase Transitions in Random Constraint Satisfaction Problems. *Journal of Artificial Intelligence Research* 12, 93–103 (2000)
21. <http://www.nlsde.buaa.edu.cn/~kexu/benchmarks/benchmarks.html>

Application of XML Data Mining in GUI Run-Time State Clustering

Jing Feng^{1,2} and Tingjie ShangGuan²

¹ Beijing University of Aeronautics and Astronautics 100083 Beijing, China

² Institute of China Electronic System Engineering Corporation 100141 Beijing, China
fj_mail@126.com, sgtj73@sina.com

Abstract. Run-time state clustering can be available to solve many software problems. As the representations of GUI states are structural, it is a challenging task for GUI run-time states clustering. In this paper, XML documents are adopted to represent and retrieve GUI states and a structural XML-based approach is proposed for clustering GUI states. This approach takes into account differences of features of XML element/attribute labels in forming clusters during XML document comparison. These labels stand for windows, widgets, properties and their value in GUI state descriptions. In this approach, a new step was added in the K-means clustering process to calculate the weights of features in each cluster so that the important features of a state cluster could be identified by the weight values. Experimental results on real GUI state datasets confirmed this new algorithm performed well in practice, and experimental evaluation showed that this proposed algorithm was almost accurate.

Keywords: GUI, Run-time state, Clustering, XML documents.

1 Introduction

Run-time state clustering is available to solve software states explosion problems, which discourages the usage of state machine. This issue is especially evident for GUI (Graphical User Interface) applications because GUIs are state rich. If a button or a menu in the GUIs is pressed, the state of the application will change into another one. Of course, when we start to store in the state all various value changes we will meet the state explosion problem, but if we ignore the trivial differences of states and capture the main features of all kinds of states, the problem of states explosion will be solved. In state clustering, that means to choose a suitable clustering granularity. Only a few existing works have studied on this issue for non-GUI application. For instance, a clustering approach was employed to partition a weighted hypergraph into k clusters and introduced these clusters into software modeling to do state definition in a different way [1].

GUI run-time states clustering is a challenging task. As a common example of event-driven software (EDS), the state of a GUI at a particular time t is the set P of all the properties of all the objects O that the GUI contains [2]. The structural

representation of state makes this task much more difficult. XML has the ability that different sources of structured data easily be combined so that data mining for GUI states has been brought new chance.

In this paper, we proposed a structural XML-based approach for clustering GUI states, taking into account differences of features of windows, widgets, properties and their value in GUI state descriptions. The remaining of this paper is organized as follows. The definition of similarity of GUI states is introduced in section 2. In section 3, we detail our proposal. Section 4 presents Case study results, and during experiments we choose a suitable clustering granularity to validate our approach. Section 5 concludes our work and snapshots our future work.

2 Definition of GUI State Similarity

2.1 GUI State

A GUI state can be modeled as a set of widgets $W=\{w_1, w_2, \dots, w_n\}$ that the GUI contains, a set of properties $P=\{p_1, p_2, \dots, p_m\}$ of these widgets, and a set of values $V=\{v_1, v_2, \dots, v_n\}$ of the properties, which is described as a set of 3-tuple $S=\{(w_i, p_j, v_k)\}$, where $w_i \in W, p_j \in P, v_k \in V$. As this description presents a multidimensional view of data, in order to read document easily and parse syntax, we adopted XML document to describe a complete state of the GUI. For example, consider an Open GUI shown in Fig.1 (a), the state of the Open GUI, partially shown in Fig.1 (b), is represented as a XML document. Fig.1(c) is a tree view of descriptions for Open GUI states.

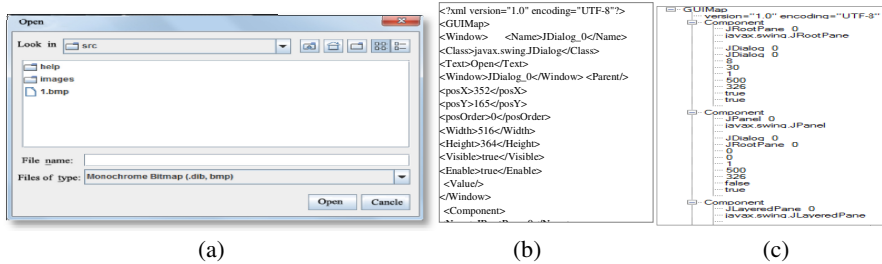


Fig. 1. (a) A simple GUI application (b) The GUI's partial state file (c) Tree view of state file

These GUI states can be captured by several tools. For example, runtime assertion checkers and runtime monitoring tool can record window sessions of GUI application.

2.2 GUI Application State Similarity

The GUI application state at a particular time t can be modeled as a set of GUIs states (to avoid confusion here named windows). If we use a XML document to express the state information of one window, then the GUI application state can be recorded by

several XML documents. Here a complex node in the schema represents a widget. Each leaf element in the XML document represents one property and corresponding value of this property. Therefore the similarity between two GUI application states must rely on two elements: similarity of windows and similarity of widgets in these windows. Formally, we give below several definitions:

Definition: (Similarity of simple widgets)

Simple widget similarity (Sim_{wid}) is the similarity of leaf properties' values. According to different data types of value, different methods can be used to calculate similarity of simple widget.

Definition: (Windows similarity)

Window similarity (Sim_{win}) is the aggregation of all widgets' similarities of related properties.

Definition: (States similarity)

States similarity ($\text{Sim}_{\text{State}}$) is the aggregation of all windows' similarities of related widgets.

3 The GUI Application State Clustering

3.1 GUI State Capture

The GUI state Capture is a process that, given an executing GUI, returns the current values of all the properties in the complete set for the GUI. There are several different approaches that can be used to automate the process of extracting actual GUI state information in a form. Two possible approaches are as follows:

Screen scraping

Screen scraping is a technique used to selectively remove information from an application's screen/terminal interface for reuse. Typically, the information is accessed by using low-level, terminal-specific system calls. The bitmaps/text obtained are analyzed to determine the correctness of the executing GUI. Although useful for determining exactly what is visible to the user, non-visible properties cannot be verified using screen scraping.

Querying

Querying the GUI's software is a technique to determine the values of all the properties present in the GUI, including non-visible and visible properties. Although the results of the querying technique are more complete than screen scraping, querying requires access to the GUI's code, possibly modifying the code to access the values of properties.

3.2 GUI States Similarity Calculation

The problem of comparing two states' similarity can be converted into the problem of comparing two XML documents' similarity. Many existing works have previously studied XML clustering [3]. For example, the similarity between $\text{Vector}_1(I_{11}, I_{12}, \dots, I_{1n})$ and $\text{Vector}_2(I_{21}, I_{22}, \dots, I_{2n})$ in XML document can be the average of similarities

of all coupled values [4]. Following this way, suppose $widget_i$ and $widget'_i$ are the same widget element in different state document, the similarity of them is computed as follows:

$$Sim_{wid}(widget_i, widget'_i) = \frac{1}{n} \sum_{k=1}^n w_k \times Sim_{pk}(I_k, I'_k) \tag{1}$$

where n is the number of properties of $widget_i$, p_k is the k th property of $widget_i$, w_k is a given weight of p_k and I_k is the value of p_k .

In this work, we use 16 properties to describe a Widget include “Name”, “Class”, “Parent”, “Position”, “Color”, ”Font”, etc. 14 of these properties are proposed as primary information for state comparison during test oracle procedure in [2]. From the tree view of XML document, the widget as an element in XML document has the same structure, and because the properties Name and Class are used to identify a widget, only 14 properties are taken into account as comparing the same widget in different state documents.

For the sake of convenience in calculating, let

$$Sim_{pk}(I_k, I'_k) = \begin{cases} 0 & I_k \neq I'_k \\ 1 & I_k = I'_k \end{cases} \tag{2}$$

The similarity of the same window win_j and win'_j in different state document can be calculated as follows in the same way:

$$Sim_{win}(win_j, win'_j) = \frac{1}{m} \sum_{i=1}^m Sim_{wid}(widget_i, widget'_i) \tag{3}$$

where m is the number of widgets contained in window win_j .

We can provide the similarity definition of two states can be calculated as follows:

$$Sim_{state}(S, S') = \frac{1}{t} \sum_{j=1}^t Sim_{win}(win_j, win'_j) \tag{4}$$

Where t is the maximum of windows recorded in both different state documents. If win_j is not exist in either of the states, let $Sim_{win}(win_j, win'_j)=0$. This similarity cannot show the importance of different windows, widgets and properties (features). For this, we will use another distance definition in the process of clustering.

3.3 Clustering Algorithm

We chose K-means algorithm to realize clustering. To enable the K-means algorithm to identify differences of features in GUI state clusters, the basic algorithm is extended to allow calculation of a weight for each feature, and the importance of a window, widget, property or value can be identified by the weight value. Mathematically, this new weighting K-means algorithm minimizes the following objective function:

$$D(\Omega, C, P) = \sum_{l=1}^k \sum_{j=1}^n \sum_{i=1}^m \omega_{l,j} \times p_{l,i} \times [1 - Sim_{state}(C_{l,i}, S_{j,i})] \tag{5}$$

subject to

$$\begin{cases} \sum_{l=1}^k \omega_{l,j} = 1, & 1 \leq j \leq n \\ \omega_{l,j} \in \{0,1\}, & 1 \leq j \leq n, \quad 1 \leq l \leq k \\ \sum_{i=1}^m p_{l,i} & 0 \leq p_{l,i} \leq 1, \quad 1 \leq l \leq k \end{cases}$$

where $k(\leq n)$ is the number of clusters; $\Omega = [\omega_{l,j}]$ is a $k \times n$ integer matrix, which means: if the l th cluster cannot contain all of the features, then $\omega_{l,j} = 0$, otherwise $\omega_{l,j} = 1$; $C = [C_1, C_2, \dots, C_k] \in S_{k \times m}$ are the k cluster centers of state space; $P = (p_1, p_2, \dots, p_k)$ is the set of weight vectors for all clusters in which each $P_l = (p_{l,1}, p_{l,2}, \dots, p_{l,m})$ is a vector of weights for m features of the l th cluster; $Sim_{state}(C_{l,i}, S_{j,i}) (\geq 0)$ is a similarity measure between the j th state document and the center of the l th cluster on the i th feature. The base algorithm presented as follow is a weighted K-means (WK-means) algorithm which takes the number of desired clusters as input:

```

WK-means Algorithm
1:  $S_{k \times m} \leftarrow$  s of all states  $\{S_0, S_1, \dots, S_n\}$  used for clustering;
2:  $k \leftarrow$  desired cluster number;
3:  $\Omega^{(0)} \leftarrow$  weights of initial state;
4:  $C^{(0)} \leftarrow$  initial cluster center of 1st clusters, and  $C_1^{(0)} \in S_{k \times m}$ ;
5:  $P^{(0)} \leftarrow$  with all entries equal to  $1/m$ ;
6:  $K_{set} \leftarrow \emptyset$ ;
7: While  $S_{k \times m} \neq \emptyset$ 
8: { for  $l=1$  to  $k$ 
9:   for  $j=1$  to  $n$ 
10:    for  $i=1$  to  $m$ 
11:     { compare  $D(\Omega, C, P)$  with  $D(\Omega^{(0)}, C^{(0)}, P^{(0)})$ ;
12:      find minimal  $D_{min}(\Omega, C, P)$  of cluster  $l$ ;
13:       $S_x \leftarrow$  state corresponding to  $D_{min}(\Omega, C, P)$ ; }
14:    $K_{set}^{(0)} = K_{set}^{(0)} \cup S_x$ ;
15:   remove state  $S_x$  from  $S_{k \times m}$ ; }
16  $K_{set}$  is the final clustering set;
    
```

Fig. 2. WK-means algorithm

3.4 Clustering Evaluation

The class labels of the run-time state used in the experiments can actually obtained during GUI testing; therefore we adopted two existing measures Accuracy and FScore to objectively assess the clustering performance.

The measures Accuracy are defined as follows. Given a set of XML documents(to describe different run-time states) R in k categories $C_1, C_2, \dots, C_i, \dots, C_k$, we use a clustering algorithm to cluster it into k clusters $S_1, S_2, \dots, S_j, \dots, S_k$. Let n_i, n_j be the numbers of documents in category C_i and cluster S_j respectively, $n_{i,j}$ be the number of documents appearing in both category C_i and cluster S_j , and n be the total number of documents in R . Then the accuracy of a clustering can be calculated by

$$Accuracy = \sum_{i=1}^k n_{i,i} / n \quad (6)$$

FScore measure [5] is a trade-off between two popular information retrieval metrics, precision and recall. It can be defined as follow:

$$FScore = \sum_{i=1}^k \frac{n_i}{n} \times \max_{1 \leq j \leq k} \left\{ \left(2 \times \frac{n_{i,j}}{n_i} \times \frac{n_{i,j}}{n_j} \right) / \left(\frac{n_{i,j}}{n_i} + \frac{n_{i,j}}{n_j} \right) \right\} \quad (7)$$

4 Case Study

4.1 Subject Applications

In this section, a GUI application was used in the experiments, named *TerpSpreadSheet3.0*, which is an Excel-like GUI application included in the *TerpOffice 3.0* [6]. Its functions include data management, data sorting, data store, data printing, data process, and graph generation.

We executed 5096 test cases to capture the run-time states of this GUI application. These test cases were generated by the traversal in the event-flow graphs of the application[7][8][9]. The entire state space of this GUI application was described by values of certain properties of widgets in all windows after the last event of a test case.

4.2 State Datasets

First of all, 5096 states documents were randomly divided into 20 groups, and there were 20 experimental datasets for clustering. During running these test cases, we found the total number of window types invoked are 11(include main window). 139 faulty versions of *TerpSpreadSheet* were downloaded from [10], all faults of which could be detected by 5096 test cases.

In order to test the clustering quality, first, all the datasets must be manual analyzed to label class of run-time states, which decided by the type of windows invoked and the faults detected by corresponding test cases. The results attained by manual analyzing are shown in tables1 and 2.

Table 1. Results of window types analysis

#of Data Set	States	Types of window	#of Data Set	States	Types of window
1	254	8	11	254	8
2	254	5	12	254	5
3	254	6	13	254	6
4	254	9	14	254	9
5	254	4	15	254	4
6	254	6	16	254	6
7	254	7	17	254	5
8	254	9	18	254	8
9	254	5	19	254	6
10	254	7	20	270	8

Table 2. Results of detected faults analysis

#of Data Set	States	Detected faults	#of Data Set	States	Detected faults
1	254	26	11	254	46
2	254	18	12	254	24
3	254	17	13	254	32
4	254	19	14	254	16
5	254	29	15	254	25
6	254	34	16	254	13
7	254	28	17	254	15
8	254	42	18	254	23
9	254	23	19	254	17
10	254	25	20	270	35

4.3 Clustering Quality

WK-Means clustering algorithms were applied to the 20 datasets, where k is decided according to number of window types in Table 1 and detected faults in Table 2 respectively. The clustering results evaluated with the two evaluation measures in Subsection 3.4 are given in Table 3 and Table 4.

Table 3. Clustering results of K decided by window types

#of Data Set	Accuracy	Fscore	#of Data Set	Accuracy	Fscore
1	0.886	0.861	11	0.867	0.843
2	0.890	0.853	12	0.903	0.907
3	0.908	0.850	13	0.922	0.950
4	0.892	0.855	14	0.899	0.855
5	0.906	0.948	15	0.849	0.848
6	0.846	0.849	16	0.892	0.854
7	0.915	0.947	17	0.906	0.851
8	0.829	0.858	18	0.897	0.851
9	0.844	0.846	19	0.856	0.891
10	0.895	0.857	20	0.878	0.841

Table 4. Clustering results of K decided by detected faults

#of Data Set	Accuracy	Fscore	#of Data set	Accuracy	Fscore
1	0.634	0.546	11	0.720	0.646
2	0.631	0.651	12	0.649	0.551
3	0.515	0.550	13	0.634	0.549
4	0.724	0.752	14	0.731	0.464
5	0.621	0.656	15	0.665	0.551
6	0.588	0.554	16	0.692	0.543
7	0.592	0.658	17	0.604	0.656
8	0.625	0.747	18	0.713	0.544
9	0.834	0.846	19	0.513	0.750
10	0.631	0.651	20	0.423	0.646

Table 3 and Table 4 illustrate that the performance of WK-Means. As we can see, when K is decided by window types, the results are always better than which decided by detected faults. This result is not hard to explain. If more than one window invoked in state S_i is different from the other states, S_i is completely different from the other states. But for the faults detected in testing, the same state may include more than one fault modes. From the above experiment results, we have two observations : (1) We can find which one is a proper principle for clustering by comparing the two clustering solution, and the result is consistent with practice. (2) Using proper principle, our proposed WKMeans algorithms is almost accurate with FScore ranging between 84% and 93%, i. e. it is almost accurate.

5 Summaries and Future Work

In this paper, we proposed a structural XML-based approach for clustering GUI states, taking into account windows, widgets, properties and their value in GUI state descriptions. To our knowledge, this was the first attempt to apply weighted K-means XML (structured data) document clustering approach on GUI states clustering. Experimental results on real GUI state datasets confirmed this new algorithm performed well in practice, and experimental evaluation showed that our proposed approach is almost accurate with FScore ranging between 84% and 93%.

This paper is only a preliminary study on GUI run-time state clustering. A lot of work is remained to be done in future. For example: (1) Clustering results are best shown through visualization. We will develop a visualization tool to visualize the GUI state clusters and the relationships between clusters; (2) the current weighted K-means clustering algorithm should be refined; (3) more various subject application should be investigated to discuss the applicability and effectiveness of our approach.

References

1. Yin, B.-B., Bai, C.-G., Cai, K.-Y.: A Data Mining Approach for Software State Definition. In: Proceedings of the 31st Annual International Computer Science and Application Conference, pp. 179–186 (2007)
2. Memon, A.M.: A Comprehensive Framework for Testing Graphical User Interfaces, PhD thesis, Dept. of Computer Science, Univ. of Pittsburgh (2001)
3. Bartolini, I., Ciaccia, P., Patella, M.: A framework for the comparison for Complex Patterns. In: Proc. 1st International Workshop on Pattern Representation and Management (PaRMa 2004), Crete, Greece (2004)
4. Zhao, Y., Karypis, G.: Evaluation of hierarchical clustering algorithms for document datasets. Technical Report 02-022, University of Minnesota (2002)
5. Schönauer, S.: Efficient Similarity Search in Structured Data, PhD thesis, Ludwig Maximilians Universität München (2003)
6. Yuan, X., Memon, A.M.: Iterative execution-feedback model directed GUI testing. *Information and Software Technology* 52(5), 559–575 (2010)
7. (2010), <http://sourceforge.net/projects/guitar/>
8. Memon, A.M.: (2011), <http://www.cs.umd.edu/~atif/TerpOfficeV3.0/index.html>
9. Yuan, X., Memon, A.M.: Generating Event Sequence-Based Test Cases Using GUI Run-Time State Feedback. *IEEE Trans. Softw. Eng.* 36(1), 81–95 (2011)
10. Memon, A.M.: (2011), <http://www.cs.umd.edu/~atif/Benchmarks/MD2006b.html>

Strong Reduction for Typed Lambda Calculus with First-Class Environments

Shin-ya Nishizaki¹ and Mizuki Fujii^{2,*}

Department of Computer Science,
Tokyo Institute of Technology,
2-12-1-W8-69, O-okayama, Meguro-ku, Tokyo, 152-8552, Japan
nishizaki@cs.titech.ac.jp, fujii.mizuki@lambda.cs.titech.ac.jp

Abstract. Programs contain variables, and the bindings of these variables to the corresponding values are kept in a so-called 'environment'. A first-class environment is a mechanism that the environments in programs can be treated as first-class entities, which are objects that can be passed and returned between functions and procedures. Nishizaki proposed the lambda calculus with first-class environments, called the *environment lambda calculus*, and has investigated its theoretical properties [6–8, 10]. The various systems of the environment lambda calculus are based on weak reduction, that is, application of a substitution to a lambda abstraction is postponed until an argument is applied to it. In this paper, we propose a simply-typed lambda calculus with strong reduction. We investigate several theoretical properties such as the subject reduction theorem.

Keywords: lambda calculus, type theory, first-class environment, reflective programming.

1 Introduction

In a program, variables are bound to values and referred in expressions. The correspondence between variables and their values at some point of time is called an environment, which is usually formalized by a partial function whose domain is a finite set of variables and whose co-domain is a set of denotable values, in the semantics of programming languages.

In the programming language Scheme [11], we can use two kinds of runtime objects – continuations and environments – as first-class citizens; that is, it is possible to pass such values as parameters and to return them as results. The availability of first-class continuations and environments increases the expressiveness of the programming language. In some version of Scheme, the following primitives enable environments to be treated as first-class citizens:

* The second author, Mizuki Fujii, completed this research when she was a student at Tokyo Institute of Technology. She is now with Nomura Research Institute, Ltd.

- **the-environment** is a zero-ary procedure returning a representation of the current environment in which the expression itself is evaluated;
- **eval** is a binary procedure mapping the representation of an expression and the representation of an environment into the value of this expression in this environment.

Using these primitives, a Scheme programmer can explicitly treat environments as first-class citizens in a functional programming style.

An environment does not appear explicitly in a functional program's computation expressed as reduction sequences. An environment is usually represented as a list of pairs of variables and their bound denotation, which forms an implicit computational structure of the λ -calculus.

The substitution is used as a meta-level mechanism to describe the -reduction of the λ -calculus, but it is not an object-level mechanism of the λ -calculus, since it is not an explicit operation in the λ -calculus. The idea of explicit substitutions [15] is an interesting approach to make substitutions work at object-level in the λ -calculus, and explicit substitutions are formalized as object-level substitutions using an environment in the $\lambda\sigma$ -calculus.

Although explicit substitutions allow us to treat an environment at object-level in the $\lambda\sigma$ -calculus, there is still a crucial difference between the object-level environments of the $\lambda\sigma$ -calculus and the first-class environments of Scheme. In the $\lambda\sigma$ -calculus, it is not possible to pass substitutions as parameters. For instance, the following term is not permissible in the $\lambda\sigma$ -calculus.

$$\lambda sub.x[sub],$$

where an explicit substitution is passed to the argument sub . The point to be stressed is that, in the $\lambda\sigma$ -calculus the syntactic class of explicit substitutions is precisely distinguished from its terms. If we introduce first-class environments into the $\lambda\sigma$ -calculus, we should allow $\lambda sub.(x[sub])$, as a permissible term in such an extended λ -calculus. Roughly speaking, the λ -calculus with first-class environments is an extended λ -calculus that allows environments as additional permissible terms.

In [7, 9], we proposed several kinds of typed lambda calculi with first-class environments. In all the systems, a term of substitution application to a lambda abstraction, such as

$$(\lambda x.M) \circ N,$$

is not be defined to be reduced, though a corresponding term of the λ -calculus

$$(\lambda x.M)[\sigma]$$

can be reduced to

$$\lambda x'.(M[x:=x'][\sigma])$$

with a fresh variable x' .

The reason why it is difficult to introduce this kind of reduction is related to the name sensitivity of the environment lambda calculus. For example, the term

$$\lambda x.id$$

means a function taking an argument through a formal parameter x and returns the current environment binding x . If the variable name x is renamed as x' , such as

$$\lambda x'.id,$$

then the returned environment is different to the result of $\lambda x.id$ since the binding names x and x' are different and consequently the returned environments are also different.

In this paper, we impose a restriction that each environment has a finite domain; in other words, each environment only binds a finite number of variables. Thanks to the restriction, we preserve some of the fundamental properties in the previous environment lambda calculus, such as subject reduction theorem and local confluence.

2 Simply-Typed Lambda Calculus with First-Class Finite-Domain Environments

We assumed that countable sets of term variables and of type variables are given in advance of the definition of types and terms. The simply-typed environment lambda calculus with first-class finite-domain environments, $\lambda FDEnv$, is defined by the following syntax and reduction.

Definition 1 (Terms). Terms of $\lambda FDEnv$ are defined inductively by the following grammar:

$$M ::= () \mid x \mid (MN) \mid \lambda x.M \mid (M/x) \cdot N \mid (M \circ N).$$

These are called the empty environment, a (term) variable, a function application, a lambda abstraction, an environment extension, and an environment composition, respectively.

We abbreviate

$$(x_1/x_1) \cdot (x_2/x_2) \cdots (x_n/x_n)(\)$$

to $id_{\{x_1, \dots, x_n\}}$. We write the variables occurring in M as $V(M)$.

Definition 2 (Type). Types of $\lambda FDEnv$ are defined inductively by the following grammar:

$$\begin{aligned} E &::= \{x_1 : A_1\} \cdots \{x_m : A_m\}, \\ A &::= \alpha \mid E \mid (A \rightarrow B), \end{aligned}$$

where $m \geq 0$, x_1, \dots, x_m are term variables distinct from each other, E is called an environment type, A and B are called types, and $(A \rightarrow B)$ is called a function type.

Definition 3 (Strong Reduction). A binary relation $M \rightarrow N$ between terms M and N , called the *strong reduction* of $\lambda FDEnv$, is defined inductively by the following rules:

NullL $(\) \circ M \rightarrow (\),$

Assoc $(L \circ M) \circ N \rightarrow L \circ (M \circ N),$

Dextn $((L/x) \cdot M) \circ N \rightarrow ((L \circ N)/x) \cdot (M \circ N),$

DApp $(M_1 M_2) \circ N \rightarrow (M_1 \circ N)(M_2 \circ N),$

VarRef $x \circ ((M/x) \cdot N) \rightarrow M,$

VarSkip $y \circ ((M/x) \cdot N) \rightarrow y \circ M$ where $x \neq y,$

DLam $(\lambda x.M) \circ N \rightarrow \lambda x'.(M \circ ((x'/x) \cdot (N \circ id_{V(N)}))),$ where $x' \notin V(M) \cup V(N),$

Beta $(\lambda x.M)N \rightarrow M \circ ((N/x) \cdot id_{V(M)}),$

AppL if $M \rightarrow M',$ then $(MN) \rightarrow (M'N),$

AppR if $N \rightarrow N',$ then $(MN) \rightarrow (MN'),$

Lam if $M \rightarrow M',$ then $\lambda x.M \rightarrow \lambda x.M',$

CompL if $M \rightarrow M',$ then $(M \circ N) \rightarrow (M' \circ N),$

CompR if $N \rightarrow N',$ then $(M \circ N) \rightarrow (M \circ N'),$

ExtnL if $M \rightarrow M',$ then $(M/x) \cdot N \rightarrow (M'/x) \cdot N,$ and

ExtnR if $N \rightarrow N',$ then $(M/x) \cdot N \rightarrow (M/x) \cdot N'.$

Example 1. the term $L \circ ((\lambda x.M)N)$ is not a term of the $\lambda\sigma$ -calculus, but a term of the λ_{env} -calculus: environments and terms are separated in the $\lambda\sigma$ -calculus, and hence, $(\lambda x.M)N$ is not allowed to be used as an environment in the $\lambda\sigma$ -calculus.

We next show two examples of reductions. The first is the one which returns an environment value. The reduction sequence is given as follows.

$$\begin{aligned}
 & ((\lambda y.\lambda x.id)M)N \\
 & \rightarrow ((\lambda x.id) \circ ((M/y) \cdot id))N && \text{Beta2} \\
 & \rightarrow id \circ ((N/x) \cdot (M/y) \cdot id) && \text{Beta1} \\
 & \rightarrow (N/x) \cdot (M/y) \cdot id && \text{IdL}
 \end{aligned}$$

The term $((\lambda y.\lambda x.id)M)N$ in the first line corresponds to Scheme program:

`((lambda (y) (lambda (x) (the-environment))) M) N).`

From this program, we know that in Scheme it is possible to make an environment at the meta-level into a data in object-level, by using the primitive `the-environment`.

In the following example, an environment value is passed as an argument of function application:

$$\begin{array}{ll}
(\lambda e.(y \circ e))((\lambda y.\lambda x.id)MN) & \\
\overset{*}{\rightarrow} (\lambda e.(y \circ e))((N/x) \cdot (M/y) \cdot id) & \text{similar to the first example} \\
\rightarrow (y \circ e) \circ (((N/x) \cdot (M/y) \cdot id)/e) \cdot id & \text{Beta2} \\
\rightarrow y \circ (e \circ (((N/x) \cdot (M/y) \cdot id)/e) \cdot id) & \text{Assoc} \\
\rightarrow y \circ ((N/x) \cdot (M/y) \cdot id) & \text{Var} \\
\rightarrow y \circ ((M/y) \cdot id) & \text{VarSkip} \\
\rightarrow M & \text{VarRef}
\end{array}$$

The term $(\lambda e.(y \circ e))((\lambda y.\lambda x.id)MN)$ corresponds to the following Scheme program:

```
((lambda (e) (eval 'y e))
  (((lambda (y) (lambda (x) (the-environment)))) M) N))
```

From this program, we find that the primitive `eval` makes a data in object-level into an environment at the meta-level, by using the primitive `eval`. Actually, the result of evaluating this program is equivalent to the value of the subexpression `M` in the implementations of Scheme, which corresponds to the result obtained by our reduction. To be exact, a gap between the two languages still exists; Scheme is a call-by-value language with computational effects, and our calculus is purely functional and does not assume any specific evaluation strategy.

Definition 4 (Typing). The *typing* of $\lambda FDEnv$ is given as the type judgement $E \vdash M : A$, which is read as “A term M is of type A under an environment type E .”

$$\begin{array}{c}
\frac{}{E \vdash () : \{\}} \text{Null} \quad \frac{}{\{x : A\}E \vdash x : A} \text{Var} \\
\\
\frac{E \vdash M : A \rightarrow B \quad E \vdash N : A}{E \vdash (MN) : B} \text{App} \quad \frac{E \vdash N : E' \quad E' \vdash M : A}{E \vdash (M \circ N) : A} \text{Comp} \\
\\
\frac{\{x : A\}E \vdash M : B}{\{x : C\}E \vdash \lambda x.M : A \rightarrow B} \text{Lam} \\
\\
\frac{E \vdash M : A \quad E \vdash N : E' \quad x \notin V(E')}{E \vdash (M/x) \cdot N : \{x : A\}E'} \text{Extn1} \\
\\
\frac{E \vdash M : A \quad E \vdash N : \{x : B\}E' \quad x \notin V(E')}{E \vdash (M/x) \cdot N : \{x : A\}E'} \text{Extn2}
\end{array}$$

3 Subject Reduction Theorem

The *subject reduction theorem* is a property in which, for every term M , its type is preserved during reduction. This is one of the fundamental properties in typed lambda calculi such as simply-typed lambda calculus and polymorphic lambda calculus. The type judgment is defined inductively by the following rules.

We first prepare several lemmas for proving the subject reduction theorem of $\lambda FDEnv$.

Lemma 1. For any environment type E , if $E \vdash M : A$, then there exists an environment type E' satisfying that $E' \vdash M : A$ and $V(E') \supseteq V(E)$.

This lemma is proved straightforwardly by induction on the structure of type derivation tree of $E \vdash M : A$.

Lemma 2. For any environment type E , if $E \vdash M : A$, then there exists an environment type E' satisfying that $E' \vdash M : A$ and $V(E') = V(M)$.

This lemma is also proved by induction on the structure of type derivation tree of $E \vdash M : A$.

Proof. We present only a few examples of case analysis due to space limitations.

Case of Rule Var. Assume that $\{x : A\}E_1 \vdash x : A$. If you take $\{x : A\}$ as E' , then you have $\{x : A\} \vdash x : A$, that is, $E' \vdash x : A$.

Case of Rule Lam and satisfying $x \notin V(M)$. Assume that

$$\frac{\{x : A\}E \vdash M : B}{\{x : C\}E \vdash \lambda x.M : A \rightarrow B}$$

By the induction hypothesis, we have $E' \vdash M : B$ from $\{x : A\}E \vdash M : B$. By **Lemma 1**, we know $\{x : A\}E' \vdash M : B$. Then, by rule Lam,

$$\{x : C\}E' \vdash \lambda x.M : A \rightarrow B$$

is obtained. We have

$$V(\lambda x.M) = V(M) \cup \{x\} = V(E') \cup \{x\} = V(\{x : C\}E').$$

Theorem 1 (Subject Reduction Theorem). For any E , if $E \vdash M : A$ and $M \rightarrow N$ then there exists E' satisfying that $V(E) \subseteq V(E')$ and $E' \vdash N : A$.

Proof. We only show the most important part of the proof due to limitations of space.

The proof is done by induction on the structure of type derivation tree of $E \vdash M : A$.

Case of rule DLam. Assume M be $(\lambda x.M_1) \circ M_2$ and N be

$$\lambda x'.(M_1 \circ ((x'/x) \cdot (M_2 \circ id_{V(M_2)}))).$$

For some E and E' , M is typed as

$$\frac{E \vdash M_2 : \{x : C\}E' \quad \frac{\vdots \Sigma_1 \quad \frac{\{x : A\}E' \vdash M_1 : B}{\{x : C\}E' \vdash \lambda x.M_1 : A \rightarrow B}}{\vdots \Sigma_2}}{E \vdash (\lambda x.M_1) \circ M_2 : A \rightarrow B}$$

By **Lemma 2**, there exists E'' satisfying that

$$E'' \vdash M_2 : \{x : C\}E'$$

from $E \vdash M_2 : \{x : C\}E'$ in the above tree. Then, N is type as

$$\frac{\frac{\frac{\vdots \Pi \quad \{x' : A\}E \vdash (x'/x) \cdot (M_2 \circ id_{V(M_2)}) : \{x : A\}E' \quad \{x : A\}E' \vdash M_1 : B}{\{x' : A\}E \vdash M_1 \circ ((x'/x) \cdot (M_2 \circ id_{V(M_2)})) : B}}{\vdots \Sigma_2}}{\{x' : D\}E \vdash N : A \rightarrow B}$$

where Π is

$$\frac{\frac{\{x' : A\}E \vdash x' : A \quad \frac{\{x' : A\}E \vdash id_{V(M_2)} : E'' \quad E'' \vdash M_2 : \{x : C\}E'}{\{x' : A\}E \vdash (M_2 \circ id_{V(M_2)}) : \{x : C\}E'}}{\{x' : A\}E \vdash (x'/x) \cdot (M_2 \circ id_{V(M_2)}) : \{x : A\}E'}}$$

It is trivial that $V(E) \subseteq V(\{x' : A\}E)$.

4 Concluding Remarks

We proposed a simply-typed lambda calculus $\lambda FDEnv$ with strong reduction. The difference between this calculus and the previous version of the environment lambda calculus is that the domain of each environment is limited to a fixed and finite set of term variables. Thanks to this limitation, we obtained the subject reduction theorem of $\lambda FDEnv$.

Many fundamental properties remain to be investigated, for example, confluence and strong normalization. The confluence of reduction \rightarrow , sometimes called *Church-Rosser* property is that

if $M \xrightarrow{*} M_1$ and $M \xrightarrow{*} M_2$ then there exists N satisfying that $M_1 \xrightarrow{*} N$ and $M_2 \xrightarrow{*} N$.

On the other hands, the *local confluence* means that

if $M \rightarrow M_1$ and $M \rightarrow M_2$ then there exists N satisfying that $M_1 \xrightarrow{*} N$ and $M_2 \xrightarrow{*} N$.

Consider a term $((\lambda x.y)M) \circ N$. This can be reduced in the following two ways:

$$((\lambda x.y)M) \circ N \rightarrow (y \circ ((M/x) \cdot id_{\{y\}})) \circ N \xrightarrow{*} y \circ (id_{\{y\}} \circ N),$$

$$((\lambda x.y)M) \circ N \rightarrow ((\lambda x.y) \circ N)(M \circ N) \xrightarrow{*} y \circ (N \circ id_{V(N)})$$

If we introduce an equivalence relation \approx satisfying

$$M \circ id_V \approx M \approx id_V \circ M,$$

then we have local confluence modulo the equivalence relation \approx , although we cannot explain it in detail, due to limitations of space. In future, we should introduce layered equivalence relation between terms for demonstrating more theoretical properties on our calculus $\lambda FDEnv$.

Acknowledgement. This work was supported by Grants-in-Aid for Scientific Research (C) (24500009).

References

1. Abadi, M., Cardelli, L., Curien, P.L., Lévy, J.J.: Explicit substitutions. *Journal of Functional Programming* 1(4), 375–416 (1991)
2. Curien, P.L.: An abstract framework for environment machines. *Theor. Comput. Sci.* 82, 389–402 (1991)
3. Curien, P.L., Hardin, T., Lévy, J.J.: Confluence properties of weak and strong calculi of explicit substitutions. *J. ACM* 43(2), 362–397 (1996)
4. Dowek, G., Hardin, T., Kirchner, C.: Higher-order unification via explicit substitutions, extended abstract. In: *Proceedings of the Symposium on Logic in Computer Science*, pp. 366–374 (1995)
5. Seaman, J., Iyer, S.P.: An operational semantics of sharing in lazy evaluation. *Science of Computer Programming* 27(3), 286–322 (1996)
6. Nishizaki, S.: ML with First-class Environments and its Type Inference Algorithm. In: Sato, M., Hagiya, M., Jones, N.D. (eds.) *Logic, Language and Computation*. LNCS, vol. 792, pp. 95–116. Springer, Heidelberg (1994)
7. Nishizaki, S.: Simply typed lambda calculus with first-class environments. *Publication of Research Institute for Mathematical Sciences Kyoto University* 30(6), 1055–1121 (1995)
8. Nishizaki, S.: Type inference for simply-typed environment calculus with shadowing. In: *Proceedings of the Fuji International Workshop on Functional and Logic Programming*, pp. 76–89. World Scientific (1995)
9. Nishizaki, S.: Polymorphic environment calculus and its type inference algorithm. *Higher-Order and Symbolic Computation* 13(3), 239–278 (2000)
10. Nishizaki, S., Akama, Y.: Translations of first-class environments to records. In: *1st International Workshop on Explicit Substitutions*, pp. 81–92 (1998)
11. Sperber, M., Dybvig, R.K., Flatt, M., van Straaten, A. (eds.): Revised [6] Report on the Algorithmic Language Scheme. Cambridge University Press (2010)

Reliable NoC Mapping Based on Scatter Search

Qianqi Le^{1,2}, Guowu Yang¹, William N.N. Hung³, and Wensheng Guo¹

¹ University of Electronic Science and Technology Chengdu, Sichuan, China

² Chengdu University of Technology, Sichuan, China

³ Synopsys Inc., Mountain View, California, USA

leqianqi@cduet.cn, {guowu, gws}@uestc.edu.cn,

william_hung@alumni.utexas.net

Abstract. Network-on-Chip (NoC) is a promising interconnection solution for systems on chip. Mapping Intellectual Property (IP) cores onto NoC architecture is an important phase of NoC design. It affects heavily the NoC performance. In this paper, we propose a multi-objectives optimization algorithm based on Scatter Search for NoC mapping. We introduce reliability evaluation into NoC mapping in order to achieve high performance and reliable NoC architectures. Experimental results show that our algorithm achieves low power consumption, little communication time, balanced link load and high reliability, compared to other traditional evolutionary algorithms.

Keywords: Network-on-Chip (NoC), mapping, scatter search, reliability.

1 Introduction

Network-on-Chip (NoC) adopts the idea of interconnected networks from parallel computers. High quality NoC has enjoyed increasing adoption and has attracted more and more research interests in recent years [1]. NoC mapping is one of the key problems in designing NoC structures. Mapping decisions dramatically affect the NoC performance. The mapping process is shown in Fig.1.

An IP core Graph (IG) is a directed graph $G(V, E)$ shown in Fig.1(a), where v_i represents an IP core which contains a set of tasks obtained from IP cores assignment process [2]. The weight of the edge e_{ij} represents the traffic between v_i and v_j . The NoC architecture we focus on is the regular tile-based topology shown in Fig.1(b).

NoC mapping is determining the location for every IP core in the NoC structure. If there are N tiles, a total of $N!$ Solutions exist. This kind of problem is a *NP*-hard problem [3], the method used widely is obtaining the approximate optimal solution by various optimization algorithms. We propose an effective optimization technique based on scatter search (SS). Scatter search has been found to be successful for optimization problems. It embodies principles and strategies that are still not emulated by other evolutionary methods and that prove advantageous for solving optimization problems [4]. But the application of scatter search to multi-objectives optimization problems has not been fully explored. We adopt four optimization objectives for mapping, this work involves multi-objectives optimization. So we deals with the adaptation of the basic scatter search to multi-objectives optimization according to

NoC design requirements and constrains. In order to enhance the search capabilities of the algorithm, we use crossover and mutation operators borrowed from genetic algorithm (GA) instead of the pure deterministic operations. Our algorithm utilizes the advantages of both scatter search and genetic algorithm.. The experimental results demonstrate our approach outperforms other traditional evolutionary algorithms.

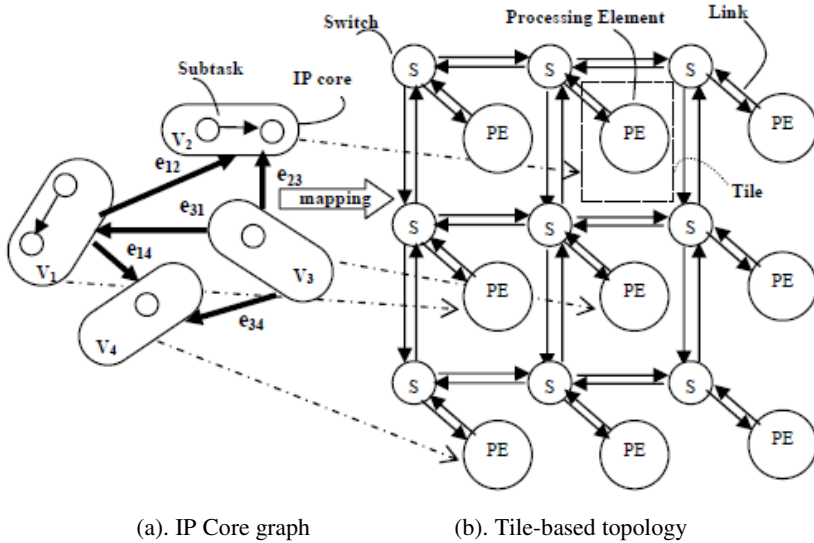


Fig. 1.

The rest of the paper is organized as follows: In Section 2, we review some related works. In Section 3, the evaluation models for NoC mapping are presented. In Section 4, we use scatter search to solve the NoC mapping problem, describing the solution representation and the optimization process. In Section 5, some experimental results are presented. Finally, we conclude in Section 6.

2 Related Work

NoC Mapping has been widely explored. In [5], Lei et al. used GA for solving the mapping problems. In [6] and [7], Hu et al. proposed a branch and bound algorithm for NoC mapping. In [8], Chou et al solved NoC mapping problems by using a linear programming approach. In [2], Nondominated Sorting GA II (NSGA II) and microGA are used for NoC mapping. In [9], Sahu, P.K et al. applied PSO to NoC mapping.

The algorithms used in the above papers can be divided into two categories: traditional optimization algorithms and evolutionary algorithms. The traditional algorithms generally have high time complexity for NoC mapping. The general evolutionary algorithms such as GA, NSGA and PSO, are easily trapped into local optimal solutions. In this paper, we propose an optimization technique based on scatter search, its search strategies prove to be efficient and not be emulated by other evolutionary algorithms.

3 Preliminaries

In practical NoC applications,, we need to coordinate multiple optimization objectives and search for Pareto optimal solutions. We adopt four indices: communication power consumption, communication time, reliability and link load balance.

3.1 Communication Power Consumption

The total power consumption of a NoC application mainly consists of the computing power consumption and the communication power consumption. The former is the power consumed when processors execute the subtasks, it is usually be valuated in the IP cores assignment phrase [2]. So in the mapping phrase, we consider the communication power consumption, which can be computed by adding up the communication power consumption of every pair of communication tiles .The power consumed by sending the data from tile i to j can be computed by (1).

$$Power_{comm}^{i,j} = P_{link} * hop_{i,j} + P_{switch} * (hop_{i,j} + 1) \quad (1)$$

Where, P_{link} is the power consumed when data pass through a link. P_{switch} is the power consumed when data traverse a switch. These two parameters are platform dependent and be set by the NoC designer [10]. $hop_{i,j}$ is the number of links between tiles i and j along a given path.

3.2 Communication Time

The total time spent in a NoC application mainly consists of the computing time and the communication time. The former is the time spent in executing the subtasks, it is usually be valuated in the IP cores assignment phrase [2]. So in the mapping phrase, we consider the latter, communication time, which is the time spent in transmitting data between tiles. Because some subtasks may be executed in parallel, we not only sum up the communication time of each subtask in the critical path, but also consider the additional time due to the congestion of parallel subtasks. The communication time model is given by (2).

$$Time_{comm} = \sum_{i,j \in critpath} (T_{link} * hop_{i,j} + T_{switch} * (hop_{i,j} + 1)) + T_{para} \quad (2)$$

Where, i and j present the subtasks in the critical path. $hop_{i,j}$ is the number of links along a given path between tiles i and j . T_{link} is the time spent in transmitting data through a link, we get this parameter from [10]. T_{switch} is the queuing time in switch buffers; its computing model is obtained from [10]. T_{para} is the penalty time added when parallel subtasks contend for the same links. If at least one of these subtasks is in the critical path, the penalty time will be added to the overall communication time.

3.3 Reliability

Switch faults because the IP cores not to communicate with the other IP cores, decreasing the system reliability. We propose a fault-tolerant mechanism which can recover the communication and balance the load. We calculate the number of redundant paths according to the traffic for every IP core, and then we build the candidate redundant paths via replacing the fault switch by its neighboring switch. Furthermore we pre-evaluate the performance of the candidate redundant paths to avoid choosing a low performance path instead of the fault path. We adopt the reliability model proposed in [10] to evaluate the structure with redundant paths, the formula is given by (3).

$$R = \prod_{\langle i, j \rangle \in IPG} \left[\sum_{s=1}^N (1 - R_s)(R_{i,j} | R_s \text{ faults}) \right] \quad (3)$$

Where, R_s is the reliability of the s -th switch in the path, i and j denote a pair of source and destination tiles in the IP Core graph. N is the number of links between tiles i and j . $R_{i,j} | R_s \text{ faults}$ denotes the redundant path when the s -th switch faults in the path between tiles i and j .

3.4 Load Balance

Link load balancing can not only relieve network congestion and queuing delay, but also avoid generating hotspots. We evaluate the load balance by computing the variance of each link load. The formula is given by (4).

$$\text{Load} = \sum_{i=1}^L \left(\text{load}_i - \sum_{i=1}^L \text{load}_i / L \right)^2 / L \quad (4)$$

Where, load_i is the load of link i , L is the total number of the links in a given NoC mapping layout. A smaller variance indicates the link load is more balanced.

4 NoC Mapping by Scatter Search

Scatter search (SS) is a class of evolutionary optimization techniques. It is based on a small population known as the reference set, which contains the best quality solutions and the most diverse solutions. These solutions are combined to construct new solutions. Based on the traditional scatter search, we propose a multi-objectives optimization SS for NoC mapping. The outline of our method is presented in Fig. 2. We discuss the implementation details of each procedure as follows:

4.1 Diversification Generation

The initial population must be a wide set of diverse solutions, so a preset number of solutions are generated randomly. We take a random permutation of $0, 1, \dots, N-1$ as

a solution, which is denoted by $(p_0, p_1, \dots, p_{N-1})$, each dimension in the solution represents a tile in the NoC topology. For example, p_i identifies the tile assigned to IP core i . N is the total number of tiles in the NoC topology.

PROCEDURE Multi-Objectives Scatter Search

```

1: BEGIN
2:   CreateDiversPopulation(Solutions);
3:   Improvement(Solutions);
4:   REPEAT
5:     RefSet = GenerateRefSet(Solutions, Ref_size);
6:     REPEAT
7:       Subset = GenerateSubSet(RefSet, Sub_size);
8:       comb_Solutions1 = CombineSolution (Subset, crossover);
9:       comb_Solution2 = CombineSolution (comb_Solutions1, mutation);
10:      Improvement(comb_Solutions2);
11:      RefSet = Update (RefSet, comb_Solution2);
12:    UNTIL (not generating new reference solutions);
13:  UNTIL (reaching iteration termination condition);
14: END

```

Fig. 2. Outline of scatter search algorithm

4.2 Improvement

Improvement is improving the solutions by local search. We select the IP core with maximum degree and arrange the IP cores communicated with it onto its adjacent tiles. Thus the length of communication path between the IP cores is reduced.

4.3 Reference Set Generation

The reference set, *RefSet*, is a collection of both high quality solutions and diverse solutions. *RefSet₁* is filled with 10 best quality solutions, while *RefSet₂* is the high diversity subset, it contains 10 solutions with the largest distance.

The quality of the solutions is decided by the Pareto domination relations of the four indices. We let the solutions undergo non-dominated sorting and rank them based on the Pareto domination relations. The solutions of the Pareto optimal front are assigned the *rank 1*, and the next set of non-dominated solutions is assigned the *rank 2*. This process is iterated until all solutions are assigned the corresponding ranks. We select 10 solutions into *RefSet₁* from the high non-dominated rank to the low.

To calculate the distance between two solutions, we compare the values of the same dimension between two solutions, if the values are same; this dimension is marked with 1, otherwise marked with 0. The sum of all dimension marks is the distance between these two solutions. The distance between solution i and *RefSet₂* is the sum of the distances between the solution i and all solutions in *RefSet₂*, its formula is given by (5).

$$distance_i = \sum_{j \in RefSet} \sum_{k=0}^{N-1} (s[i][k] \oplus s[j][k]) \quad i=1, 2 \dots popsize \quad (5)$$

Where *popsize* is the population size, *s[i][k]* is the value of the dimension *k* in the solution *i*, it represents the tile assigned to IP core *k* in the solution *i*, *s[j][k]* is the value of the dimension *k* in the solution *j*, which belongs to *RefSet*, *N* is the total number of dimensions. We select 10 solutions with maximum distance into *RefSet*₂.

4.4 Subset Generation

Subset generation is producing subsets of their solutions as a basis for creating new solutions. We adopt 2-element subsets [4], pairing each solution in *RefSet* with another solution randomly. Every pair of solutions is taken as a subset.

4.5 Solution Combination

Solution combination operates on every subset, we introduce the crossover and mutation operators into scatter search. The traditional crossover maybe generates valid solutions for NoC mapping. We propose a match crossover operator according to the NoC mapping restrictions.

Suppose there are two solutions *S*₁ and *S*₂, the start crossover dimension is *point*₁, the number of crossing dimensions is *length*, the value of dimension *point*₁ in *S*₁ is *x*. We search the same value *x* in *S*₂ and mark its dimension as *cpoint*₂, then we swap the values between dimension *cpoint*₁ and *cpoint*₂ in *S*₁. This operator continues to operate on the dimension *cpoint*₁+1, *cpoint*₁+2... until reaching the *length*. After crossover, some dimensions of the new solution are same to *S*₁, and some other dimensions same to *S*₂. For example, *S*₁=(2,5,3,8,0,7,4,1,6), *S*₂=(1,4,6,7,8,3,5,2,0), the *cpoint*₁ is 2, the *length* is 3, then the new solution *S*_{new}=(2,4,7,0,8,3,5,1,6). The new solution not only inherits partial feature from parent solutions, but also keeps its own partial feature. At the same time we ensure the new solution is a valid solution.

The mutation operates on the *RefSet* with a small probability *P*_m. We adopt swap mutation to avoid generating invalid solutions. The dimensions with mutation probability higher than the preset probability *P*_m are swapped with other random dimensions of the same solution.

4.6 Reference Set Update

We calculate four valuation indices for new solutions obtained by combination, and then compare them with the old solutions in *RefSet*. If their four indices are all superior to the old, the new solutions replace the old. Otherwise, we calculate the distances between the new solutions and the *RefSet*, the new solutions with distance greater than the old are selected into the *RefSet*, replacing the old.

5 Experimental Results

We have implemented our algorithm on the Embedded Systems Synthesis benchmarks Suite (E3S), which provide a set of NoC applications. These applications are real applications based on embedded processors from the Embedded Microprocessor Benchmark Consortium (EEMBC). We apply our algorithm to 10 typical applications of E3S and compare it with Nondominated Sorting Genetic Algorithm II (NSGA II) and Particle Swarm Optimization (PSO). The experiment results are shown in Fig.3 – Fig.6.

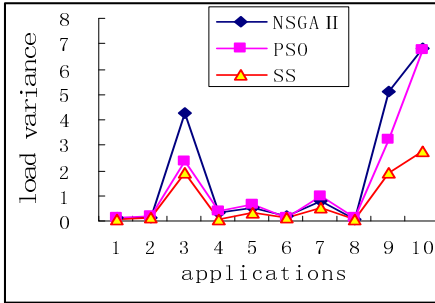


Fig. 3. Link load balance

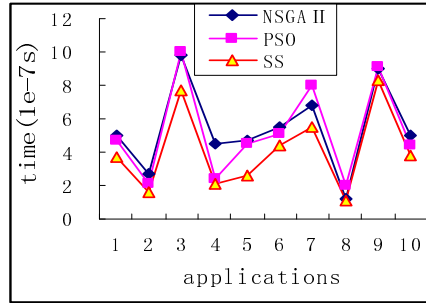


Fig. 4. Communication time

Fig.3 shows the link load balance of NoC topologies, the variance obtained by SS is less than NSGA II and PSO. This means lower network congestion and queuing delay. Fig.4 shows the communication time of applications, the time obtained by SS is less than the other two algorithms. The advantage of SS is most obvious when it is applied to the application 3 which has the most subtasks.

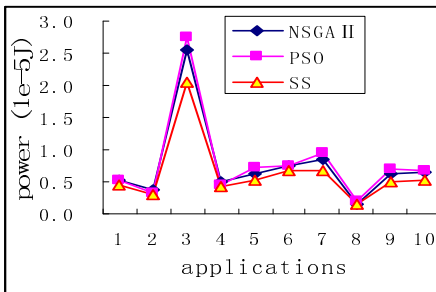


Fig. 5. Communication power consumption

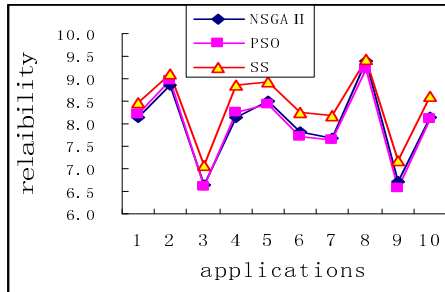


Fig. 6. Reliability

Fig. 5 shows the communication power consumption. SS improves the solutions by arranging the IP cores with heavy communication onto the neighbor tiles. This method shortens the communication path, resulting in reducing power consumption. So the power consumption obtained by SS is lower than the others. Fig.6 shows the reliability of the fault tolerance NoC architectures. The reliability of each switch is

assumed to be 0.92. The strategies of SS can avoid being trapped local optimal solutions, so the more reliable paths can be explored. As shown in the Fig.6, the reliability obtained by SS is higher than NSGA II and PSO.

6 Conclusion

We adopt four valuation indices and perform a comprehensive performance valuation for NoC mapping. We design multi-objectives optimization scatter search for NoC mapping. The search method of scatter search is based on the strategy instead of the total random search, this merit leads better solutions. The experiment results have indicated that our algorithm outperforms traditional NSGA II and PSO. We will try to apply our algorithm to the other kind of NoC structures.

Acknowledgment. This research supported in part by NSFC Program (No.60973016) of China

References

1. Marculescu, R., Ogras, U.Y., Peh, L.-S., Jerger, N.E., Hoskote, Y.: Outstanding research problems in NoC design: system, microarchitecture, and circuit perspectives. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 28, 3–21 (2009)
2. da Silva, M.V.C., Nedjah, N., Mourelle, L.M.: Power-aware multi-objectives evolutionary optimization for application mapping on NoC platforms. *International Journal of Electronic* 97, 1163–1179 (2010)
3. Garey, M.R., Johnson, D.S.: A guide to the theory of NP completeness. W. H. Freeman & Co., New York (1990)
4. Rama Mohan Rao, A., Arvind, N.: A Scatter Search Algorithm for Stacking Sequence Optimisation of Laminate Composites. *Composite Structures* 70(4), 383–402 (2005)
5. Lei, T., Kumar, S.: A two-step genetic algorithm for mapping task graphs to network on chip architecture. In: *Proceedings of the Digital System Design*, pp. 180–187 (2003)
6. Hu, J., Marculescu, R.: Energy-aware mapping for tile-Based NoC architecture under performance constraints. In: *Design Automation Conference, Proceedings of the ASP-DAC 2003*, pp. 233–239 (2003)
7. Muralimanohart, N., Modarressi, M., Tavakkol, A., Sarbazi-Azad, H.: Application-Aware Topology Reconfiguration for On-Chip Networks. *IEEE Transactions on Very Large Scale Integration Systems*, 2010–2022 (November 2011)
8. Chou, C.-L., Marculescu, R.: Contention-aware Application Mapping for Network-on-Chip Communication Architectures. In: *ICCD 2008*, pp. 164–169 (2008)
9. Sahu, P.K., Venkatesh, P., Gollapalli, S.: Application Mapping onto Mesh Structured Network-on-Chip Using Particle Swarm Optimization. In: *2011 IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*, pp. 335–336 (2011)
10. Das, R., Eachempati, S., Mishra, A.K., Narayanan, V., Das, C.R.: Design and Evaluation of a Hierarchical On-Chip Interconnect for Next-Generation CMPs. In: *HPCA 2009*, pp. 175–186 (February 2009)

Towards the Applied Hybrid Model in Decision Making: Support the Early Diagnosis of Type 2 Diabetes

Andrea Carvalho Menezes, Placido Rogerio Pinheiro,
Mirian Caliope Dantas Pinheiro, and Tarcísio Pequeno Cavalcante

University of Fortaleza (UNIFOR), Graduate Program in Applied Informatics,
Av. Washington Soares, 1321 - Bl J Sl 30 - 60.811-905, Fortaleza, Brazil
a.carvalhomenezes@gmail.com

Abstract. A hybrid model, combining Bayesian network and a multicriteria method, is presented in order to assist with the decision making process about which questions would be more attractive to the definition of the diagnosis of Diabetes type 2. We have proposed the application of an expert system structured in probability rules and structured representations of knowledge in production rules and probabilities (Artificial Intelligence - AI). The importance of the early diagnosis associated with the appropriate treatment is to decrease the chance of developing the complications of diabetes, reducing the impact on our society. Diabetes is a group of metabolic diseases characterized by hyperglycemia resulting from defects in insulin secretion, insulin action, or both.

Keywords: Bayesian network, diabetes, early diagnosis, expert systems, multicriteria.

1 Introduction

The prevalence of diabetes continues to grow and people with this chronic illness require continuing medical care and ongoing patient self-management education and support to prevent acute complications and to reduce the risk of long-term complications. On average, their medical expenditures are approximately 2.3 times higher than the expenditures would be in the absence of diabetes, not including indirect costs due to increased factors such as absenteeism, reduced productivity, and lost productive capacity due to early mortality. In this study, we are proposing an expert system in order to achieve a consistent outcome in diagnosing diabetes precocious, combining probabilities rules and the Multi-Criteria Decision Analysis [1].

2 Diabetes

Diabetes is a group of metabolic diseases characterized by high levels of sugar in blood because the pancreas does not make enough insulin (a hormone produced by

the pancreas to control blood sugar), the cells do not respond to insulin normally, or both. The chronic hyperglycemia of diabetes is associated with long-term damage, dysfunction, and failure of different organs, especially the eyes, kidneys, nerves, heart, and blood vessels. Symptoms of marked hyperglycemia include polyuria, polydipsia, weight loss, sometimes with polyphagia, and blurred vision [8].

2.1 A Public Health Problem

Diabetes set today as a global epidemic and now affects some 246 million people worldwide. It is also estimated that most people who have diabetes know their condition. There are 4 million deaths per year related to diabetes and its complications, representing 9% of world mortality [3, 19].

Diabetes is associated with a major economic and social impact for both the individual and society. Their high costs are mainly related to a high frequency of acute and chronic complications, such as higher incidence of cardiovascular and cerebrovascular diseases, blindness, kidney failure and non-traumatic amputations of lower limbs, which are causes of hospitalization, greater needs for medical care, disability, lost productivity and premature death of life [7].

2.2 Classification

There are three major types of diabetes, which causes and risk factors are different for each type [4]:

Type 1 diabetes can occur at any age, but it is most often diagnosed in children, teens, or young adults, affects some 10% of cases and is caused by the reduction of pancreatic beta cells which results in deficiency of insulin production, and often the person needs to receive daily injections of insulin;

Type 2 diabetes, which affects 90% of cases and the person has almost normal levels of insulin in the blood, but suffers a reduction in the number of receptors of this hormone in target cells, reducing the ability of these cells to absorb glucose in the blood.. It most often occurs in adulthood, but teens and young adults are now being diagnosed with it because of high obesity rates. Many people with type 2 diabetes do not know they have it. The type 2 will be addressed in this work [3].

Gestational diabetes is a preclinical stage of diabetes. A woman develops signs of hyperglycemia of varying intensity, first diagnosed during pregnancy and usually resolves in the postpartum period, but in most cases women who have gestational diabetes are at high risk of their developing diabetes after.

2.3 Criteria for the Diagnosis

High blood glucose levels can cause several symptoms, including excessive thirst, frequent urination, involuntary weight loss and hunger exaggerated and may also be other symptoms such as repeated infections in skin or mucosa, problems with blood clotting, weakness, fatigue and serious blood circulation problems. Sometimes the diagnosis is made from chronic complications such as neuropathy, retinopathy or atherosclerotic cardiovascular disease [5, 6].

Because type 2 diabetes develops slowly, some people with high blood sugar have no symptoms. Factors indicating higher risk for diabetes are obesity (Body Mass Index BMI > 25), physical inactivity, history of macrosomia or gestational diabetes, family history (mother or father) of diabetes, adults over age 45, hypertension (> 140/90 mmHg), HDL cholesterol < 35 mg/dL and/or triglycerides > 150 mg/dL, stress, previous diagnosis of polycystic ovary syndrome and cardiovascular disease, cerebrovascular or peripheral vascular set [10].

Some blood tests also can be done, such as: Fasting blood glucose level, Hemoglobin A1c test and Oral glucose tolerance test.

3 Model of Decision Making Support for Proposed

For this study we considered a diabetes data set consisted of 768 females' patient of Pima Indian. Within the set, 268 of the patients were classified as having diabetes, refers to 34.9% of all. Three attributes were considered to each patient to build the model, as body mass index (BMI), age and blood pressure, according to Table 1.

Table 1. Criteria

Criteria	Values
A - Body Mass Index (BMI)	A1 - BMI < 25 (Normal) A2 - 25 ≤ BMI < 30 (Overweight) A3 - BMI ≥ 30 (Obesity)
B - Age	B1 - Up to 30 B2 - Between 30 and 50 B3 - Starting from 50
C - Blood pressure	C1 - Hypotension C2 - Normal pressure C3 - Hypertension

We use the Bayesian networks which offer an approach to probabilistic reasoning which includes graph theory, to establish relationships between sentences and, probability theory, for assigning levels of reliability for the sentences in the knowledge base, so the system can act in situations of uncertainty. The tool Netica (<http://www.norsys.com/>) was used to build the Bayesian network, as shown in Fig. 1, where each criterion has its unconditional probability table.

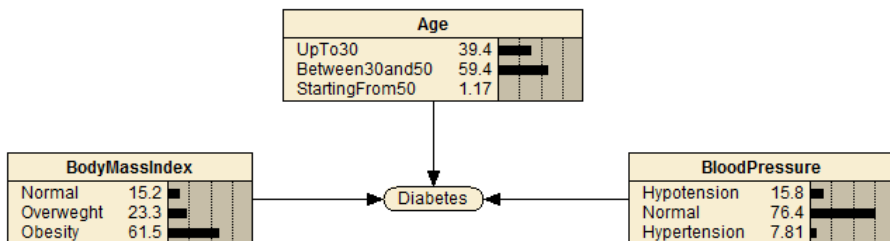


Fig. 1. Bayesian network obtained in the tool Netica

Applying the Bayes rule, we defined the conditional probability of having diabetes for each alternative, as shown in Fig. 2. For example, 17.09% is the probability of a person who is obese, aged between 30 and 50 and normal pressure (Alternative A3, B2, C3) have diabetes.

BodyMassIndex	Age	BloodPressure	Have
Normal (A1)	UpTo30 (B1)	Hypotension (C1)	0.03
Normal (A1)	UpTo30 (B1)	Normal (C2)	0.18
Normal (A1)	UpTo30 (B1)	Hypertension (C3)	0.03
Normal (A1)	Between30and50 (B2)	Hypotension (C1)	0.11
Normal (A1)	Between30and50 (B2)	Normal (C2)	0.7
Normal (A1)	Between30and50 (B2)	Hypertension (C3)	0.1
Normal (A1)	StartingFrom50 (B3)	Hypotension (C1)	0
Normal (A1)	StartingFrom50 (B3)	Normal (C2)	0.02
Normal (A1)	StartingFrom50 (B3)	Hypertension (C3)	0
Overweight (A2)	UpTo30 (B1)	Hypotension (C1)	0.12
Overweight (A2)	UpTo30 (B1)	Normal (C2)	0.8
Overweight (A2)	UpTo30 (B1)	Hypertension (C3)	0.11
Overweight (A2)	Between30and50 (B2)	Hypotension (C1)	0.48
Overweight (A2)	Between30and50 (B2)	Normal (C2)	3.12
Overweight (A2)	Between30and50 (B2)	Hypertension (C3)	0.44
Overweight (A2)	StartingFrom50 (B3)	Hypotension (C1)	0.02
Overweight (A2)	StartingFrom50 (B3)	Normal (C2)	0.11
Overweight (A2)	StartingFrom50 (B3)	Hypertension (C3)	0.01
Obesity (A3)	UpTo30 (B1)	Hypotension (C1)	0.67
Obesity (A3)	UpTo30 (B1)	Normal (C2)	4.36
Obesity (A3)	UpTo30 (B1)	Hypertension (C3)	0.61
Obesity (A3)	Between30and50 (B2)	Hypotension (C1)	2.64
Obesity (A3)	Between30and50 (B2)	Normal (C2)	17.09
Obesity (A3)	Between30and50 (B2)	Hypertension (C3)	2.39
Obesity (A3)	StartingFrom50 (B3)	Hypotension (C1)	0.09
Obesity (A3)	StartingFrom50 (B3)	Normal (C2)	0.58
Obesity (A3)	StartingFrom50 (B3)	Hypertension (C3)	0.08

Fig. 2. Conditional probability of having diabetes for each alternative

The Multi-Criteria Decision Making system (MCDMs) implemented through software Hiview, improve the way that complex problems are viewed and presented to decision makers, helping the choice of the best decisions in an uncertain setting. Using the method MACBETH (Measuring Attractiveness by a Categorical Based Evaluation Technique) [2, 6], implemented by Hiview, the decision maker quantifies the relative attractiveness of options, comparing two alternatives at a time, using entirely verbal judgments, in order to create a robust and consistent decision model. The difference in attractiveness between two alternatives is judged as no difference, very weak, weak, moderate, strong, very strong or extreme. For example, the alternative (A3, B2, C2) has very strong

attractiveness compared to the alternative (A1, B1, C1). As judgments are entered into the software, it automatically verifies their consistency, in order to compose a robust matrix of judgment. Finally, when the matrix is complete, the numerical scale is generated, consistent with all the decision maker's judgments. The alternatives will be exported to the tool Expert SINTA for composing the knowledge base of an expert system, which represents the information that is used by an expert, to aid in the diagnosis of diabetes [13].

4 Expert Systems

Expert systems are computer programs that solve certain problems as a human expert would solve, under certain conditions [11, 12]. The Expert SINTA was the computational tool used to simplify the implementation work of the expert system for diagnosis of diabetes. The software Expert SINTA was created by a group of scholars at the Federal University of Ceará (UFC) and the State University of Ceará (UECE), called Group SINTA (Applied Intelligent Systems). This tool uses a model of knowledge representation based on production rules, with conditions in the style IF... THEN..., with the possibility of including logical connectives relating the attributes in the scope of knowledge and the use of probabilities, using an inference engine shared and probabilistic treatment of the production rules. The Expert SINTA use the architecture described in fig. 3.

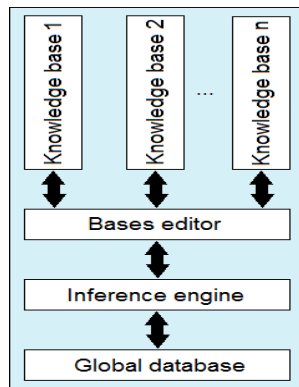


Fig. 3. Architecture of Expert SINTA

The first step was defining the variables and their values according to the criteria defined in table 1 and the objective. Then the questions were registered to be performed by the system to the user, regarding the given attribute. The next step was to define the rules to model human knowledge, according to the alternatives in Fig. 2, for example, in rule number one, if the user is obese, his age between 30 and 50 years and normal pressure, as in the worst case, it is possible to say that he has diabetes. These steps can be seen in fig. 4.

An expert system seeks to achieve conclusions for certain objectives asking questions to the user. Whenever one of these objectives is reached, or when they exhaust all the possibilities, the Expert SINTA display a window with the results, and how they came to that conclusion, as shown in fig. 5.

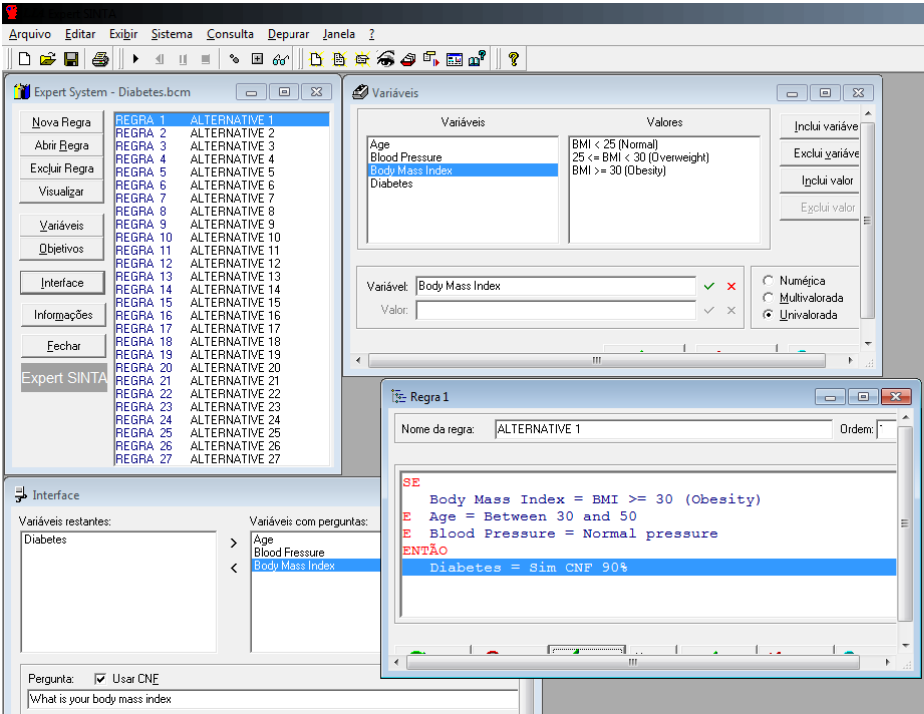


Fig. 4. Expert Sinta

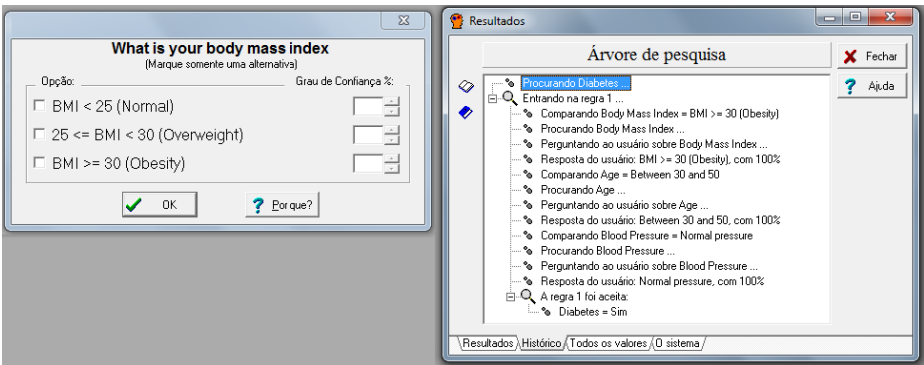


Fig. 5. Example of data entry of variables and the analysis of the result

5 Conclusion and Future Works

The present study demonstrates that the probability rules and the degrees of confidence defined in the software Hiview can be useful to compose the knowledge base that will be processed by the machine inference of the expert system, which uses this information to make the diagnosis. This article is part of a study with the proposal to contribute to the development of automated diagnosis with quality. Such methodology induces the decision maker to establish decision criteria to assess the control events relative importance, relying on the judgments of a panel of experts. The methodology described in this paper will be applied to another data set with more evaluation criteria. The model here proposed needs major investigation regarding the correct value of the threshold which ultimately determines the quality of the results. We also intend to implement a new analysis of the model through the inclusion of values in the utility node of the influence diagram. This way, it will be possible to make inferences on the network aiming the definition of the diagnosis by the application of only the most attractive questions. Other methodologies may be applicable in the definition of the diagnosis of diabetes type 2 [14, 15].

Acknowledgments. The second author is thankful to the National Council of Technological and Scientific Development (CNPq) for the support received.

References

1. American Diabetes Association, <http://www.diabetes.org/>
2. Bana e Costa, C.A., Ensslin, L., Correa, E.C., Vansnick, J.C.: Decision support systems in action: integrates application in a multicriteria decision aid process. *European Journal of Operational Research* 133, 315–335 (1999)
3. BRASIL. Ministério da Saúde. Secretaria de Atenção à Saúde. Departamento de Atenção Básica. Diabetes Mellitus: Caderno de Atenção Básica - n.º 16. Série A. Normas e Manuais Técnicos. Brasília, DF (2006)
4. de Castro, A.K.A., Pinheiro, P.R., Pinheiro, M.C.D.: An Approach for the Neuropsychological Diagnosis of Alzheimer's Disease: A Hybrid Model in Decision Making. In: Wen, P., Li, Y., Polkowski, L., Yao, Y., Tsumoto, S., Wang, G. (eds.) *RSKT 2009*. LNCS, vol. 5589, pp. 216–223. Springer, Heidelberg (2009)
5. de Castro, A.K.A., Pinheiro, P.R., Pinheiro, M.C.D.: A Hybrid Model for Aiding in Decision Making for the Neuropsychological Diagnosis of Alzheimer's Disease. In: Chan, C.-C., Grzymala-Busse, J.W., Ziarko, W.P. (eds.) *RSCTC 2008*. LNCS (LNAI), vol. 5306, pp. 495–504. Springer, Heidelberg (2008)
6. Figueira, J., Greco, S., Ehrgott, M.: *Multiple Criteria Decision Analysis: State of the Art Surveys*. Springer, Boston (2005)
7. Harris, M.I.: Diabetes in America: Epidemiology and scope of the problem. *Diabetes Care* 21(sup. 3), 11–14 (1998)
8. Holmes, D.M.: The person and diabetes psychosocial context. *Diabetes Care* 9(2), 194–206 (1986)
9. Larichev, O.I., Moshkovich, H.M.: *Verbal decision analysis for unstructured problems*. Kluwer Academic Publishers, The Netherlands (1997)

10. Nunes, L.C.: A Hybrid Model for Supporting Diagnosis of Psychological Disorders. Master Thesis - Graduate Program in Applied Informatics, University of Fortaleza (2010)
11. Tamanini, I.: Improving the ZAPROS Method Considering the Incomparability Cases. Master Thesis - Graduate Program in Applied Informatics, University of Fortaleza (2010)
12. The Expert Committee on the diagnosis and classification of diabetes mellitus. Report of the Expert Committee on the diagnosis and classification of diabetes mellitus. *Diabetes Care* 20, 1183–1197 (1997)
13. World Health Organization. Definition, diagnosis and classification of diabetes mellitus and its complications: report of a WHO consultation. Geneva, World Health Organization, 59 p. (1999)
14. Roy, B.: Paradigms and Challenges. In: Figueira, J., Greco, S., Ehrgott, M. (eds.) *Multiple Criteria Decision Analysis: State of the Art Surveys*. Series: International Series in Operations Research & Management Science, vol. 78, XXXVI, pp. 3–24 (2005)
15. Atkinson, M.A., Maclaren, N.K.: The pathogenesis of insulin dependent diabetes. *N. Engl. J. Med.* 331, 1428–1436 (1994)

Replacement Study of the Retention Time in Chromatography Economic Analysis

Ping Shen^{1,*} and Shibing You²

¹ Biological Engineering School, Wuhan Polytechnic, Wuhan, China, 430074

² Economics and Management School, Wuhan University, China, 430072
sping99@126.com, sbyou@whu.edu.cn

Abstract. The principle of the chemical chromatography analysis method was introduced into the chromatography economic analysis in order to show the role of the retention time in economic research. Based on the plate theory, the economic replacement of the concept of the retention time was carried out. With the assumption of separation, the bonus behaviors of the listed companies were studied to model the separation process. The results show it is feasible to classify the complex economic phenomena using the chromatography analysis method.

Keywords: Chromatography economic analysis, Retention time, Separation simulation.

1 Introduction

In our previous work, the chemistry chromatography analysis was successfully applied to the social science field by You et al [1]. The foundational work was carried out and the results indicated that the distribution ration was key role in composition separation [2]. The difference of distribution can be arrived by adjusting the retention time. Therefore, how to replace the retention time is a key step in chromatography economic analysis research.

For the gas chromatography, the chromatographic column will be divided into the same plates in plate theory. In the plates, a part of spacing is occupied by stationary phase and the other part of spacing is occupied by mobile phase as the carrier gas [3]. When the mixtures to be separated flow into the plates with the pulse of carrier gas, if the diffusion of the components along chromatographic column can be ignored and the isothermal distribution of any component in stationary phase and mobile phase is linear, the total volume and export outflow of the component in any column plate are shown in Table 1. The number of theoretical plates per column (n) is 6 ($n = 6$). r is the plate number and the value is 0,1,2, ..., $n-1$. The distribution ratio of the sample to be separated is $p:q$ ($k' = p/q$).

When two different compositions go through the same column, the peak will be obtained with different volumes of carrier gas, namely the time to get to the peak will

* Corresponding author.

be different [4]. If there are thousands of the theoretical plates in the chromatography column and the difference in distribution ration of the components is very small, the completely separated peaks can be obtained in order to realize the complete separation of the all components [5-8]. Therefore the different components can be collected at various times and the completed separation of the components will be realized with different distribution ration. Under a certain chromatography, all the substances have exclusive retention time. When the process separation finishes, the category of the substance is determined according to the values of the retention time, namely, the separation and qualitative analysis for the samples can be conducted using the different retention time [9-11].

Table 1. The component distribution in the chromatographic column

Volume of plate carrier gas	r						Outflow at column export
	0	1	2	3	4	5	
1	p	q	0	0	0	0	0
2	p2	2pq	q2	0	0	0	0
3	p3	3p2q	3pq2	q3	0	0	0
4	p4	4p3q	6p2q2	4pq3	q4	0	0
5	p5	5p4q	10p3q2	10p2q3	5pq4	q5	0
6	p6	6p5q	15p4q2	20p3q3	15p2q4	6pq5	q6
7	p7	7p6q	21p5q2	35p4q3	35p3q4	21p2q5	7pq6
8	p8	8p7q	28p6q2	56p5q3	70p4q4	56p3q5	28p2q6

2 Classification of the Listed Company Group by Using the Chromatography Analysis Method

Similar to the definition of the retention time in chemistry, the retention time in the chromatography economic analysis can be defined as the time (tR) that the subjects begin to interact until the peak occurs. Here, there is something to explain. First, the target is chosen according to the separation of the components in specified field. The optimal choice can deduce the internal mechanism and realize the goal to separate effectively. Second, the “specified time” is not also the concept of the general time, but also the specific factors or variables in expand economics. In this paper, the cash bonus behaviors of the listed company are considered as the research object to replace the retention time. Based on the plate theory in chemistry, the type of the listed companies will be separated according to a certain assumption.

2.1 Assumption of the Separation Process

Assumption 1: There are many independent companies in a group of similar listed ones. In each period, all companies obtained net profit. There are profits available for the distribution to the shareholders after the companies removed the statutory provident fund and statutory public welfare fund, which meets the conditions of the bonus. At the same time, at each period, some companies have dividends, the other have no dividends.

Assumption 2: Assume the stationary phase as the retained profits trend, the mobile phase as the company's profit-sharing trend. When a group of similar listed companies enters into a chromatography column with pulsed mobile phase, if the listed companies are chosen to share bonus, the company will enter into next plate, otherwise, the company will retain in the primary plate until the bonus is selected. Here, the ratio of the amount of the companies to choose bonus and the amount of the companies having no bonus is constant. The constant value is distribution ration (k') of listed companies group. If the distribution ration of a particular type of listed companies is $p:q$ ($k'=p/q$), the proportion of the dividends will be determined as $p/(p+q)$.

Assumption 3: the amount of the period to share bonus is determined based on the amount of the plates. From assumption 1, at each period, the listed companies obtain the distributable profits to shareholders. Therefore, the bonus behavior can be selected in the company at each period, namely, the distribution behavior will be carried out in the stationary phase and mobile phase. Hereby, if each period is considered as one plate, the period can be self-selected according to the research needs.

Assumption 4: when the companies obtain the distributable profits to shareholders, they meet the condition to share bonus to the shareholders. Assuming the profits as carrier gas, it is considered as the internal mechanism to progress the bonus behaviors. Because the profits can't be obtained continuously and are obtained from different department similar to a pulse type, the form of the profits is in agreement with the characteristics of the carrier gas in chemistry.

2.2 Simulation of the Separation Process

The following simulated process is based on the assumption of separation process above. Assuming the plate number of the chromatographic column is 6, the sample to be classified is a mixture of A and B, which are two different type listed companies. The distribution ration of A-type companies is 1:1; B-type is 1:4. The starting number of the listed company group in A and B is one unit ($m_A=1$, $m_B=1$). The specific target is tested at column export when the 6 bonus behaviors have been carried out. Here, the target is pointed to the last number of the company to share bonus. When the mixed sample enters into the chromatographic column, the simulated results to the capital distribution in two type listed companies are shown in table1.

Under the situation of 6 times cash bonus distribution and $t = 6$, the form indexes of A-type listed companies turn out the sample, when $t = 11$, it comes the sample peak value. B-type ones also show sample, when $t = 7$, it comes the peak. After they come to the peak, the numbers of the emerged samples gradually reduce and present an inverted U type law. Based on the definitions of retention times in the chromatographic economic analysis method, the retention time of A-type firms is 11 ($t_{RA}=11$), the retention time of B-type is 7 ($t_{RB}=7$). We get the outflow curve of public companies group by testing the emerging companies' number and corresponding to the time. (fig. 1).

In conclusion, there are different types of distribution ratio of different listed companies groups after the emerging sample peak time of the same column. Under the situation of thinking chromatographic analysis, a kind of type listed companies group only get a kind of retention time. The different groups have different time. Therefore, we can separate different types of listed companies groups by the different retention time of listed companies outflow, mixing the sample of the same column. Meanwhile, we can confirm characteristic of listed companies on the basis of determining the retention time. That is available for retention time to separate and analyze qualitatively the components in the mixed samples.

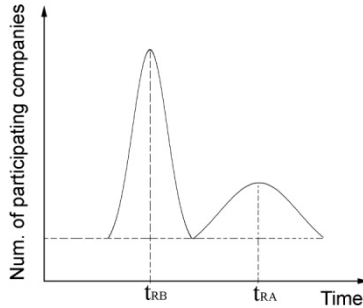


Fig. 1. Outflow curves of the mixtures

3 Significance of Retention Time in the Chromatographic Economic Analysis

In the field of chemistry, retention time is from the emerging sample to a peak value, retention time is to display the external phenomena of economic law and other in the field of economic.

The tower plate numbers depend on the research need to choose their selves. When the theoretical plate number of the column is 6, it appears the phenomenon of cross-peaks in the flow curves of different types of distribution ratio of different groups of listed companies. When the numbers of the theory tower plate are increased, we can get the complete separated peaks of listed companies group (Fig. 2) to enhance the separation effect of different types of listed companies group. Of course, we should increase the observed periods if it is necessary to increase the number of theoretical plates, that is, the separate accuracy needs many observations results. The bigger of the observed periods set, the more of the tower plate, especially, the difference of the distribution ratio is amplified evidently when the observations go on to make the separate result better.

There are 2 important significances in the formula (1). First, if the distribution ratio of listed companies group is k' , we calculate the retention time t_R . Second, we access to retention time t_R by detecting the out flowing status of listed companies group at the column outlet to figure out the distribution ratio.

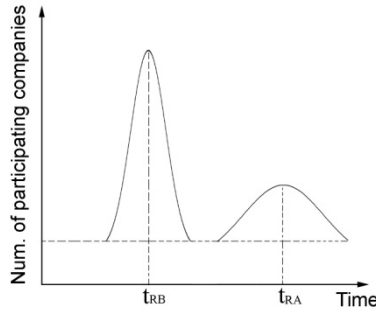


Fig. 2. The outflow curve of listed companies groups

In above example, the retention time of A-type Company is 11, B-type is 7, and it is different the retention time because of different distribution ratio. The difference of company group has different retention time, and the different retention time has different distribution ratio, therefore, the retention time is a very good point to separate the listed companies group. We confirm the rate of return and separate the 2 types of companies by the outflow curve of retention time. That supply a important referred guideline for the market investor to choose good invest cases.

4 Thinking the Case about Replacing the Retention Time

4.1 Operation

Whether dividend or not is a common phenomenon, we can choose it by the relevant information and data in the stock .On the one hand, we can get the real retention time to do the relevant empirical test and confirm the research in the future. On the other hand, we can discover the difference between theoretical and empirical to analyze the reasons. They also support the further improvement of the chromatographic method of economic analysis.

4.2 Change of Retention Time by Actual Situation

However, we can set the content and quantity of the reasonable theoretical plate number in the field of economics based on the practical need to pursuit a good separation. For example, the number of set theoretical plates is for listed companies the time of continuous period of distributing bonus, if investors believe it is acceptable to divide 6 times in the 8 years, the number of theoretical plates can be set 6. And under that condition we do the separated validity test. If investors believe it is acceptable to divide 8 times in the 8 years, the number of theoretical plates can be set 8. In practice analysis, the number of the plates can be set flexibly based on the characteristics of the samples and effectiveness and aims of the separations. This is the greatest advantage of the chromatography economic analysis. Under certain

chromatography, anything has a sole retention time. Therefore, the separation and qualitative analysis of the mixed samples can be conducted using the retention time simultaneously. This is another advantage, which is scarce in the economic analysis method.

4.3 The Use of Retention Time Is Not Sole

The most importance of the retention time is to separate and qualitative analyze the mixed samples according to the retention time. In above cases, it can be found that the retention time plays other role in economics. On hand, we can classify the company for long-period or short-period ratio of return based on the retention time. On the other hand, for investors to stock market, the retention time is a crucial parameter to select the investment period except for the long-period and short-period return. According to the estimation of the retention time, the investors will select more listed companies to invest.

5 Conclusions

First, on the base of tower plate theory we classify listed companies groups to achieve the replacement of chromatography retention time in the field of economic. Although the performance form is different between the economic models of plate theory and the chromatography plate theory, both of them has the same nature, that is, the separated key is the difference of retention times led by the different distribution ratio. Second, this model is unique and innovative itself and has the potential practical significance. The separation of listed companies plays an important refer role for the market investors in the selection of invest component. Last, we will go on researching further for practical significance of chromatographic economic analysis.

References

1. You, S., Wu, B., Shen, P., Mei, M., Su, Z.: Theoretical prospect of the complex economic phenomenon classification method innovation – reference and thinking based on the chemical “chromatography analysis method”. *Tong Ji Yu Jue Ce* 7 (2011)
2. You, S., Wu, B., Mei, M.: Foundation research on applying the chromatography analysis principle in the economic field – taking the financial securities investment and fast moving consumer goods industry as examples. *Tong Ji Yu Jue Ce* 11 (2011)
3. Shen, P., Zhang, P., Mao, K., Li, G., You, S.: Chromatography economy analysis method replacement series research, distribution ratio. *Tong Ji Yu Jue Ce* 17 (2011)
4. Su, L., Zheng, Y.: Chromatography analysis method. Tsinghua University Press, Beijing (2009)
5. Liu, Z.-L., Li, C.-J.: Quantitative Analysis on the Creation of Dualistic Interdisciplinary Science. *Studies in Dialectics of Nature* 20 (2004)
6. Zhou, Y., Wan, A.T.K., Xie, S., et al.: Wavelet analysis of change-points in a non-parametric regression with heteroscedastic variance. *Journal of Econometrics* 159 (2010)

7. Ren, Q., Li, S., Qiao, D., et al.: Application of key factor analysis method for elastic coefficient of highway freight transportation, Chengdu, China (2010)
8. Noorossana, R., Eyvazian, M., Amiri, A., et al.: Statistical monitoring of multivariate multiple linear regression profiles in phase i with calibration application. *Quality and Reliability Engineering International* 26 (2010)
9. Yu, F.J., Tsou, C.S., Huang, K.I., et al.: An economic-statistical design of \bar{x} control charts with multiple assignable causes. *Journal of Quality* 17 (2010)
10. Florez Lopez, R.: Effects of missing data in credit risk scoring. A comparative analysis of methods to achieve robustness in the absence of sufficient data. Houndmills, Basingstoke, Hants, RG21 6XS, United Kingdom (2010)
11. You, S., Bao, L., Zhong, S., Guan, X., Wang, L.: Chromatography economy analysis method replacement series research, trays theory. *Tong Ji Yu Jue Ce* 1 (2012)

Study on Frustration Tolerance and Training Method of College Students

Naisheng Wang

Department of Social Sciences Education, Shandong Traffic and Transport University,
Jina, 250032, China
naishengwang@126.com

Abstract. The frustration tolerance is the character structure of the essential factors, and it is the mental health scale. It has a direct impact on students learning motivation in the process, belief, self consciousness and self motivated. And it also affects the students' intelligence, ability of the normal play and promotion; and it even affects the students' psychological health level and personality perfection. This paper discusses the ability of setbacks and failure reasons; and next, a more detailed analysis of the path to guide college students to overcome the setback.

Keywords: The College Students, Frustration Tolerance, Training Path.

1 Introduction

Good endurance to frustration is the mental health standard. If one's frustration tolerance is weak, he would be easily controlled in the face of setbacks by negative emotional, and easily to fall into bad habits and inextricably bogged down in mud; and he could be easily negative, blame everyone and everything but not oneself, even commit suicide in the face of difficulties. The person whose frustration tolerance is weak always his psychology is negative, and unhealthy. Because the university students' life experiences are shallow, the ideal color is thick and self evaluation is high, and thus the life expectancy is also high. It is easy to make them feel the real life with their subjective expectations gap. After lots of investigation and interview, we found that most college students can't objectively vision of the future, but as the idea of a lyric that the idea of " the future". Then, such common phenomenon appears in the university campus: the students action, effort, or receive serious rebuff. Because of the subjective to the setbacks is more sensitive, so the more intense frustration response. Therefore, discussion of contemporary college students' setback 's general performance ,and study how to train the students' frustration endurance, it will be of great advantage to strengthen college students' psychological health education, improve psychological quality and health level.

2 Frustration and Frustration Tolerance

Human social life practice shows, if one exists, it will produce a variety of needs, and will be due to unmet need or goals cannot be achieved and the frustrations come.

In psychology, the frustration feeling is a kind of motivation drive below, in achieving the goals of the action in the process, he encounters insurmountable or unable to overcome obstacles and interference, making its motive cannot come true, As the result, the tension and emotional reaction generate. From some kind of meaning, the setback is a social component, it is a partner in life, people are likely to encounter frustration whenever and wherever possible. For example, several years of hard study, the dream into the ideal university, but failed to do so; a very powerful player, wins the gold medal in the competition, because the game play is not good and fail in official nowhere; a scientist, after hundreds of, thousands of times, even thousands of experiments, has not found the ideal results; an entrepreneur in the overnight dissipate one's fortune; a millionaire suddenly encountered natural calamities and man-made misfortunes and lost property or relatives. These are the blow in life that can not be avoided. Therefore, understanding setback, learning how to face the setbacks and actively resolving them, are each person's life-long tasks. However, how to deal with setbacks situation as well as to the ability of bearing frustration, the differences between people are not the same. Some people can suffer serious setbacks, never yield in spite of reverses, rises in adversity, until to succeed; some people have slightly setbacks that depression, they would be unable to get up after a fall; some people can endure life, study, work in a serious setback, but cannot bear pride was a bit hurt and so on. All of these are decided by the individual and the frustration tolerance. The so-called frustration tolerance, also known as the ability to endure setbacks, refers to the individual to setbacks can endure, with an acceptable degree of size. It is for people to adapt to the setback, resistance to deal with setbacks one kind of ability.

3 The Analysis of the Psychological Reasons Caused the Frustration

Most of the college students' psychological frustrations are caused by the various contradictions which exist between the subjective desire and objective reality that are about individual needs and motivation. Therefore, analysis of the contradictions between the college students' subjective desire and objective reality will find that college students have a sense of frustration causes. The college students' frustration psychological causes are diverse, but they mainly lie on two aspects: the extrinsic and intrinsic factors.

3.1 The Extrinsic Factors Causing the College Students' Frustration Psychology

The extrinsic factors causing the college students' frustration psychology mainly refer to the environmental factors, which are comprehensive factors of natural and social conditions and they do not transfer in college students' subjective desire. With the rapid social change and development, the college students are difficult to adapt to the individual and social environment changes, unable to meet individual and psychological setbacks. Such as natural disasters, the death of a loved one, family

adversity, unexpected trouble, social corruption, social security, a large gap between rich and poor natural environment and social environment factors may lead to frustration generation. A student of the enormously proud of one's success to grind to continue their education, but the family financially forced their work, in order to reduce the burden on the family; and the student who wants to give outstanding performance in the school sports meeting, but because of the unexpected injury, he as a specialty students can not achieve the goal of sports etc.

The school environmental factors are also a major source to cause college students psychological frustrations. The college students will produce a sense of difference after they entered the school, finding with their desired a larger gap. And there are intramural competitions in the school life such as the party joining, the student cadre competition, and various scholarship assessment etc. They may be easy to make the college students not adapt to the situation. In addition, with the expansion of higher education, the school authorities only pay special attention to the scale, while neglecting the educational quality of colleges and universities, education lags behind the phenomenon exists, and have not created a good quality or space and atmosphere for the students' ability training. For the students, the ideal and the reality of a strong contrast, the psychological imbalance, resulting in spiritual loneliness, loneliness and strongly not to adapt. All of these lead to frustration appeared. In addition, the lives of students and learning are mainly in the school. The school management level, teaching effect quality, campus cultural atmosphere, interpersonal differences, these are the great influence on the college students. In the management of the system, most schools also use the traditional management method, with the authority, control, punishment method of management of college students. The school forms a management system that leads the serious conflicts with the students' motivation. On the teaching contents and methods, the students generally feel their learned knowledge is not enough or practical to the society. Many students feel that the class of monotone, dull, hollow, rigid .The main reason is the curriculum setting, teaching materials, teaching content, teaching method, basically the same, the test form a single, and the teachers can not adapt to the requirements of the actual teaching, echoing what the books say, shutting the door to teaching, and leading the separation between theory and practice. In the aspect of interpersonal relations, some college students, due to the lack of necessary learning, training, and environmental constraints, can not reach normal interpersonal and their sports ability is limited, and leading to the result of isolating themselves. In the course of time, they are shy to participate in group activities and begin to doubt their ability, and then they can not get the correct guidance, and at last form a kind of psychological setback. For example: there is a first grade student in a university in Beijing has, because there a fellow student with the dormitory can play the flute, many people admire, but he felt no proficiency in a particular line, and others could despise himself. And the he many times ask the stage director to let him take a year left, but did not tell the class director the reason. The class director again and again asked him the reason why he should do so, he secretly told the chairman, he wanted to go home to learn flute, lest the others despised him.

3.2 The Intrinsic Factors of the College Students' Frustration Psychology

The intrinsic factors of the college students' frustration psychology mainly refer to an individual's physiological, psychological and other factors caused by interference or disorder. Young college students are in the process of development of psychology coming from immature to mature. Their will, emotion and sense of self are in a stage of development. In the face of emotional entanglements, will frustration, the changes in his life, they would appear to be confound and emotional fluctuations, so that they can not quickly adjust the stability state of mind to adapt to the environment, thereby forming the degree of frustration. Its basically

3.2.1 The Individual Feature Constraints

As a result of the college students' individual characteristics constraints, so that some college students produce frustration because of the individual characteristics of dissatisfied. They could complain themselves, as the body too fat, not engaged in the dance or to participate in an exercise and; and because the character is too introverted, withdrawn by the students, not gregarious, inadvertently left out in the cold, and the depression often with negative emotions can affect the students' learning and living, and these will be prone to frustration, if it is serious enough it still may form a vicious spiral.

3.2.2 The Conflict of the Motivation

In the daily study and life, the college students often simultaneously produce two or more than two kinds of motivations. But due to some restrictions to "not both fish and bear's paw", then the students' motivational conflict and a fierce ideological struggle would come, they can not simply choose a certain motivation than another. If the psychological contradiction lasts too long or only a motivation is met and other motivation is blocked, and then it will cause frustration. Such as taking part in the entrance exams for postgraduate schools and anterior good jobs are always prone to conflict; wanting to like heterosexual friends, but again being afraid delay time, and their studies abandoned and so on, are easy to make college students get into trouble.

3.2.3 High Expectations

The college students as a result of the thought too simple, too simple, lacking some experience, can correctly estimate their own ability and level and develop higher or unable to reach the goal. If the target can not be realized they will have a strong sense of frustration. Particularly when they failed to realize soberly this ambition, leading higher expectations, more frustrations. At present, the "employment first, and career second, undertaking last" idea is from the front guide students in employment during the early to maintain normal mentality, they face social employment pressure, making reasonable employment goals and expectations, alleviating obtain employment contradiction and setbacks. These are the effective methods.

4 To Guide the College Students to Correctly Deal with Setbacks

The purpose to study on the causes of the College Students' frustration is to regulate, constraint behavior among the college students. The students undergo the setbacks, will generally produce nervous mood. At the same time, they will have to eliminate or alleviate the tense state of behavioral responses, that setback adaptation. The ability of adaptability to the setbacks always differs from man to man. In general, a strong man than the weak with higher frustration tolerance, after repeated failures, or repeatedly beating people with higher frustration tolerance, extrovert than introverts, frustration tolerance to. In addition, but also by individual ideological level, intellectual level difference effect. Frustration adaptation can be roughly divided into two categories: positive and negative adaptations. The university students who are actively adapting could take effective action, facing the reality, to relieve tension. Negative adaptation also known as maladaptive, negative adaptation of university students the action probably be of no avail, can alleviate pain and frustration, they generally take on attack, defense, frustration object damage or destroy, will give personal or social cause greater harm. Therefore, ideological and political workers, if they can correctly guide the college students who are suffering the setbacks, encourage the positive adaptation, eliminate its negative adaptation. These kinds of doing often can play a multiplier effect. The method has the following kinds:

4.1 A Correct Understanding of Their Own

Young college students should establish correct world outlook, outlook on life and values, to found their own advantages, and they should be sure of their ability and value, but also to be comprehensive and correct evaluation of their own shortcomings and defects. On this basis, they also should enhance the life of courage and confidence, with correct, positive view of society and life of the people. The correct cognition formed on the frustration event, accepting defeat with common state of mind, and treating the frustration as to promote personal development opportunity, these kinds of attitudes can reduce the pressure caused by the frustration.

4.2 The Rational Spiritual Catharsis

The reasonable spiritual vent can dissolve bad mood and alleviate the psychological pain. They can talk their grievances and feelings to the intimate friend or respectable teacher after frustrations. The college students should learn to talk, soothe the soul and emotional wounds are healing; on the other hand it is also a help to find a solution to the problem. Because of the reasonable talk can be stuffed in the heart melancholy, pain and injustice, most incisive to pour out, to get others to understand and help clean the soul, haze, regain the balance of your mind and life fulcrum. In addition, it can also travel to change the environment, to participate in various activities to divert attention for the negative influence of dilute the setbacks.

4.3 Enhancing the Frustration Tolerance

The frustration tolerance refers to the ability of individual resistance to attack or failure. It is an important component of psychological qualities. The personal frustration resistance ability is not constant. It is related to their own physiological quality, cognitive factors, personality factors, frustration and social support influence frequency, and to their frustration experience directly as well. The social changes require the young college students get the life and frustration experiences as soon as possible and keep calm and maintain higher frustration endurance in the new face of setbacks. The college students who are mental healthy, should recognize the setbacks are part of life in real life from theory and experience of life, and they are a normal phenomenon which can be inevitable or escape. In the face of setbacks, they ought to swallow humiliation and bear a heavy load and keep the striving spirit, and try best to overcome setbacks situation to eliminate the negative factors caused by the frustration. And take efforts in changing the frustration into motivate and positive factors, so as to consciously cultivate frustration tolerance ability, and in order to achieve good social adaptation.

References

1. The "Social Psychology," 134th pages. Renmin University of China Press (1986)
2. Chen, Y.: Internet's Influences on the University Students. Guangzhou University Journal (Social Sciences Edition) 07 (2002)
3. Information on
<http://wenku.baidu.com/view/196c0060caaedd3383c4d3e8.html>
4. Information on
<http://gfgyjx.blog.163.com/blog/static/59598415200819114626627/>
5. Management Dictionary of Psychology, 176th pages. The People's Liberation Army Press (1990)

On Statistical Analysis and Management Countermeasures of Occupation Burnout for College Teachers

Youcai Xue and Jianhui Pan

School of Science, Zhejiang University of Science and Technology, 318 Liuhe Road,
Hangzhou 310023 China
xueyoucai@126.com

Abstract. Occupation burnout has recently been a fascinating topic in research fields like management and psychology. Due to work pressure and other reasons, college teachers' occupation burnout catches more concerns than ever. As teachers' occupation burnout appears generically and yet not being optimistic in Zhejiang Province, we conduct recently a questionnaire survey on teachers' occupation burnout at a local college, and then manipulate the SPSS software to carry out the statistical analysis. The final results reveal that teachers' occupation burnout in this college is quite negative and alarming. Furthermore, we would propose several corresponding management countermeasures for occupation confidence reconstruction.

Keywords: College Teachers, Occupation Burnout, Deindividuation, Statistical Analysis, Management Countermeasures.

1 Introduction

Much research display that teaching is one of the most stressful occupations, and teachers are among the high risk population of occupation burnout. Relevant survey conducted in mainland of China indicates that parts of teachers are lack of the sense of occupation identity and the sense of belonging. Occupation burnout therefore has become a realistic and alerted problem, and has since resulted in great losses of teachers and widespread tired teaching phenomenon in teaching positions, which arguably has become a negative factor attributing to impacts in the institutional development of education[1,2].

With survey and statistical analysis on working pressure and occupation burnout of teachers at a local college in Zhejiang province, this research analyzes in difference both occupation burnout and working pressure through various dimensions by making contrastive comparison of college teachers with different background such as gender, age and professional title. The obtained result suggests the current state of college teachers' occupation burnout is not rather optimistic. It furthermore puts forward some management countermeasures for coping with college teachers' occupation burnout.

2 Objectives of Research and Methodology

Sample used in this study is randomly selected from teachers working at the same college. Moreover, total 200 questionnaires have been delivered, and recycled 191 questionnaires whereupon 191 questionnaires are valid. It therefore could evidently estimated that the recovery rate of questionnaire delivery is 95.5% .

Study uses the anonymous Q&A style, and collects statistical datum after delivering and recollecting questionnaires at the scene in time. Study uses a similar questionnaire refereed to Maslach's questionnaire form in 1986 specifically designed for teachers' occupation burnout, and just make some appropriate adaptations and maximum limit obtains a desirable one. The formalization of questionnaire essentially consists of 15 items, and moreover judges occupation burnout in three dimensions, namely, the emotional exhaustion (in 5 questions), a low sense of achievement (in 5 questions) and the personality (in 5 questions). Underlying questionnaires take the method of seven points scoring, namely, scoring from 1 to 7, wherein between 11 and 15 entitled the reverse scoring. All summed score in each dimension amounts to the total number corresponding to scores of five relevant questions within one group. According to the literature [4], each dimension has a critical value, namely, assigns a threshold value of 25 for emotional failure, assigns a threshold value of 11 for deindividuation, and assigns a threshold value of 16 for low sense of achievement. Occupation burnout is defined as mild provided one summed score exceeds one critical value in any of preceding three dimensions. Occupation burnout is similarly defined as moderate provided two summed scores exceed critical values in some two of preceding three dimensions. Extremely, Occupation burnout is similarly defined as highly weary provided all summed scores exceed critical values in three dimensions. The study analyses all datum through the application of statistical software SPSS17.0.

3 Results and Related Analysis

General situation in this college teachers' occupation burnout

From table one, the ratio for occupation burnout of teachers at this college evidently exceeds 55%, and is a lot higher than 33.12% proportionally accounted for occupation burnout of senior professionals, a data published by Economic Information Daily on 2006, and is as well higher than 37.44% accounted for occupation burnout of personals in Zhejiang province.[5] According to datum at table two, generally speaking, one could deduce many teachers have lost enthusiasm in teaching positions due to over-loaded working pressure at college from obviously poor performance in the dimension of deindividuation. In the dimension of emotional exhaustion, only a small fraction of teachers could not experience fun and happiness of working. Teachers at this college are often lack of working enthusiasm because they have been imposed in overriding stress and frustration. Some 25.1% of college teachers investigated in this research have revealed a low sense of achievement from datum at table two.

Table 1. General situation in occupation burnout

	Total numbers(N)	Not any (O.B.)	Mild (O.B.)	Moderate (O.B.)	High (O.B.)
General situation in Occupation Burnout(O.B.)	191	86	68	33	4
Percentile (%)		45.0%	35.6%	17.3%	2.1%

Table 2. General situation in various dimensions

Dimension	Total number(N)	Exceeds threshold value	Percentile(%)
Emotional exhaustion	191	14	7.3%
Deindividuation	191	90	47.1%
Low sense of achievement	191	48	25.1%

Occupation burnout of College teachers and difference analysis of demographic variables

Occupation burnout of college teachers and gender differences

The study on gender differences in occupation burnout of college teachers should require T test on independent samples, and mainly make the comparison between male college and female college teachers.(see table 3 for further details) From table three, we could find out that occupation burnout of male college teachers is slightly higher than that of female college teachers. Between these two sexes, there is a significant difference in terms of variances but, there is no significant difference in terms of means. On the issue of deindividuation, female college teachers are less serious in this condition than male college teachers. In this dimension, variances between male and female college teachers exist significant differences but, means between these two are ultimately no significant difference. On the issue of low sense of achievement, male and female college teachers are almost on neck and neck. There is no significant differences in terms of variances and means.

Occupation burnout of college teachers and age differences

For college teachers at different ages, our study on occupation burnout applies variance analysis of a single factor, and takes one decade as a successive unit. From table four, in the general situation of the dimension of occupation burnout, college teachers over 55 years old are most obvious in occupation burnout and, college teachers whose age between 36 and 45 years old have minimal scores in terms of means. In the dimension of low sense of achievement, college teachers over 55 years old and college teachers between 25 and 35 years old reveals obviously on its impact but, college teachers whose age between 36 and 45 are less obvious at it. The result is consistent with the reality. In fact, at the age from 36 to 45 years old, this fraction of college teachers has already adapted to the role of teaching and, at the same time, in this age college teachers usually could carry out most ingenious researches and reap other fruitful benefits and rewards. On the other hand, due to working pressure, it is too easy to generate the emotional exhaustion because at this age college teachers long for a strong sense of personal achievement. This is also consistent with underlying results in table five. However, in the dimension of deindividuation, college teachers whose ages are between 25 and 35 years old show more enthusiasm than other age groups. In the dimension of emotional exhaustion, college teachers whose ages are between 25 and 35

years old have a maximal mean, that could plausibly be related to prepare inadequately because they have just started to understand how to teach and conduct researches. In the difference analysis, only in the dimension of deindividuation, college teachers of different ages have a significant difference in terms of means. However, three items such as: emotional exhaustion, low sense of achievement and occupation burnout do not exist any significant difference in terms of means.

Table 3. Inspection on Effects of different genders for occupation burnout of college teacher through three dimensions

Dimension	Gender	N	M	SD	F	sig	t	sig
Emotional exhaustion	male	111	18.31	5.687	5.782	0.017	-0.256	0.798
	female	80	18.50	4.739				
Deindividuation	male	111	12.40	4.942	6.028	0.015	1.743	0.083
	female	80	11.31	3.651				
Low sense of achievement	male	111	15.21	4.376	2.403	0.123	1.242	0.216
	female	80	14.43	4.179				
Occupation Burnout	male	111	45.78	11.094	11.06	0.001	1.063	0.289
	female	80	44.28	8.508				

Note: N denotes by total number, M denotes by mean, SD denotes by standard deviation

Occupation burnout of college teachers and differences of teaching periods

Through variance analysis, differences in teaching periods (see table 5 for details) indicate that in three dimensions and in general situations of occupation burnout there is no significant difference in terms of means. Relatively speaking, in the dimension of general situations, college teachers whose teaching periods are between 11 and 15 years have much better performance than other college teachers. In dimension of low sense of personal achievement, college teachers whose teaching periods are amid 4 and 6 years, 7 and 10 years could enjoy working joys much better than other college teachers. For them, without frustration, one could reward a good sense of achievement from teaching at college.

Occupation burnout of college teachers and differences in professional titles

In the dimension of emotional exhaustion, datum in table six enable us to find out lecturers are obviously impacted, while in the dimension of deindividuation, datum also tell us lecturers and assistant professors are better in good condition than associate professors and full professors. In the dimension of low sense of achievement, in the highest mean score, full professors exhibit a more severe reduction on personal accomplishment, while datum for college teachers with lower professional titles are relatively good.

Occupation burnout of college teachers and subject differences

From table seven, it should in the dimension of emotional exhaustion be not hard to discover college teachers of Humanity and Arts department are much more affected than college teachers from departments of Engineering, Science and other fields. However, in dimensions of deindividuation and occupation burnout, teachers from Engineering department are not obvious than teachers working in other fields. In dimension of low sense of achievement, teachers from Science department are in serious and problematic state. The P value indicates that in three dimensions and general condition of occupation burnout, teachers from different fields are not in significant difference in terms of means.

From table nine, that teachers in different states of work pressure have different occupation burnout. College teachers who earn a highest degree of work pressure have exactly a total score of 51.47, and convey that they have more serious problems in occupation burnout than college teachers from other fields. Moreover, college teachers with less and no pressure have a better performance in teachings than those under the pressure of other three types.

Table 4. F test in three dimensions for occupation burnout of college teachers in different ages

Dim. Age	N	M	S.D.	F	P
Emotional exhaustion	25-35	58	18.50	0.033	0.992
	36-45	61	18.21		
	46-55	49	18.45		
	>55	23	18.43		
Deindividuation	25-35	58	10.98	2.760	0.044
	36-45	61	11.51		
	46-55	49	13.20		
	>55	23	12.83		
Low sense of achievement	25-35	58	15.24	0.774	0.510
	36-45	61	14.31		
	46-55	49	14.78		
	>55	23	15.70		
Occupation Burnout	25-35	58	44.69	0.709	0.548
	36-45	61	44.05		
	46-55	49	46.22		
	>55	23	46.96		

Table 5. F test in three dimensions for occupation burnout of college teachers in different teaching periods

Dim. Age	N	M	S.D.	F	P
Emotional exhaustion	3<	27	17.07	0.709	0.587
	4-6	53	18.58		
	7-10	41	18.00		
	11-15	36	19.11		
	>16	34	18.82		
Deindivi dulization	3<	27	11.04	1.564	0.186
	4-6	53	12.32		
	7-10	41	11.15		
	11-15	36	13.28		
	>16	34	11.62		
Low sense of achievement	3<	27	15.30	1.021	0.398
	4-6	53	14.45		
	7-10	41	14.07		
	11-15	36	15.81		
	>16	34	15.21		
Occupation Burnout	3<	27	43.44	1.254	0.290
	4-6	53	45.28		
	7-10	41	43.27		
	11-15	36	47.92		
	>16	34	45.65		

Table 6. F test in three dimensions for occupation burnout of college teachers with different professional titles

Dimension	N	M	S.D.	F	P	
Emotional exhaustion	Ass. Pro.	22	16.73	6.072	1.095	0.353
	Lec.	73	19.03	5.185		
	Asso. Pro.	49	18.35	4.965		
	Ful. Pro.	47	18.21	5.413		
Deindividuation	Ass. Pro.	22	10.68	4.110	1.904	0.130
	Lec.	73	11.36	3.928		
	Asso. Pro.	49	12.76	5.134		
Low sense of achievement	Ful. Pro.	47	12.60	4.548	0.130	0.942
	Ass. Pro.	22	14.59	3.554		
	Lec.	73	14.77	5.026		
	Asso. Pro.	49	15.18	3.296		
Occupation Burnout	Ful. Pro.	47	14.87	4.426	0.978	0.404
	Ass. Pro.	22	41.91	10.419		
	Lec.	73	45.16	9.351		
	Asso. Pro.	49	46.29	10.378		
	Ful. Pro.	47	45.47	10.752		

Table 7. F test in three dimensions for occupation burnout of college teachers with different subjects

Dimension		N	M	S.D.	F	P
Emotional exhaustion	Engineering	45	17.22	6.190	1.313	0.338
	Science	59	18.37	5.378		
	Humanity & Art	59	19.08	4.717		
	Others	28	18.82	4.651		
Deindividuation	Engineering	45	11.44	3.487	0.242	0.867
	Science	59	12.12	5.285		
	Humanity & Art	59	12.08	4.550		
Low sense of achievement	Other-s	28	12.07	3.962	0.792	0.500
	Engineering	45	14.64	4.749		
	Science	59	15.59	4.403		
	Humanity & Art	59	14.49	4.309		
Occupation Burnout	Others	28	14.57	3.202	0.732	0.534
	Engineering	45	43.22	11.049		
	Science	59	45.97	10.509		
	Humanity & Art	59	45.66	10.046		
	Other-s	28	45.46	7.471		

Table 8. The statistical table on pressure distribution

Pressure distribution	Very	Moderate	Normal	Less normal	Not any
Quantity	19	67	66	35	4
Percentile(%)	10.0	35.1	34.6	18.3	2.0

Table 9. The statistical table for different stress conditions in group

		N	M	S.D.
Occupat ion Burnout	very	19	51.47	8.072
	moderate	67	44.03	9.576
	normal	66	45.44	10.462
	Less normal	35	43.51	10.755
	Not any	4	43.50	6.758

4 Management Countermeasures for Coping with Teachers' Occupation Burnout

Coping strategies for Occupation burnout of college teachers

Face pressure, self-extraction [6]. The pressure in working settings always exist, and there are occupation pressure in every professionals. When pressure generated, one should think at first about self- extraction, and take positive coping strategy instead of negative avoidant, and alleviate pressure through various ways.

Maintain a good attitude and control the mood effectively. Pressure in everyday life and working settings is inevitable, and not solvable to each individual. Nevertheless, all that you could make is to maintain a positive personal attitude. So one should keep a good personal attitude, think over beautiful and remarkable things, and learn to take effective control of personal emotions and especially not to take negative personal moods to working circumstances.

Keep personal role and strengthen the sense of social responsibility. One always does personally not forget the social role and sense of social responsibility, and could effectively resist to depersonalization issues. One should be aware of his social role, and should take upon relevant social responsibility.

Excavation working fun and cultivate working enthusiasm as much as possible[7]. One should fully excavate meaningful aspects of work, and cultivate his interest and enthusiasm in work, and positively face difficulties during working. These tactics are often helpful to cope with the pressure, difficulty and discomfort in working circumstances.

Organizational intervention countermeasures for occupation burnout of college teachers

In response to cope with occupation burnout, it is evidently not enough relying solely on individual's effort. Organization and management should in both provide positive and effective measures to prevent and interfere with occupation burnout [7].

Under people-oriented goals, school leaders should ultimately make positive and effective policies, and promote development of teachers' profession. The predominant consideration should be taken into account to teachers who are working in the first-line so that they could feel the concern and love from their schools, and acquire the considerate and concern from their organizations. Policy makers therefore should require every attempt to improve college teachers' sense of achievement and sense of honor, and implement actively on the prevention of occupation burnout.

School should regularly invite psychological experts host ‘psychology in the classroom’ for teachers. It should target to help teachers to release psychological pressure, formulate a positive attitude on work, and as well assist to develop healthy psychology. It should create, develop and cultivate a harmonious campus and relevant cultural circumstance respectively.[9] Culturally harmonious circumstance should basically include cultivation of common values, spiritual style cultivation, moral and behavioral development, and creation of common values. Campus culture could play several motivated functional such as coagulation function, incentive function, radiation function and standard function, etc. It could sufficiently make cohesions among college teachers, and thus lead them to concern and love their working settings. Incentive function reflects in the arousal of teachers' enthusiasm, initiative and creativity. Radiation function at last mainly influences behaviors of teacher, strengthens individual's role and sense of responsibility, and reduces individual behaviors through different kinds of cultural activities.

School should improve management habits such as cultural motivation, spiritual encouragement and material reward in order to stimulate the enthusiasms of college teachers. Caring teachers should become an important school management, and hence personally set an example by taking part in management. It as well should intend to and understand teachers' ideological and psychological fluctuation, and take into accounts of college teachers' psychological health, ideological transition, physical condition, family life, etc. Only when college really concerns its staffs, teachers would in return heartly consider college as their household, and put their own heart and soul into powerful contribution of college's development. College teachers would furthermore to a relative degree alleviate occupation burnout.

References

1. Cedoline, A.J.: Job burnout in public education: Symptoms, causes, and survival, pp. 17–22. Teachers College Press, New York (1982)
2. Yao, L.: Teachers' stress management, pp. 38–45. Zhejiang University Press (2005)
3. Wu, X.: Teachers' occupation burnout prevention, pp. 30–45. China Light Industry Press (2008)
4. Hong, L.: Teachers' job stress Management, pp. 120–130. China Light Industry Press (2008)
5. Emotional exhaustion, diminished personal accomplishment... job burnout diffused in working settings, <http://news.hsw.cn> the economic information daily (July 18, 2006) the data from <http://www.chinaHRD.net>
6. Meng, H.: Occupation psychology, pp. 88–105. China Light Industry Press (2009)
7. Sun, L.: Job burnout problem in competition times, vol. 3, pp. 77–100. China Economic Publishing House, Beijing (2009)
8. Jian, X., Li, J.: On the Study of sport culture. Journal of Zhejiang University of Science and Technology 2, 59–62 (2006)
9. Xue, J.: On sp: University Teachers's Professional Burnout in Xinjiang. Journal of Bintuan Education Institute 2, 43–47 (2011)
10. Angerer, J.M.: Job burnout. Journal of Employment Counseling 9, 98–107 (2009)

Development of a RFID Multi-point Positioning and Attendance System Based on Data Comparison Algorithm

Wenyu Zhao¹, Jun Gao², Xiaotian Liu¹, and Yaping Wu³

¹ Department of Electronic Informaiton, Northeastern University at Qinhuangdao 066004, Qinhuangdao, China

² Electronics and Information Technology Institute, Northeastern University at Qinhuangdao 066004, Qinhuangdao, China

³ Liren College Experimental Center, Yanshan University 066004, Qinhuangdao, China
gaojunq@163.com

Abstract. In order to avoid the misreading of electronic tags in the multi-point positioning and attendance system based on RFID, the application of data comparison algorithm between the backup data and collected data is put forward in this paper. Basic principles and implementation method of this algorithm are specified. Also, system architecture and test result are discussed. The measurement results show that the algorithm can improve the reliability of the multi-point positioning and attendance system, thus it has great value to promote.

Keywords: RFID, data comparison algorithm, positioning, attendance system.

1 Introduction

At present, the technology of Internet of Things (IOT) has been widely promoted, and the multi-point positioning and attendance system is one typical application of IOT technology, as a new generation of attendance system, its main feature is that it can real-time record the attendance automatically, which avoids the drawbacks of the traditional attendance systems, only recording the personnel state of 'in' and 'out'. The multi-point positioning and attendance system can be more fully utilize the advantages of RFID technology. In the scope of this system, there are several areas which are under monitoring at the same time. The attendance system uses positioning technology to determine the target in which specific regions, and then completes attendance. The adjacent positions may be quite close to each other. The location of the three specific positions in Figure 1 illustrates the application scenarios of the system.

For the previous multi-point positioning and attendance systems, the reader reads the information of the electronic tags and sends to the back-end database to set it as attendance results. However, there is such a problem that it is easy to read other neighbor labels as its own, because the RFID positioning is not definitely accurate. : In order to avoid such a problem, this paper puts forward the application of data comparison algorithm between the backup data and collected data.

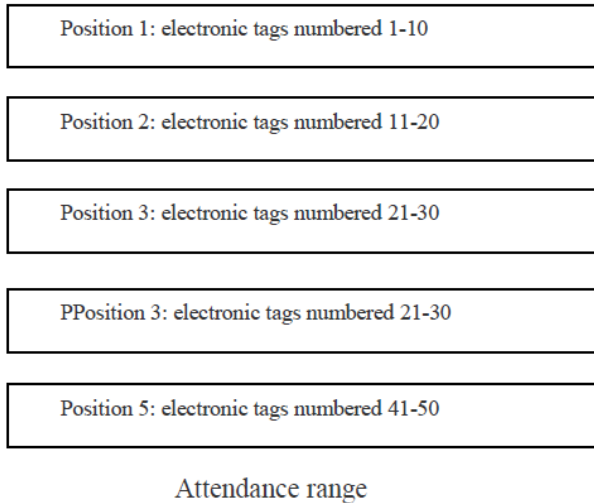


Fig. 1. System application scenarios

2 System Architecture

As is shown in Figure 2, the dynamic attendance system consists of four main parts: the RFID tags, RFID readers, the transmission network and the server. The RFID tags are made of microchips, it can transmit signals that can be detected by the RFID reader, the RFID reader is responsible for reading IDs of the electronic tags and storing them, while the transmission network is responsible for the transmission of data which the RFID reader has stored, and the server is responsible for data processing and database updating after the attendance is over. The system is easy and cheap to establish for the devices are not complicated.

In school, it needs to record the attendance result in every classroom, and every student is required to wear his/her own RFID tag. When one classroom needs to record attendance, the server sends 'action' orders to invoke the RFID reader through the transmission network. Then the RFID reader reads the RFID tags in the classroom and stores the IDs of the corresponding labels, after the attendance is over, the reader sends the data back to the server, and the server will process it background and update the database.

Although the dynamic attendance system can do the attendance automatically, it has several problems at the same time. For one thing, the adjacent classrooms are so close to each other that the RFID reader may read inaccurately, for another, the students' tags in other classrooms may also be read and recorded. Therefore, some algorithm is needed to exclude this error. Through many researches, the date comparison algorithm has been proposed to deal with these problems.

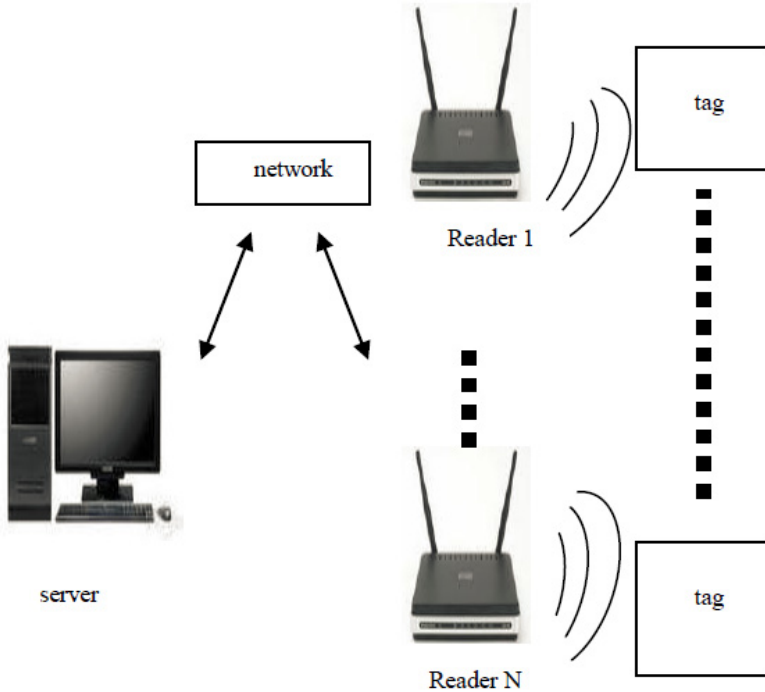


Fig. 2. The dynamic attendance system

3 Data Comparison Algorithm

The core of data comparison algorithm is that the electronic tags ID in each legitimate attendance scope are stored in a central database in advance. After the reader finishes reading, and sends the data of ID to the server, the server compares the tag ID transmitted by the reader with the tag ID stored in the database, discarding the interfere tags ID. To introduce data comparison algorithm and the basic principles, take the attendance system used in classrooms for example.

The server sends 'action' orders to invoke the RFID readers to record attendance, in the meantime, the server loads the 'should-be' student list to the cache from the central database,

After recording is finished, the RFID reader sends the 'real-be' student list back to the server, then the server compares the two lists. If a student's electronic tag ID exists both in the 'real-be' student list and the 'should-be' student list, then the server determines the student has attended; however, if a student's electronic tag ID only exists in the 'real-be' student list, exclude in the 'should-be' student data list, then the server determines the student is not belonged to the classroom, which is the interference, discards it. If a student's electronic tag ID only exists in the 'should-be' student data list, exclude in the 'real-be' student data list, then the server determines the student is absent,

Transfer the results from the cache to the attendance database to finish database updating, accomplishing attendance recording.

Flow chart of the algorithm is shown in Figure 3.

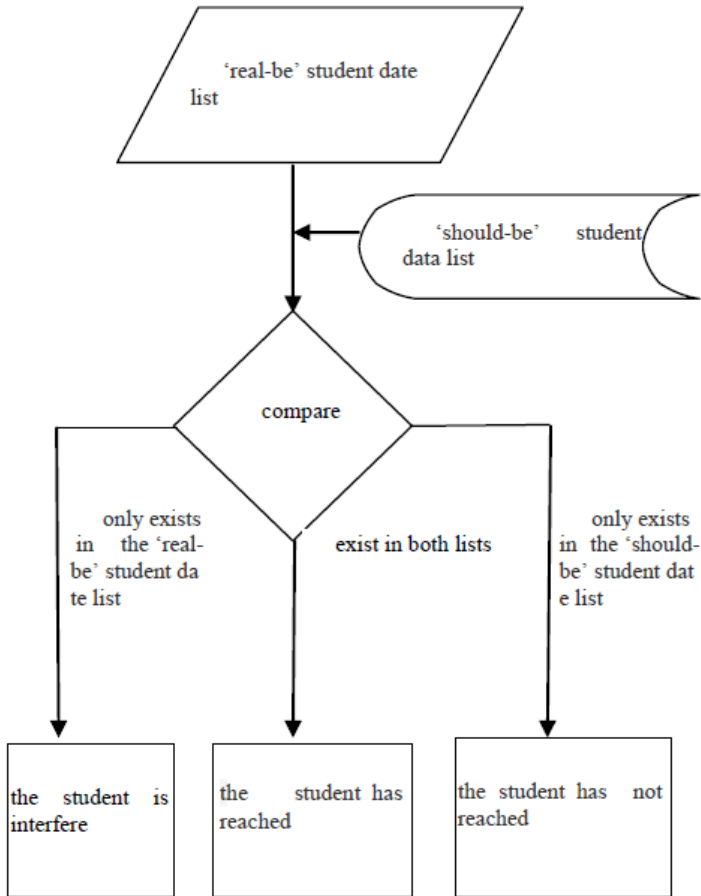


Fig. 3. Flow chart of the algorithm

Partial code of the algorithm:

```

Private void distinguish (int [] a, the string c)
// this function is used to distinguish the attendance in
position "c", where the reader locates
{
int [] buffer_dis;
try
{
buffer_dis = return_NO (c);

```

```

// this function returns an array, the return value is
'c' position's 'should-be' student list
foreach (int m in a)
// this code is used for comparing
{
// If it is the should-be student, call 'set_attend'
function
foreach (int n in buffer_dis)
if (m == n)
{
set_attend (m);
// to record the current location of card number 'm'
}
else
{
;
// Do nothing (for future possible modifying)
}
}
}
catch
{
;
// Set aside to handle exceptions
}
}

```

Even if the readers of the classroom returns the 'real-be' student list at the same time, the server can still complete comparing very fast for C language runs quickly.

Since the ID's length of electronic tag is 32-bits, so large, in order to speed up the comparison, it takes partition processing; each partition consists of 8-bits.

When the comparison is on, it traverses the first partition firstly, if it equals with the data in the database, then traverses the second partition, and then the third partition, the fourth last. Else, if it does not equal with the data in the database if one partition, then finishes comparing and determines the state of the student. This solution will greatly increase the data throughput of the server.

Besides, while one data transmitted by the RFID reader is compared with the data stored in the central database, it takes the binary tree traversal search algorithms to increase the searching speed.

4 Test Results

To test the multi-point positioning and attendance system based on RFID, we have utilized the system in a classroom which can accommodate 200 students. The results show that the system can complete recording attendance in about 10 seconds, and can

effectively eliminate the interference of other RFID tags outside the classroom, efficiently monitoring the students within a specific scope. To further enhance the accuracy of the attendance result, the system is set to record attendance three times. If a student's ID emerges two times in three, then determine the student has attended. After 1000 times testing, such an attendance system proves to record attendance accurately.

With regard to the comparison efficiency, we use the partition algorithm and the normal traversal algorithm to compare the same data and record the spending time. The result shown in Table.1 proves that the partition algorithm can be ten times faster than the normal algorithm.

Table 1. Speeding time of different algorithm

Algorithm sort	Speeding time(second)
Partition algorithm	4.1s
Normal algorithm	42.3s

5 Conclusion

The attendance system based on RFID and data comparison algorithm proposed in this paper has resolved the electronic tags misreading of the past multi-point positioning and attendance system. The test results show that the data comparison algorithm can effectively exclude attendance data errors caused by inaccurate RFID positioning. This new type of multi-point positioning and attendance system has more reliability and practicality. And it can be applied to schools, conferences, and factories that have many scopes of monitoring attendance.

References

1. Lahiri, S.: RFID Sourcebook. IBM Press, New Jersey (2006)
2. Shepard, S.: RFID Radio Frequency Identification. MacGraw-Hill (2005)
3. Yorozu, Y., Hirano, M., Oka, K., Tagawa, Y.: Electron spectroscopy studies on magneto-optical media and plastic substrate interface. *IEEE Transl. J. Magn. Japan* 2, 740–741 (1987)
4. Kassem, A., Hamad, M., Chalhoub, Z., El Dahdaah, S.: An RFID attendance and monitoring system for university applications. In: *IEEE ICECS*, pp. 851–854 (March 2011)
5. Lim, T.S., Sim, S.C., Mansor, M.M.: RFID based attendance system. In: *IEEE ISIEA*, pp. 778–782 (2009)
6. Yeop Sabri, M.K., Abdul Aziz, M.Z.A., Mohd Shah, M.S.R., Abd Kadir, M.F.: Smart Attendance System by suing RFID. In: *APACE*, pp. 1–4 (2007)
7. Qaiser, A., Khan, S.A.: Automation of Time and Attendance using RFID Systems. In: *IEEE ICET*, pp. 60–63 (2006)
8. Wahab, M.H.A., Mutalib, A.A., Kadir, H.A., Mohsin, M.F.M.: Design and development of portable RFID for attendance system. In: *IEEE INFRKM*, pp. 172–178 (2010)

9. Chen, W.-D., Chang, H.-P.: Using RFID technology to develop an attendance system and avoid traffic congestion around kindergartens. In: IEEE UMEDIA, pp. 568–572 (July 2008)
10. Chalasani, S., Boppana, R.V.: Data architectures for RFID transactions. IEEE Trans. Ind. Informat. 3(3), 246–257 (2007)
11. Sarac, A., Absi, N., Dauzère-Pérès, S.: A literature review on the impact of RFID technologies on supply chain management. Int. J. Prod. Econ. 128, 77–95 (2010)

High Throughput Constraint Repetition for Regular Expression Matching Algorithm

Kunpeng Jiang, Julong Lan, and Youjun Bu

National Digital Switching System Engineering
& Technological Research Center
Zhengzhou, Henan, China
bjay371@163.com

Abstract. To investigating high throughput pattern matching of regular expressions, This paper present a novel NFA-based architecture. In this paper, two theorems were proved and were used to prove the correctness of the algorithm. Our approach was based on three basic module to construct NFA which easily were reused in a FPGA or ASIC. Our approach is able to process many symbols per one clock cycle, and to run at high frequency. Due to FPGA constraints, the throughput of our approach achieved 512Gbps. The latency of our approach is lower than 2ns.

Keywords: Network, FPGA, NFA, Regular expression.

1 Introduction

It is well known that a deterministic finite automata (**DFA**) or a nondeterministic finite automata (**NFA**) is generated by a regular expression and a regular expression is accepted by a finite automaton, deterministic or nondeterministic [2]. However, since the Field-programmable gate array (FPGA) or Application-specific integrated circuit (ASIC) has a lot of computing resources, the running time of NFA can be reduced to $O(1)$ by using the computing resources. Variety of FPGA (or ASIC) applications utilize the NFA algorithm to process regular expressions.

This paper is based on the research of Clark, and extended the idea of his decoder-character to generate a decoder matrix. According to the symbol of pattern, a vector is selected from the decoder matrix. After processing vectors, generate the results of regular expression matching.

2 Vector-And Algorithm

2.1 The Base of Algorithm

Let $L(M)$ be a regular expression over an alphabet $\Sigma \cup \{(\,,\,), \emptyset, \cup, *\}$ ($\Sigma = \{s_1, s_2, \dots, s_{|\Sigma|}\}$) is matched inside a text $T_n = t_0 t_1 \dots t_{n-1} \in \Sigma^*$, and

$P_m = p_0p_1 \dots p_{m-1}$ is the symbols collocation in the $L(M)$, where $n \gg m$. Considering the characteristics of the network, without loss of generality, we assume that use 8-bits to represent a symbol, that is, $|\Sigma|=256$. A text $T_l = t_0t_1 \dots t_{l-1} \in \Sigma^*(l \leq n)$ is accepted in each clock cycle. Taking into account the match will use data accepted in two or more clock cycles, we can define the last data as $t_{-j}t_{-j+1} \dots t_{-1}$ and t_{-j} is the $(\lceil \frac{j}{l} \rceil \cdot l - j)$ symbol in the $\lceil \frac{j}{l} \rceil$ cycle.

we start with the notion of

Definition 1 (Enable function)

$$E_i^j = \begin{cases} 1 & \text{if } i = 0 \text{ or } t_{j-i}t_{j-i+1} \dots t_{j-1} = p_0p_1 \dots p_{i-1} \\ 0 & \text{if } i \neq 0 \text{ and } t_{j-i}t_{j-i+1} \dots t_{j-1} \neq p_0p_1 \dots p_{i-1} \end{cases}$$

for $0 \leq i < m, 0 \leq j < l$.

and the other notion of

Definition 2 (Result function)

$$R_i^j = \begin{cases} 1 & \text{if } t_{j-i}t_{j-i+1} \dots t_j = p_0p_1 \dots p_i \\ 0 & \text{if } t_{j-i}t_{j-i+1} \dots t_j \neq p_0p_1 \dots p_i \end{cases}$$

for $0 \leq i < m, 0 \leq j < l$.

The algorithm proposed by the paper consists of the following modules:

Decoder Matrix A. Because the symbol assumed by the paper is represented by 8-bits, the text $T_l = t_0t_1 \dots t_{l-1}$ which is a $l * 8$ bits matrix can be decoded as a $l * 256$ bits matrix A. the bits of every row is 256 bits, that is, the element of the Decoder Matrix A is:

$$A_{i,j} = \begin{cases} 1 & \text{if } t_i = j \\ 0 & \text{if } t_i \neq j \end{cases}$$

for $0 \leq i < l, 0 \leq j < 256$.

Symbol Comparing. The module transform two column vectors to a column vector. The one column C_i is coming from the Decoder Matrix according to the symbol p_i . Let $C_i = [A_{j,k}]$, where $0 \leq i < m, 0 \leq j < l, k = p_i$.

Theorem 1 (Symbol comparing). If $C_i = [A_{j,k}]$, for $0 \leq i < m, 0 \leq j < l, k = p_i$ and Decoder Matrix A then

$$C_i^j = \begin{cases} 1 & \text{if } t_j = p_i \\ 0 & \text{if } t_j \neq p_i \end{cases}$$

for $0 \leq i < m, 0 \leq j < l$.

Proof. It is known that $\forall i, j (0 \leq i < m, 0 \leq j < l) C_i^j = A_{j,k} (k = p_i)$ and

$$\left\{ \begin{array}{l} C_i^j = A_{j,k} \\ k = p_i \\ A_{j,k} = \begin{cases} 1 \text{ if } t_j = k \\ 0 \text{ if } t_j \neq k \end{cases} \end{array} \right\} \Rightarrow C_i^j = \begin{cases} 1 \text{ if } t_j = p_i \\ 0 \text{ if } t_j \neq p_i \end{cases}$$

for $0 \leq i < m, 0 \leq j < l$.

The other column vector is defined in the Enable function(Definition 1). The algorithm transform the above two column vectors to a column vector R_i^j defined as Result function(Definition 2).

Theorem 2 (Resulting). $R_i^j = E_i^j \wedge C_i^j$ for $0 \leq i < m, 0 \leq j < l$, " \wedge " is bit and operating.

Proof. $\forall i, j(0 \leq i < m, 0 \leq j < l)$

(a) **Suppose** $E_i^j = 0 \Rightarrow E_i^j \wedge C_i^j = 0 \xrightarrow{\text{Definition 1}} t_{j-i}t_{j-i+1} \dots t_{j-1} \neq p_0p_1 \dots p_{i-1} \Rightarrow t_{j-i}t_{j-i+1} \dots t_j \neq p_0p_1 \dots p_i \xrightarrow{\text{Definition 2}} R_i^j = 0 \Rightarrow R_i^j = E_i^j \wedge C_i^j = 0$.

(b) **Suppose** $E_i^j = 1$ and $C_i^j = 0 \Rightarrow$

$$\left\{ \begin{array}{l} E_i^j = 1 \xrightarrow{\text{Definition 1}} \left\{ \begin{array}{l} E_i^j \wedge C_i^j = C_i^j \\ t_{j-i} \dots t_{j-1} = p_0 \dots p_{i-1} \end{array} \right\} \\ C_i^j = 0 \xrightarrow{\text{Theorem 1}} t_j \neq p_i \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} E_i^j \wedge C_i^j = C_i^j = 0 \\ t_{j-i}t_{j-i+1} \dots t_j \neq p_0p_1 \dots p_i \end{array} \right\} \xrightarrow{\text{Definition 2}} \left\{ \begin{array}{l} E_i^j \wedge C_i^j = 0 \\ R_i^j = 0 \end{array} \right\} \xrightarrow{\text{Definition 2}} R_i^j = E_i^j \wedge C_i^j$$

(c) **Suppose** $E_i^j = 1$ and $C_i^j = 1$

$$\left\{ \begin{array}{l} E_i^j = 1 \xrightarrow{\text{Definition 1}} \left\{ \begin{array}{l} E_i^j \wedge C_i^j = C_i^j \\ t_{j-i} \dots t_{j-1} = p_0 \dots p_{i-1} \end{array} \right\} \\ C_i^j = 1 \xrightarrow{\text{Theorem 1}} t_j = p_i \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} E_i^j \wedge C_i^j = C_i^j = 1 \\ t_{j-i}t_{j-i+1} \dots t_j = p_0p_1 \dots p_i \end{array} \right\} \xrightarrow{\text{Definition 2}} \left\{ \begin{array}{l} E_i^j \wedge C_i^j = 1 \\ R_i^j = 1 \end{array} \right\} \xrightarrow{\text{Definition 2}} R_i^j = E_i^j \wedge C_i^j$$

therefore conclude $R_i^j = E_i^j \wedge C_i^j$

Vector Transforming. The module transform a column vector to a column vector. The column R_i^j is coming from the result of Symbol comparing.

Definition 3 (Result-Enable function)

$$RE_i^j = \begin{cases} R_i^{j-1} & \text{if } 0 < j < l \\ R_i^{l-1} & \text{if } j = 0 \end{cases}$$

RE_i^0 is one clock cycle later then RE_i^k ,
for $0 \leq i < m, 0 \leq j < l, 0 < k < l$.

In order to achieve the effect of which RE_i^0 is one clock cycle later then RE_i^k , RE_i^0 is generated in FPGA by which R_i^{l-1} pass a register.

2.2 The Vector-And Algorithm

Regular expressions can be divided into **string** and **Kleene star**. Two part of the implementation was elaborated as the following.

String Matching. A **string** over an alphabet Σ is a finite sequence of symbols from the alphabet [2].

Now, we introduce the construct of the algorithm. In this case, the pattern string is expressed as $P_m = p_0p_i \dots p_{m-1}$. The text is expressed as $T_n = t_0t_1 \dots t_{n-1}$, where $n \gg m$. The entered text per one clock cycle is expressed as $T_l = t_0t_1 \dots t_{l-1} (l \leq n)$. For entered T_l at any clock cycle, Decoder Matrix A can be generated at the same clock cycle. First, let $E_0 = [1, 1, \dots, 1]^T$, $C_0 = [A_{j,k}] (0 \leq i < m, 0 \leq j < l, k = p_0)$, E_0 and C_0 though module 2 generate R_0 , that is, $R_0 = E_0 \wedge C_0$. R_0 though module 3 generate RE_0 , that is, $RE_0^j = R_0^{j-1}$. Second, for each subsequent $m - 1$ characters are constructed using the same method. $E_i = RE_{i-1}$, $C_i = [A_{j,k}]$, $R_i = E_i \wedge C_i$, $RE_i^j = R_i^{j-1} (0 < i < m, 0 \leq j < l, k = p_i)$. Finally, the result of algorithm is R_{m-1} .

Proof: Mathematical induction is able to prove that the Vector-And algorithm can meet the any length string matching by induction on the length of the pattern. the length of the pattern is expressed as $|P_m| = m$.

(a) *Basis Step.* Let the length of the Pattern is 1, that is, $m = 1$. Then $P_m = p_0$. For entered T_l any clock cycle, Decoder Matrix A can be generated at the same clock cycle. According to the Theorem 1, we know that

$$C_0^j = \begin{cases} 1 & \text{if } t_j = p_0 \\ 0 & \text{if } t_j \neq p_0 \end{cases}$$

for $0 \leq j < l$.

on the other hand, $E_0 = [1, 1, \dots, 1]^T$. By the Theorem 2:

$$R_0^j = E_0^j \wedge C_0^j = C_0^j = \begin{cases} 1 & \text{if } t_j = p_0 \\ 0 & \text{if } t_j \neq p_0 \end{cases}$$

for $0 \leq j < l$.

so, the elements of the result (R_0) of the algorithm are exactly the matching results. if $R_0^j = 1$, then $t_j \in T_l, t_j = p_0$, else, $t_j \in T_l, t_j \neq p_0$.

- (b) *Induction Hypothesis.* Let $m > 1$, and suppose that the Vector-And algorithm can meet the any length string matching provided that $|P| \leq m$. That is, the elements of the result (R_{m-1}) of the algorithm are exactly the matching results, if $R_{m-1}^j = 1$, then $t_j \in T_l, t_{j-m+1}t_{j-m+2} \dots t_j = p_0p_1 \dots p_{m-1}$, else, $t_j \in T_l, t_{j-m+1}t_{j-m+2} \dots t_j \neq p_0p_1 \dots p_{m-1}$.
- (c) *Induction Step.* Let $|P| = m + 1$. Then $P_{m+1} = p_0p_i \dots p_m$. By the induction hypothesis, the elements of the result (R_{m-1}) of the algorithm are exactly the matching results. that is:

$$R_{m-1}^j = \begin{cases} 1 & \text{if } t_{j-m+1}t_{j-m+2} \dots t_j = p_0p_1 \dots p_{m-1} \\ 0 & \text{if } t_{j-m+1}t_{j-m+2} \dots t_j \neq p_0p_1 \dots p_{m-1} \end{cases}$$

for $0 \leq j < l$.

According to construction methods and Definition 3. Therefore,

$$E_m^j = RE_{m-1}^j = R_{m-1}^{j-1} = \begin{cases} 1 & \text{if } t_{j-m}t_{j-m+1} \dots t_{j-1} = p_0p_1 \dots p_{m-1} \\ 0 & \text{if } t_{j-m}t_{j-m+1} \dots t_{j-1} \neq p_0p_1 \dots p_{m-1} \end{cases}$$

for $0 \leq j < l$.

By the Theorem 1. The matching between the current clock cycle entered text T_l and the last symbol of the pattern string p_m is expressed by:

$$C_m^j = \begin{cases} 1 & \text{if } t_j = p_m \\ 0 & \text{if } t_j \neq p_m \end{cases}$$

for $0 \leq j < l$.

According to the above two equations and Theorem 2. Therefore,

$$R_m^j = E_m^j \wedge C_m^j = \begin{cases} 1 & \text{if } t_{j-m}t_{j-m+1} \dots t_j = p_0p_1 \dots p_m \\ 0 & \text{if } t_{j-m}t_{j-m+1} \dots t_j \neq p_0p_1 \dots p_m \end{cases}$$

for $0 \leq j < l$.

so, the elements of the result (R_m) of the algorithm are exactly the matching results. if $R_m^j = 1$, then $t_j \in T_l, t_{j-m}t_{j-m+1} \dots t_j = p_0p_1 \dots p_{m-1}$, else, $t_j \in T_l, t_{j-m}t_{j-m+1} \dots t_j \neq p_0p_1 \dots p_{m-1}$.

In summary. The Vector-And algorithm proposed by the paper can meet the any length string matching, and the result is exactly R_{m-1} .

Regular Expression Matching. The **regular expressions** over an alphabet Σ^* are all strings over the alphabet $\Sigma \cup \{ (,), \emptyset, \cup, * \}$ [2]. In this paper, We only introduce the construction and proof of kleene star.

Now, we introduce the construction of the algorithm. In this case, This paper describes only the pattern is similar to the $P_m = p_0p_i \dots p_{m-1}(p_{m1}p_{m1+1} \dots p_{m2-1})^*p_{m2}p_{m2+1} \dots p_{m-1}$. Other types of kleene star is similar to the above type. The text is $T_n = t_0t_1 \dots t_{n-1}$, where $n \gg m$. The entered text per one clock cycle is $T_l = t_0t_1 \dots t_{l-1} (l \leq n)$. First P_m is divided into three parts,

$p_0p_i \dots p_{m_1-1}$, $p_{m_1}p_{m_1+1} \dots p_{m_2-1}$, and $p_{m_2}p_{m_2+1} \dots p_{m-1}$. Second, three parts are constructed by the above method as String matching. Finally, let $E_{m_1} = RE_{m_1-1} \mid RE_{m_2-1}$, $E_{m_2} = E_{m_1}$, and the result of algorithm is R_{m-1} .

Proof: We have proved the Vector-And algorithm proposed by the paper can meet the any length string matching. On this basis, to complete the rest of the proof.

- (a) *Basis.* " $p_0p_i \dots p_{m_1-1}$ " is a string. The previous proof has been proven on " $p_0p_i \dots p_{m_1-1}$ " can be matched by the algorithm, and the result is exactly R_{m_1-1} . Similarly, " $p_{m_1}p_{m_1+1} \dots p_{m_2-1}$ ", and " $p_{m_2}p_{m_2+1} \dots p_{m-1}$ " can also be considered separately as a string. So, they can be matched by the algorithm.
- (b) *Kleene star.* In order to prove kleene star, need to use the following definitions.

Definition 4 (String Power). [2] For each string w and each natural number i , the string w^i is defined as

$$w^0 = \phi, \text{ the empty string}$$

$$w^{i+1} = w^i \circ w, \text{ for each } i \geq 0$$

According to construction methods, $E_{m_1} = RE_{m_1-1} \mid RE_{m_2-1}$, in which $RE_{m_1-1}^j = R_{m_1-1}^{j-1}$, and $RE_{m_2-1}^j = R_{m_2-1}^{j-1}$. The following equation can be deduced:

$$E_{m_1}^j = \begin{cases} 1 & \text{if } t_{j-m_1} \dots t_{j-1} = p_0p_1 \dots p_{m_1-1} \\ & \text{or } \exists i, t_{j-m_1-i \cdot m_2} \dots t_{j-i \cdot m_2-1} (t_{j-m_2} \dots t_{j-1})^i \\ & = p_0p_1 \dots p_{m_1-1} (p_{m_1}p_{m_1+1} \dots p_{m_2-1})^i \\ 0 & \text{else} \end{cases}$$

for $0 \leq j < l, i \in \mathbb{N}$ (Natural numbers set).

proof: Though the previous subsection proof, we know the result of " $p_0p_i \dots p_{m_1-1}$ " is

$$R_{m_1-1}^j = \begin{cases} 1 & \text{if } t_{j-m_1+1} \dots t_j = p_0p_1 \dots p_{m_1-1} \\ 0 & \text{if } t_{j-m_1+1} \dots t_j \neq p_0p_1 \dots p_{m_1-1} \end{cases}$$

for $0 \leq j < l$.

According to construction methods and Definition 3. Therefore,

$$RE_{m_1-1}^j = R_{m_1-1}^{j-1} = \begin{cases} 1 & \text{if } t_{j-m_1}t_{j-m_1+1} \dots t_{j-1} = p_0p_1 \dots p_{m_1-1} \\ 0 & \text{if } t_{j-m_1}t_{j-m_1+1} \dots t_{j-1} \neq p_0p_1 \dots p_{m_1-1} \end{cases}$$

for $0 \leq j < l$.

Similar to the previous subsection, mathematical induction is able to prove that the following equation are true for any $i \in \mathbb{N}$ by induction on the i . the

result of $p_0p_i \dots p_{m-1}(p_{m1}p_{m1+1} \dots p_{m2-1})^*$ is:

$$\begin{aligned}
 & E_{m1} = RE_{m1-1} \mid RE_{m2-1} \Rightarrow \\
 R_{m2-1}^j = & \begin{cases} 1 & \text{if } \exists i, t_{j-m1-i*m2} \dots t_{j-i*m2-1} (t_{j-m2} \dots t_{j-1})^i \\ & = p_0p_1 \dots p_{m1-1}(p_{m1}p_{m1+1} \dots p_{m2-1})^i \\ 0 & \text{else} \end{cases} \\
 & \text{for } 0 \leq j < l, i \in \mathbb{N}.
 \end{aligned}$$

Therefore, conclude the equation needed to be prove is true.

- (c) *Conclusion.* According to construction methods, $E_{m2} = E_{m1}$, and the previous subsection proof, we can deduce the following equation is the result of $p_0p_i \dots p_{m1-1}(p_{m1}p_{m1+1} \dots p_{m2-1})^*p_{m2} \dots p_{m-1}$:

$$\begin{aligned}
 R_{m-1}^j = & \begin{cases} 1 & \text{if } \exists i, t_{j-m-m1-(i-1)*m2+1} \dots t_{j-m-(i-1)*m2} \\ & (t_{j-m+m1+1} \dots t_{j-m+m2})^i \\ & t_{j-m+m2+1} \dots t_j \\ & = p_0p_1 \dots p_{m1-1}(p_{m1}p_{m1+1} \dots p_{m2-1})^i \\ & p_{m2}p_{m2+1} \dots p_{m-1} \\ 0 & \text{else} \end{cases} \\
 & \text{for } 0 \leq j < l, i \in \mathbb{N} \cup \{0\}.
 \end{aligned}$$

so, the elements of the result (R_{m-1}) of the algorithm are exactly the matching results. if $R_{m-1}^j = 1$, then $t_j \in T_l, \exists i \in \mathbb{N} \cup \{0\}, t_{j-m+1}t_{j-m+2} \dots t_j = p_0p_i \dots p_{m1-1}(p_{m1}p_{m1+1} \dots p_{m2-1})^ip_{m2} \dots p_{m-1} = P_m$, else, $t_j \in T_l, \forall i \in \mathbb{N} \cup \{0\}, t_{j-m+1}t_{j-m+2} \dots t_j \neq p_0p_i \dots p_{m1-1}(p_{m1}p_{m1+1} \dots p_{m2-1})^ip_{m2} \dots p_{m-1} = P_m$.

In summary. The Vector-And algorithm proposed by the paper can meet the any regular expression matching, and the result is exactly R_{m-1} .

3 Experimental Results

This section presents experimental results of the Vector-and algorithm proposed. In our experiments, the throughput and the averages logical cells (LEs) per character are used as the performance measures of the Vector-and algorithm proposed, which contains a subset of SNORT [1] rules. The size of the subset is specified by the total number of symbols of the rules in the subset. the throughput indicates the maximum number of bits per second which the circuit can process.

The **LE** is the basic building block for the FPGA devices. One LE contains a 4-input look-up table(LUT), a 1-bit register, and additional carry and cascade logic, so the average number of LEs per symbol can be the criteria for evaluating the area cost.

Table 1 list some kinds of the numbers of symbols per one clock cycle, the corresponding throughput, and the LEs per symbol. Due to the tools of FPGA and the time, we did not do more experiments to prove that our approach can handle more symbols per one clock cycle. It can be observed from Table 1:

- the number of symbols per one clock cycle may be any natural number.
- the frequency of the proposed approach is the same as the highest frequency of FPGA.
- the number of LEs per symbol roughly linearly increase with the number of symbols per one clock cycle. This is because we construct all circuit to match text with similar construction.

At last, it is easy to conduct that the latency of the vector-and algorithm proposed by the paper is lower than one clock cycle. In the existing FPGA, the frequency of FPGA is no more than 500MHz (megahertz). the latency of our approach is lower than 2ns.

4 Conclusion

This paper presents a novel NFA-based regular expression algorithm which is a powerful technique for high throughput and low latency regular expression matching.

With the regular expressions matching demand of high-speed networks, this paper proposed a parallel regular expression matching algorithm. The experiment proved that our approach achieved the theoretical throughput to achieve wire-speed processing regular expression matching purposes which is clock frequency multiplied by the width of entered symbols per one clock cycle, and a Stratix II series EP2S180 has ten thousand regular expression-level processing power.

References

1. SNORT Network Intrusion Detection System, www.snort.org
2. Lewis, H.R., Papadimitriou, C.H.: Elements of the theory of computation, 2nd edn. Prentice-Hall, Inc.
3. Floyd, R.W., Ullman, J.D.: The Compilation of Regular Expressions into Integrated Circuits. *Journal of ACM* 29(3), 603–622 (1982)
4. Sidhu, R., Prasanna, V.K.: Fast Regular Expression Matching Using FPGAs. In: Field-Programmable Custom Computing Machines (FCCM 2001), pp. 227–238 (2001)
5. Hutchings, B.L., Franklin, R., Carver, D.: Assisting Network Intrusion Detection with Reconfigurable Hardware. In: 10th Annual IEEE Symposium on Field-Programmable Custom Computing Machines (FCCM 2002), p. 111 (2002)
6. Clark, C.R., Schimmel, D.E.: Scalable Pattern Matching for High Speed Networks. In: Field-Programmable Custom Computing Machines (FCCM 2004), pp. 249–257 (2004)
7. Brodie, B.C., Taylor, D.E., Cytron, R.K.: A Scalable Architecture For High-Throughput Regular-Expression Pattern Matching. *SIGARCH Comput. Archit. News* 34(2), 191–202 (2006)
8. Lin, C.H., Huang, C.T., Jiang, C.P., Chang, S.C.: Optimization of Pattern Matching Circuits for Regular Expression on FPGA. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 15(12), 1303–1310 (2007)

Interference Mitigation Based on Multiple SINR Thresholds in 60GHz Wireless Networks

Weixia Zou, Fang Zhang, Guanglong Du, and Bin Li

Wireless Network Lab, Beijing University of Posts and Telecommunications
Key Lab of Universal Wireless Communications, MOE Beijing, China
byrzhangfang@126.com

Abstract. To investigate spatial interference statistics for multi-Gigabit indoor networks operating in the unlicensed 60 GHz “millimeter (mm) wave” band, power control was studied with a non-cooperative game-theoretic approach. Considering different data rates corresponding to different signal-to-interference plus noise ratio (SINR) thresholds in the standard of 60 GHz wireless networks, the Piecewise Utility-based Transmit Power Control (PUTPC) was proposed to improve the system performance. The proposed method is related to the SINR threshold and the obtained SINR. A detailed analysis of the existence and uniqueness of Nash equilibrium for the above non-cooperative game is presented. Numerical results demonstrate that the new algorithm can efficiently improve the system capacity and can lower average outage percentage effectively with saved power.

Keywords: 60GHz mm wave, power control, piecewise utility function.

1 Introduction

The amount of relatively unused spectrum available in the unlicensed 60 GHz “mm wave” band has generated recent interest to develop standards to form numbers of standardization bodies in this band[1-7]. In most standards, the main application scenarios are oriented toward transferring a high-quality video stream over a wireless channel in home and office environments. It is envisioned that the network topology will be a greatly dense one. As everyone knows, there are only three non-overlapping communication channels available in most countries over the 7 GHz wide spectra. Unfortunately, it creates problems that the interference between neighboring networks becomes harder to manage. Therefore, it is necessary to find interference mitigation (IM) techniques for 60 GHz wireless networks.

In recent years, the game-theoretic with low complexity is widely used in power control to suppressing interference. In [8] and [9], power control is modeled as a non-cooperative game in which users change their transmit powers to maximize their net utilities. Based on the traditional Sigmoid function, the utility function is pretty flexible on the demand of SINR threshold value. In [10], it provides a costing function which can gain the optimum global solution, but the operation is quite

complicated. Reference [11] investigates the mutual interference of coexistence scenarios which is comprised of two different physical layer Monte Carlo simulations.

In short, the contributions of this paper are: This paper borrows basic concept of the utility function based on the Frame Success Rate (FSR) defined in [12]. Due to the characteristics of 60 GHz communication systems having multiple SINR thresholds, this paper smartly uses the idea of combining the utility function with piecewise function, which improves the performance of traditional algorithm.

2 System Model

This paper considers a typical dense office environment including 9 cubicles adjacent to each other as shown in Fig. 1 (a). Each cube has a single link defined by two randomly placed devices within the cube as shown in Fig. 1 (b) (Red denote transmitters, Blue denote receivers). Furthermore, we assume that all 9 cubes operate on the same channel for the worst case scenario.

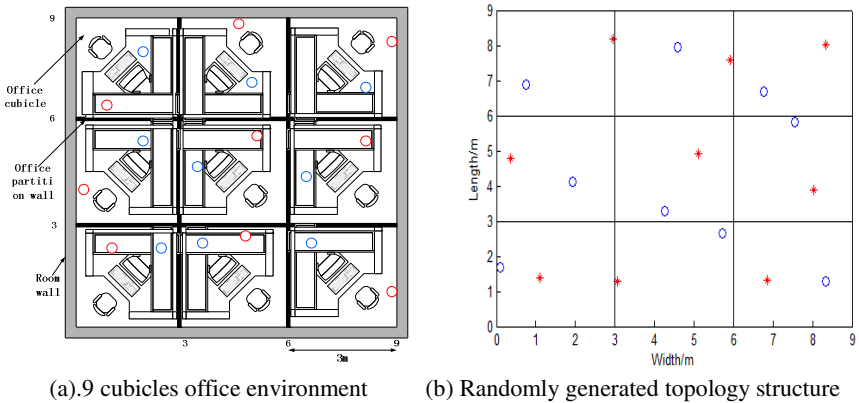


Fig. 1. The simulation network

Instead of using the ideal isotropic radiation pattern as in [13], each antenna element is modeled as a realistic patch antenna which has a cardioids radiation pattern. The patch antenna gain $\rho(\theta)$ is defined in the way similar to [14]. In order to lower the complication, assume all devices in the networks use smart antenna arrays with codebook of IEEE 802.15.3 [5]. The path loss between a transmitter and a receiver adopt the definition of [14].

3 PUTPC Definition

Borrowing the concepts from [12] and [13], we will reformulate the definition of utility function according to the feature of 60GHz communication (see Tab.1).

Table 1. MCS and SINR threshold

MCS mode	MCS1	MCS2	MCS3
Modulation	QPSK	QPSK	16QAM
Code rate	1/2	2/3	2/3
SINR threshold	5.5dB	13dB	18dB
Data rate	0.952Gbps	1.904Gbps	3.807Gbps

In line with Table 1, there are different SINR threshold values for different data rates. Incorporating the idea of piecewise function into the utility function defined in [13]. We can formulate the goodput function as

$$V_i(p_i, p_{-i}) = R_i f(\gamma_i) \tag{1}$$

Where R_i is the data rate shown in table 1, γ_i denotes the obtained SINR and $f(\gamma_i)$ denotes the frame success rate ranging from [0,1] which can be expressed as

$$f(\gamma_i) = \frac{1 - e^{(-\gamma_i)}}{1 + e^{(\tau_i - \gamma_i)}}, \tau_i \text{ denotes the SINR threshold values as shown in table 1.}$$

If a link gets the SINR lower than 5.5dB, its goodput utility function will borrow the goodput utility function corresponding to the SINR threshold equaling to 5.5dB to reduce the outage probability of the link. Therefore, the specific forms of goodput utility function can be reformulated as

$$V_i(p_i, p_{-i}) = \begin{cases} \frac{0.952e^9 (1 - \exp(-\gamma_i))}{1 + \exp(10^{0.55} - \gamma_i)} & \gamma_i < 13dB \\ \frac{1.904e^9 (1 - \exp(-\gamma_i))}{1 + \exp(10^{1.3} - \gamma_i)} & 13dB \leq \gamma_i < 18dB \\ \frac{3.807e^9 (1 - \exp(-\gamma_i))}{1 + \exp(10^{1.8} - \gamma_i)} & \gamma_i \geq 18dB \end{cases} \tag{2}$$

Where,

$$\gamma_i = \frac{p_i h_{ii}}{\sum_{j \neq i} p_j h_{ij} + \sigma^2} \tag{3}$$

The specific cost function can be considered as the charge for the high difference between the SINR the link got and its corresponding SINR threshold value. It can be expressed as $C(p_i, p_{-i}) = \lambda_i (\gamma_i - \tau_i)^2$, λ_i is a constant pricing coefficient associating with the channel gain. τ_i is the SINR threshold values as shown in table 1. If a link gets the SINR lower than 5.5dB, τ_i equals zero. Hence, the specific forms of cost function can be reformulated as

$$C(p_i, p_{-i}) = \begin{cases} \lambda_i (\gamma_i)^2 & \gamma_i < 5.5dB \\ \lambda_i (\gamma_i - 10^{0.55})^2 & 5.5dB \leq \gamma_i < 13dB \\ \lambda_i (\gamma_i - 10^{1.3})^2 & 13dB \leq \gamma_i < 18dB \\ \lambda_i (\gamma_i - 10^{1.8})^2 & 18dB \leq \gamma_i \end{cases} \quad (4)$$

Therefore, the utility function can be expressed as

$$U_i(p_i, p_{-i}) = V_i(p_i, p_{-i}) - C(p_i) \quad (5)$$

Thus, the algorithm of interference mitigation based on multiple SINR threshold can be performed as each coexisting links maximize its own net utility according to (5) by tuning the transmit power in response to other’s action, i.e.

$$\begin{aligned} & \max_{p_i \geq 0} U_i(p_i, p_{-i}) \end{aligned} \quad (6)$$

4 Feasibility and Flowchart of PUTPC

4.1 Nash Equilibrium

According to (5), each link uses lower transmitting power to maximize its utility. Specifically, the maximum of $U_i(p_i, p_{-i})$ can be obtained following the first order condition of (4) that

$$\frac{\partial U_i}{\partial p_i} = (R_i f'(\gamma_i) - 2\lambda_i (\gamma_i - \tau_i)) \times \frac{\partial \gamma_i}{\partial p_i} \quad (7)$$

Then let $\frac{\partial U_i}{\partial p_i} = 0$ to derive the best transmitting power p_i^0 .

Thus, the net utility $U_i(p_i, p_{-i})$ achieved at γ_i^0 denoted by $U^0(p_i, p_{-i})$

$$U^0(p_i, p_{-i}) = V(\gamma_i^0) - V'(\gamma_i^0)(\gamma_i^0 - \tau_i) / 2. \quad (8)$$

Here let the $U^0(p_i, p_{-i}) = 0$ achieved at $\underline{\gamma_i^0}$ defined as the turn-off point, which is inherently determined by the goodput function $V_i(p_i, p_{-i})$. It should be noted that if $U^0(p_i, p_{-i}) > 0$ at any P then the $\underline{\gamma_i^0}$ equals the start value of the piecewise function according to (5).

Under the constraint of $p_i^0 \leq p_{\max}$, the value of λ_i is

$$\lambda_i \geq \frac{V'(\gamma_i^0)}{2\left(p_{\max} \frac{h_{ii}}{Q(P_{-i})} - \tau_i'\right)} \tag{9}$$

Where $Q(P_{-i}) = \sum_{j \neq i} p_j h_{ij} + \sigma^2$ Then, we can gain the partial derived function of (6) related to p_j is

$$\frac{\partial^2 U_i}{\partial^2 p_i p_j} = \frac{1}{p_i} \times \frac{\partial \gamma_i}{\partial \gamma_i} \times [R_i [f''(\gamma_i) \gamma_i + f'(\gamma_i)] - \lambda_i \gamma_i + 2\lambda_i \tau_i'] \tag{10}$$

Where, $f''(\gamma_i) = \frac{(1 + e^{\tau_i'}) (e^{\tau_i - \gamma} - 1)}{e^{\gamma} (e^{\tau_i - \gamma} + 1)^3} \leq 0$, then

$$R_i (f''(\gamma_i) \gamma_i + f'(\gamma_i)) - 4\lambda_i \gamma_i + 2\lambda_i \tau_i' \leq 0$$

As well, $\frac{\partial \lambda_i}{\partial p_j} = (-1) \frac{p_i h_{i,i} h_{j,i}}{Q(P_{-i})^2} \leq 0$, according to (10)

$$\frac{\partial^2 U_i}{\partial^2 p_i p_j} \geq 0, \quad \forall j \neq i. \tag{11}$$

According to (11), the (5) meet the Nash equilibrium of Super-modular Games. Hence, there is the only one p_i^0 can meet (6).

5 Performance Analysis

The experimental network is a typical dense office environment specified in section 2. Each cubicle has a dimension of $3m \times 3m$ (width \times length). The initial transmitting power is 10dBm. In this paper, all the devices are equipped with same 2-D antenna array having the same number of antennas, i.e. $M_t = M_r = 4 \times 4$. Generally, there are several metrics used to evaluate the performance of the power control such as network data rate and outage probability. Note that referenced algorithm 1 refers to [13], referenced algorithm 2 refers to [12]. We run the simulation 1000 times with random network topologies.

5.1 The SINR Distribution

At first, we will study the distribution of SINR shown in Fig. 2. After adopting the proposed PUTPC algorithm, the obtained SINRs are distributed mainly on the right

side of the SINR thresholds, particularly most are greater than 18dB threshold, which are corresponding to (6) whose max utility are obtained at the right side of the SINR thresholds. Reference [12] has the Poisson distribution of SINR due to the unique utility function. Although, the SINR distribution of reference [13] is similar to PUTPC, its most SINRs are distributed on the right of minimum SINR

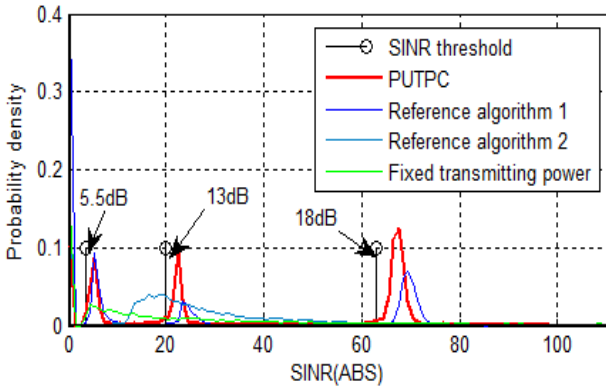


Fig. 2. The distribution of SINR

5.2 Network Data Rate and Outage Probability

The Complementary Cumulative Distribution Function (CCDF) of the aggregated data rate and number of interrupted links are shown in Fig. 3. Simulation result shows that referenced algorithms are obviously not suitable to 60 GHz mm wave communication systems having multiple SINR threshold values. In [12], it reduces the asymmetric performance degradation due to the mutual interference between two type of Ultra-wideband (UWB) systems. In [13], the goodput function is based on sigmoid function, which is also a piecewise function, but it built a unified pricing function varying linearly with the channel gain. There is an obvious performance improvement for the proposed schemes.

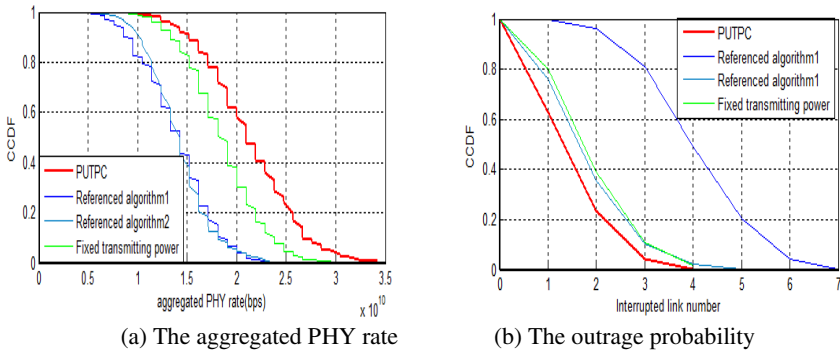


Fig. 3. Performance comparison of aggregated PHY rate and outage probability

6 Conclusion

In this paper, we have developed a simulation model for 60 GHz mm wave systems with a realistic 2-D antennal model and propagation model and then proposed a piecewise utility-based transmit power control algorithm. The novelty in the proposed PUTPC lies in the pricing function which is relative to the SINR threshold and the obtained SINR. In addition, the different obtained SINRs under the max transmitting power result in the differentiated pricing coefficients. The feasibility of PUTPC is analyzed and the Nash equilibrium is proved under the feasible condition. The simulation results reveal that PUTPC can always outperform the conventional utility-based TPC in multiple SINR threshold systems, but it still not well-developed. How to improve the ratio of energy to rate are all topics requiring further research.

Acknowledgments. This work was supported by NSFC (61171104), the Fundamental Research Funds for the Central Universities (G470270,G470415), Important National Science & Technology Specific Projects(2009ZX03006-006/-009),the BUPT excellent Ph.D. students foundations (CX201122)

References

1. Yong, S.-K., Xia, P., Valdes-Garcia, A.: 60GHz technology for Gbps WLAN and WPAN, from theory to practice, pp. 21–35. A John Wiley and Sons, Ltd. (2011)
2. Guo, N., Qiu, R.C., Mo, S.S., Takahashi, K.: 60-GHz millimeter wave radio: Principle, technology and new results. EURASIP Journal on Wireless Communications and Networking 2007(1) (January 2007)
3. Geng, S., Kivinen, J., Zhao, X., Vainikainen, P.: Millimeter-wave propagation channel characterization for short-range wireless communication. IEEE Transactions on Vehicular Technology 58(1), 3–13 (2009)
4. Singh, H., Yong, S.-K., Oh, J., et al.: Principles of IEEE 802.15.3c: Multi-Gigabit Millimeter-Wave Wireless PAN. In: The 18th International Conference on Computer Communications and Networks, pp. 1–6 (August 2009)
5. IEEE 802.15, WPAN Millimeter Wave Alternative PHY Task Group 3c (TG3c), <http://www.ieee802.org/15/pub/TG3c.html>
6. Wireless HD 1.0, <http://www.wirelesshd.org/>
7. IEEE 802.11, Very High Throughput in 60 GHz TaskGroup ad (TGad), http://www.ieee802.org/11/Reports/tgad_update.html
8. Xia, W., Qi, Z.: Power control for Cognitive Radio Base on Game Theory. In: 2007 International Conference on Wireless Communications, Networking and Mobile Computing, pp. 1256–1259. IEEE, Shanghai (2007)
9. Uykan, Z., Koivo, H.N.: Sigmoid-Basis Nonlinear Power-Control Algorithm for Mobile Radio Systems. IEEE Transactions on Vehicular Technology 53(1), 265–271 (2004)
10. Aein, J.M.: Power balancing in systems employing frequency reuse. COMSAT Tech. Review 3(2), 277–299 (1973)
11. Zhang, Y., Wu, H., Zhang, Q., Zhang, P.: Interference mitigation for coexistence of heterogeneous ultra-wideband systems. EURASIP Journal on Wireless Communications and Networking 2006(2) (April 2006)

12. Li, J., Xue, F., Yang, C., Li, W., Shi, H.: Fast adaptive power control approach for cognitive radio network. *Journal of Xidian University* 37(2) (April 2010) (in Chinese)
13. Vilzmann, R., Bettstetter, C., Hartmann, C.: On the impact of beamforming on interference in wireless mesh networks. In: *Proc. IEEE Workshop on Wireless Mesh Networks*, Santa Clara, CA (September 2005)
14. Park, M., Gopalakrishnan, P.: Analysis on Spatial Reuse and Interference in 60-GHz Wireless Networks. *IEEE Journal on Selected Areas in Communications* 27(8) (October 2009)

Analysis of the Coupling Action in Nonlinear Harmonic Vibration Synchronization System

Xiaohao Li* and Zhenwei Zhang

School of Mechanical Engineering and Automation, Northeastern University,
Shenyang 110819, China

{xhli, zhenwzhang}@me.neu.edu.cn

Abstract. The nonlinear harmonic vibration system that synchronously excited by two eccentric rotors is a kind of typical synchronous vibration machine which possesses the typical coupled motion characteristics. Because of the nonlinear vibration motion of the vibrating body, which make the two eccentric rotors to reach the coupling synchronous rotary motion, and the coupling action will result in the change of the motion characteristics and movement patterns for the two eccentric rotors, eventually to achieve the coupling synchronous motion. Based on the nonlinear vibration theory, the paper has established the electromechanical coupling nonlinear dynamics equations of the nonlinear harmonic vibration system. By analyzing the coupling factor of the two eccentric rotors, the influence to the equilibrium state of the nonlinear vibration system because of the nonlinear coupling strength has been discussed, and the synchronous movement evolution procedure of the two eccentric rotors has been researched further. Based on the research result, the engineering method that to design the structural parameters of the nonlinear harmonic synchronization vibration machine which excited by two eccentric rotors has been deduced in the paper also.

Keywords: Nonlinear system, Harmonic vibration, Two-eccentric rotors, Coupling action.

1 Introduction

Vibration behavior of the nonlinear synchronous system not only depends on the movement of the vibration characteristics of the system, but also depends on the interaction between the vibrating system itself and the exciting sources [1-5]. The strength of coupling action between exciting sources could directly affect the vibration characteristics of the complex systems movement and state of motion [7, 10]. To research the coupling action characteristics of the exciting source of the nonlinear vibration system [9], and to analyze the influence of the dynamic behavior of the nonlinear vibration system because of the coupling characteristics is an important part of science issue [11, 15]. The vibration synchronization theory for eccentric rotors which driven by two exciting motors had been studied deeply by

* Corresponding author.

Wen Bangchun and other scholars [8]. In the literature [6, 16], the vibration synchronization problem of hydraulic motors which used to drive the eccentric rotors had been researched.

Based on the vibration synchronization dynamic model that shown as Figure 1, the interaction of every independently driven eccentric rotor's rotary movement had been researched in the paper, and based on the coupling effect, the evolutionary process of two separately rotors' rotation speed status automatically consistent synchronization motion state had been discussed also.

2 Coupling Motion Description of the Two Eccentric Rotors

According to the literature [1], the dynamic model of the synchronization vibration system was shown as Fig. 1.

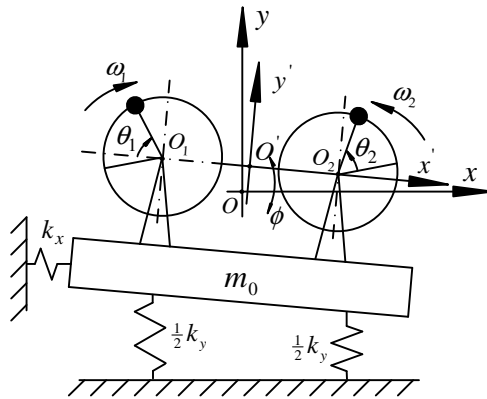


Fig. 1. Model of the nonlinear vibration synchronous system

To calculate the kinetic energy, potential energy, and energy dissipation of the synchronization vibration system, and based on the Lagrange equations, the eccentric rotor motion equations and the steady-state plasmid solution of the vibrating body can be described as

$$J_i \ddot{\theta}_i + c_i \dot{\theta}_i + m_i r_i (\ddot{y} \cos \theta_i - \ddot{x} \sin \theta_i) - m_i r_i l \varphi^2 \sin(\theta_i - \varphi) = T_i - m_i r_i g \cos \theta_i \quad (1)$$

where, J_i is the eccentric vibrating motor inertia. θ_i is the vibrating motor eccentric back corner. c_i is the vibrating motor for the rotary damping. m_i is the eccentric vibrating motor block quality. r_i is the exciting motor eccentric eccentricity. x , y and φ is the displacement in x , y direction and torsion angle in φ direction of the nonlinear vibration system. T_i is the vibrating motor output torque. Based on nonlinear vibration energy method [13], the torsion angle φ can be obtained as Eq.(2)

$$\varphi = -\frac{(m_1 r_1 \sin \theta_1 + m_2 r_2 \sin \theta_2)}{m_0 l^2 \left(\frac{\theta_1 + \theta_2}{2}\right)^2 - k_\varphi} \tag{2}$$

where, k_φ the spring stiffness coefficient for the vibration system in φ direction.

Set \bar{v} is the average speed for the two vibrating motors, then

$$\theta_i = \bar{v}t + \Delta\tau_i \tag{3}$$

To take the derivative about Eq.(1), and to calculate the average 2π interregional, then Eq.(1) can be expressed as

$$\frac{d^2\Delta\tau_i}{dt^2} + \frac{c_i}{J_i\bar{v}}\left(\frac{d\Delta\tau_i}{dt} + 1\right) + \frac{H \sin(\Delta\tau_1 - \Delta\tau_2)}{2J_i\bar{v}^2} = \frac{T_i}{J_i\bar{v}^2} \tag{4}$$

where,

$$H = \frac{m_1 m_2 r_1 r_2 \cdot \Delta k}{m_0^2} \tag{5}$$

Know from the Eq.(5), the equation a double excitation motor coupling term . Where H describes the coupling action between the excitation characteristics of the two eccentric rotors.

3 Coupling Motion Characteristics Analysis of Two Eccentric Rotors

In the actual engineering systems that correspond with Fig. 1, the structural properties of two eccentric rotor is very close to the corresponding difference between the structural parameters, the difference between the corresponding damping coefficient is small. Set $i=1$ and 2 to make a difference using the two equations, the angle phase difference $\Delta\tau = \Delta\tau_1 - \Delta\tau_2$ variable equations of the two exciting motors can

be gotten, so set $\vartheta = \frac{d\Delta\tau}{dt}$ then

$$\frac{d}{dt}\vartheta = \frac{c_1}{J_1\bar{v}}\vartheta + \frac{H}{J_1\bar{v}^2}\left(\frac{\Delta T}{H} - \sin \Delta\tau\right) \tag{6}$$

In Eq.(6), to set

$$\frac{\Delta T}{H} = \kappa \tag{7}$$

Based on Eq.(5), Eq.(7) to derivation of the deformation $\Delta\tau$, then

$$\frac{d\vartheta}{d\Delta\tau} = -\frac{c_1'}{J_1\bar{v}} + \frac{H}{J_1\bar{v}^2\vartheta}(\kappa - \sin \Delta\tau) \tag{8}$$

To carry out phase plane analysis of Eq.(8), if $\kappa = \sin \Delta\tau$ is true, then

$$\begin{cases} \Delta\tau = \arcsin \kappa \\ \text{or } \Delta\tau = \pi - \arcsin \kappa \end{cases} \tag{9}$$

At this point the equilibrium state solution $\Delta\tau = \arcsin \kappa$ of Eq.(8) has the stable.

According to Eq.(6), Eq.(7) and Eq.(9), range of values of κ can be obtained which can be used to reflect the coupling strength of two motors

$$0 \leq \kappa < 1 \tag{10}$$

When $\kappa = 1$, know the stability characteristic value of Eq.(9), $\Delta\tau = \frac{\pi}{2}$ is the critical steady state of Eq.(8).

4 Coupling Motion State Evolution of the Two Exciting Rotors

Know from Fig. 1, there is elastic support system, due to vibrating motion of the system, the two rotor rotary movement coupled with the coupling strength changes, the rotor rotation speed and phase angle also changes, resulting in system evolved different sports structure.

To lead the transformation Eq.(11)

$$\begin{cases} \frac{d\Delta\tau_1}{dt} = X_1 \\ \frac{d\Delta\tau_2}{dt} = X_2 \end{cases} \tag{11}$$

For the two- rotor excitation, when the realization of resonant synchronous rotation, the balance due to the phase plane

$$\frac{dX_1}{dt} = \frac{dX_2}{dt} = 0 \tag{12}$$

Based on two rotors nonlinear vibration system excitation conditions $m_1 = m_2$, $r_1 = r_2$ and $c_1' = c_2'$, to deduce Eq.(11), Eq.(12) and Eq.(8), then Eq.(13) can be obtained

$$\begin{cases} X_1 = \frac{1}{c_1 \dot{\omega}} \left[-2 \left(c_1 \dot{\omega} - \frac{T_1 + T_2}{2} \right) + H(\kappa - \sin \Delta \tau) \right] \\ X_2 = \frac{1}{c_1 \dot{\omega}} \left[-2 \left(c_1 \dot{\omega} - \frac{T_1 + T_2}{2} \right) - H(\kappa - \sin \Delta \tau) \right] \end{cases} \quad (13)$$

Know from Eq.(3) and Eq.(11), that Eq.(13) reflects the two rotor speed and excitation system vibration coupling strength between .

In Eq.(13), within the coupling strength value κ that if $\kappa - \sin \Delta \tau = 0$ then, know from Eq.(13), that $X_1 = X_2$. Namely know from Eq.(3), the two vibrating motors possess the same rotor speed $\dot{\theta}_1 = \dot{\theta}_2$, the two eccentric rotors can realize the vibration coupling of the vibration body such as speed synchronous movement. And based on Eq.(9) of the analysis process knowledge, then the vibration system has a stable equilibrium.

Know that, the value $(\kappa - \sin \Delta \tau)$ determines the ability to achieve double-excitation coupled synchronous motor sport.

Based on Eq.(5), Eq.(7) and Eq.(10) to get that

$$\frac{|T_1 - T_2|}{(m_1 r_1) \cdot (m_2 r_2)} < \frac{|k_x - k_y|}{m_0^2} \quad (14)$$

The double excitation of nonlinear vibration systems the kinetic parameters of rotor synchronization conditions can be obtained from Eq.(14). In electrical engineering can be used as double resonance excitation synchronous system architecture an important basis for the design parameters.

5 Engineering Example

To select the structural parameters as followed for the nonlinear synchronization vibration model which driven by two exciting motors that shown as Fig. 1.

$$m_0 = 2550\text{kg} , m_1 = m_2 = 85\text{kg} , r_1 = r_2 = 0.10\text{m} , c_1 = c_2 = 0.15\text{N} \cdot \text{m} \cdot \text{s}/\text{rad} , |k_x - k_y| = 304.5\text{kN}/\text{m} , |T_1 - T_2| = 2.87\text{N} \cdot \text{m} .$$

Lead the parameters into Eq.(14) then $\frac{|T_1 - T_2|}{m_1 r_1 m_2 r_2} = 0.0397$, $\frac{|k_x - k_y|}{m_0^2} = 0.0468$. To

compared with Eq.(14) conditions, that it can be satisfied perfectly. Take the initial condition $\Delta \tau = 5^\circ$ to carry out the experiment and the simulation experiment results can be found as Fig. 2.

The rotary speed of two eccentric rotors change from the non-synchronous state ($\dot{\theta}_1 \neq \dot{\theta}_2$) to the harmonic synchronous state ($\dot{\theta}_1 = \dot{\theta}_2$) process had been shown in Fig. 2 by the speed difference and phase difference phase plane trajectory change

simulation, which can be used to verify theoretical harmonic synchronization conditions for the two eccentric rotors which were used to synchronously excite the nonlinear harmonic vibration machine that is shown in Fig. 1. The practical application data showed that the theoretical analysis results obtained in the paper could be used to deal with the harmonic vibration synchronization problem of the two eccentric rotors perfectly.

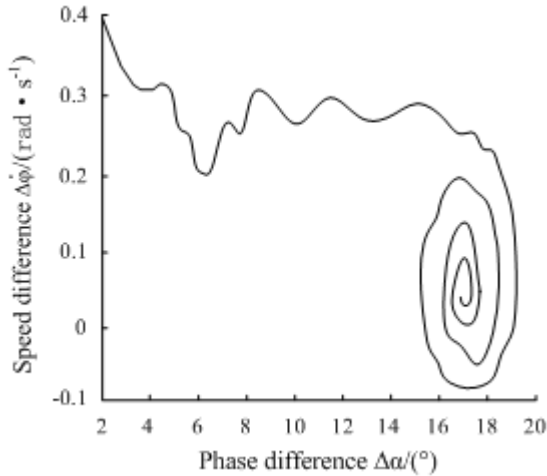


Fig. 2. Phase plane trajectory simulation of the synchronization process of the two eccentric rotors

6 Conclusion

Based on the nonlinear vibration motion of the vibrating machine the two eccentric rotors can realize the mutual coupling actions. By analyzing the coupling strength factor of the two eccentric rotors, the equilibrium stability of the nonlinear vibration machine can effectively be discussed, and it can be used to solve the coupling synchronization problem of the eccentric rotors rotary motion. The coupling factor parameters discussion of the nonlinear vibration system that carried out in the paper can be used to provide important reference frame for the structural parameters selection of the nonlinear vibration systems, which the two eccentric rotors can be used to realize the harmonic synchronization motion.

Acknowledgments. This work was supported by the Young Scientists Fund of the National Natural Science Foundation of China (Grant No. 51105066), and the Fundamental Research Funds for the Central Universities (Grant No. N110403004).

References

1. Wen, B.C., Zhao, C.Y., Su, D.H.: *Vibration Synchronization and Control Synchronization of the Machine System*. Science Press, Beijing (2003)
2. Blekhman, I.I., Fradkov, A.L., Tomchina, O.P.: Self-synchronization and Controlled Synchronization: General Definition and Example Design. *Math. Comput. Simul.* 58, 367–384 (2002)
3. Ke, L.L., Wang, Y.S., Yang, J., Kitipornchai, S.: Nonlinear free vibration of size-dependent functionally graded microbeams. *Int. J. Eng. Sci.* 50, 256–267 (2012)
4. Huang, X.L., Jia, X.L., Yang, J., Wu, Y.F.: Nonlinear vibration and dynamic response of three-dimensional braided composite plates. *Mech. Adv. Mater. Struct.* 15, 53–63 (2008)
5. Kanchan, R.S., Gopakumar, K., Kennel, R.: Synchronized Carrier-based SVPWM Signal Generation Scheme for the Entire Modulation Range Extending up to Six-step Mode Using the Sampled Amplitudes of Reference Phase Voltages. *IET Electr. Power Appl.* 1, 407–415 (2007)
6. Ebrahimi, F., Rastqoo, A., Bahrami, M.N.: Investigating the thermal environment effects on geometrically nonlinear vibration of smart functionally graded plates. *J. Mech. Sci. Technol.* 24, 775–791 (2010)
7. Ghayesh, M.H., Kazemirad, S., Reid, T.: Nonlinear vibrations and stability of parametrically excited systems with cubic nonlinearities and internal boundary conditions: A general solution procedure. *Appl. Math. Modell.* 36, 3299–3311 (2012)
8. Saranqi, S.K., Ray, M.C.: Active damping of geometrically nonlinear vibrations of laminated composite shallow shells using vertically/obliquely reinforced 1-3 piezoelectric composites. *Int. J. Mech. Mater. Des.* 7, 29–44 (2011)
9. Yacmini, R., Smith, K.S., Ran, L.: Monitoring Torsion Vibrations of Electromechanical Systems Using Stator Currents. *ASME J. Vib. Acoustics* 120, 72–79 (1998)
10. Zhang, G.C., Hu, D., Chen, L.Q., Yang, S.P.: Galerkin method for steady-state response of nonlinear forced vibration of axially moving beams at supercritical speeds. *J. Sound Vib.* 331, 1612–1623 (2012)
11. Shooshtari, A., Khadem, S.E.: A multiple scales method solution for the free and forced nonlinear transverse vibrations of rectangular plates. *Struct. Eng. Mech.* 24, 543–560 (2006)
12. Chen, Y.S., Cao, D.Q., Wu, Z.Q.: Recent Developments in Nonlinear Dynamics: Theory and Its Applications in Mechanical Systems. *J. Astronautics* 28, 794–804 (2007)
13. Chorfi, S.M., Houmat, A.: Nonlinear free vibration of a moderately thick doubly curved shallow shell of elliptical plan-form. *Int. J. Comput. Methods* 6, 615–632 (2009)
14. Alijani, F., Arnabili, M., Bakhtiari, N.F.: Thermal effects on nonlinear vibrations of functionally graded doubly curved shells using higher order shear deformation theory. *Compos. Struct.* 93, 2541–2553 (2011)
15. Younesian, D., Cao, D.Q., Wu, Z.Q.: Nonlinear vibration of a three-dimensional moving gantry crane subjected to a travelling trolley hoisting a swinging object. *Trans. Can. Soc. Mech. Eng.* 34, 333–350 (2010)
16. Guo, P.F., Lang, Z.Q., Peng, Z.K.: Analysis and design of the force and displacement transmissibility of nonlinear viscous damper based vibration isolation systems. *Nonlinear Dyn.* 67, 2671–2687 (2012)

Performance Analysis of Hierarchical Modulation with Higher Spectrum Efficiency for a Higher Date Rate T-DMB System

Linlin Dong, Lixin Sun, Xiaoming Jiang, and Na Zhu

School of Computer Science and Telecommunication Engineering Jiangsu University,
Jiangsu, China
donglinlin3390@163.com

Abstract. In order to improve current transmission capacity of the ancillary data in terrestrial digital multimedia broadcasting (T-DMB) system, the hierarchical modulation is adopted. The hierarchical modulation consists of Higher-priority (HP) data, which is the same as in the legacy T-DMB, and Low-priority (LP) data, which carries the additional data for the upgrade system, while the backward compatibility is guaranteed. Both QPSK modulation and 4ASK modulation are good candidates for LP data modulation. So in this paper we adopt coherent detection with comb-type pilot arrangement, the Least Squares (LS) estimate of pilot signals and the piecewise-linear interpolation channel estimation. Simulation results show that LP-QPSK provides a better choice than LP-4ASK for achieving a higher data rate in the Advance T-DMB (AT-DMB).

Keywords: AT-DMB, LP-4ASK, LP-QPSK, Hierarchical modulation.

1 Introduction

The digital multimedia services (video, audio, data, picture, etc) will become an increasing important aspect for all kinds of applications in our dairy life. Now T-DMB is a more important one. T-DMB is based on the Eureka-147 Digital Audio Broadcasting (DAB) structure [1], which is a novel audio broadcasting system intended to supersede the existing analog audio systems. Comparing to DAB, H.264/MPEG-4 Advance Video Coding (AVC) in source coding, additional Reed-Solomon code and byte interleaver in channel coding are introduced [2], [3].

So as to maximize the throughput in the T-DMB system, we adopt the hierarchical modulation. For the main goal of the hierarchical modulation scheme is to provide higher data rate through 1.536 MHZ system bandwidth. It is proposed to provide different classes of data to users with different reception conditions [4], [5]. So that the T-DMB with hierarchical modulation can transmit different bit rates according to the different modulating constellations, and error correction of different coding rates. Hierarchical modulation can be effectively used for satisfying the increase of the data rate and the assurance of backward compatibility with legacy receivers. It contains HP data and LP data. LP data can be modulated by different modulations. But we choose

QPSK and 4ASK as good candidates for they have better performances, respectively. LP-4ASK uses amplitude coefficients of the channel transfer function with the use of amplitude pilots for channel estimation [6]. The CDD estimation method may be used for coherent detection of LP-QPSK [7]. Only considering the performance of the system, we can choose the two modulation modes. Here we use the same channel estimation method to compare their performance. In this paper, we adopt coherent detection with comb-type pilot arrangement [8], the LS estimate of pilot signals and the linear interpolation channel estimation for LP data.

2 System Model

2.1 Hierarchical Modulation

Hierarchical modulation is to offer different levels of protections to transmit signals according to their significance. The transmitted stream in the source mode is classified into two subsets, HP data and LP data in terms of their significance. HP and LP are the original bits and the additional bit respectively. It is one of the techniques for multiplexing and modulating multiple data stream into one signal symbol stream, where the HP symbols and LP symbols are synchronously superimposed together before being transmitted. When hierarchical modulation is employed, users with good reception and advanced receiver can demodulate more than one data streams. For a user with poor reception, it may be able to demodulate data streams in original data [9].

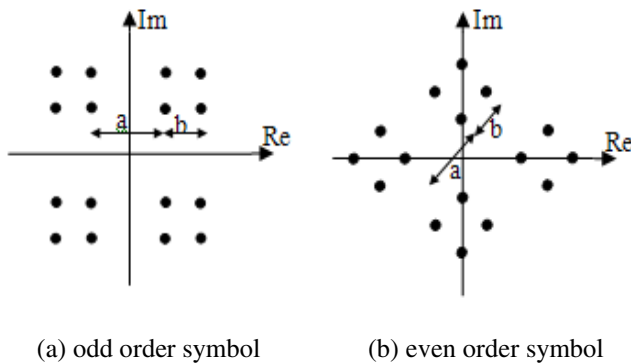


Fig. 1. (DQPSK/QPSK mode) Hierarchical constellation

The hierarchical modulation symbols are modulated by the OFDM modulator for transmission. Fig. 1 shows DQPSK /QPSK hierarchical constellation. Here λ_1 is a hierarchical parameter is given by $\lambda_1 = a/b$ Where a represents the minimum distance between two symbols in adjacent quadrances and b represents the distance between two neighboring symbol within one quadrant. Fig.2 shows DQPSK/4ASK hierarchical constellation. λ_2 is a hierarcheal parameter is given by $\lambda_2 = c/d$ Where

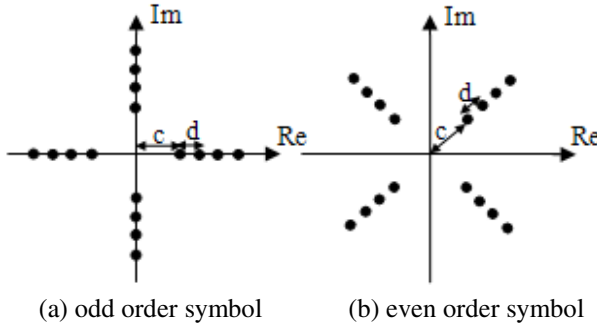


Fig. 2. (DQPSK/4ASK mode) Hierarchical constellation

c is the minimum distance between the original and the LP symbols, and d is the minimum distance between two adjacent LP symbols. In order to show the same degradation in legacy, LP-QPSK with $\lambda_1=3$ and LP-4ASK with $\lambda_2=2$ is applied for computer simulation.

2.2 HP Demodulation for AT-DMB with Hierarchical Modulation

At the transmitter of the AT-DMB system, HP data and LP data after channel encoded, are A-DPSK (LP-4ASK) modulation and QAM (LP-QPSK) modulated, respectively. Then the complex symbols are fed to an OFDM modulation, finally transmitted. In the receiver the corresponding inverse operation shaves to be carried out. The channel estimation and the frequency-domain equalizer are used for signal correction. Assuming that cycle prefix completely eliminates ISI and ICI, and then the received signal $Y_{n,k}$ in the frequency domain can be given by

$$Y_{n,k} = X_{n,k} \cdot H_{n,k} + N_{n,k} \tag{1}$$

Where $X_{n,k}$ is a transmitted frequency domain symbol, $H_{n,k}$ is a channel coefficient and $N_{n,k}$ is additive noise correspond to the k th subcarrier in the n th OFDM symbol. At the receiver, the received symbol $Y_{n,k}$ is itself used for HP demodulation. In order to produce soft-decision output, we adopt an optimal soft-decision method in HP-DQPSK demodulation corresponding to log-like-hood rate (LLR) as follows in [10].

$$\begin{aligned} LLR(b_o^H) &= \text{Re}\{Y_{n,k} \cdot Y_{n,k-1}^*\} \\ LLR(b_l^H) &= \text{Im}\{Y_{n,k} \cdot Y_{n,k-1}^*\} \end{aligned} \tag{2}$$

Where superscript $*$ denotes complex conjugate b_0^H and b_1^H are the most significant bit of $Y_{n,k}$ and the least bit of $Y_{n,k}$ for HP symbols, respectively. Demodulation of the LP symbols can be effectively performed by comb-type pilot arrangement.

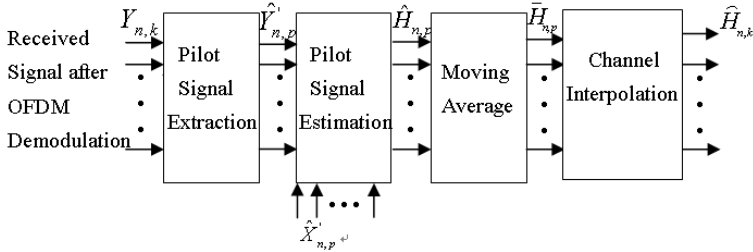


Fig. 3. Channel estimation for LP signal correction based on pilot

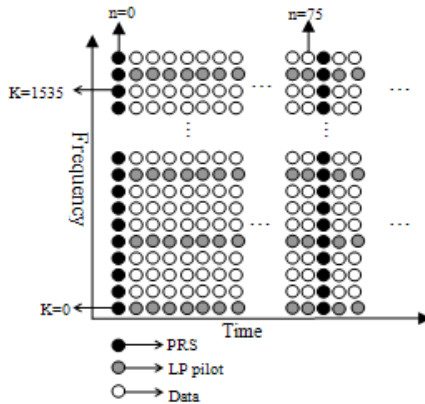


Fig. 4. A comb-type LP pilot arrangement for the hierarchical modulation

2.3 LP Demodulation for AT-DMB with Hierarchical Modulation

For coherent detection of LP-4ASK or LP-QPSK signals, the same channel estimate method is adopted as shown in Fig.3. The equalization process can be realized by a multiplier bank at the FFT output in the receiver. For comb-type pilot subcarrier arrangement, the N_m pilot signal $\hat{X}'_{n,p}$, $p=0, 1 \dots N_m-1$ are uniformly inserted into $X_{n,k}$ in Fig.4. That is the total subcarriers N are divided into N_m groups, each with $L=N/N_m$ adjacent subcarriers. The first pilot subcarrier is used to transmit pilot signal. subcarrier is used to transmit pilot signal. The OFDM signal modulated on the k th sub-carrier can be expressed as

$$X_{n,k} = X_{n,pL+l} = \begin{cases} \hat{X}'_{n,p}, & l = 0 \\ LP \text{ data}, & l = 1, 2 \dots l-1 \end{cases} \quad (3)$$

Fig.3 shows a block diagram of channel estimation for LP signal correction based on pilot. where $Y_{n,k}$ is the received symbol $\hat{Y}_{n,p}$ is the received pilot symbol extracted from $Y_{n,k}$, $\hat{X}'_{n,p}$ is a decided LP pilot symbol. We use a LS estimator, so the pilot signal response is estimated by:

$$\hat{H}_{n,p} = \hat{Y}_{n,p} / \hat{X}_{n,p} = \hat{H}'_{n,p} + \hat{N}'_{n,p} / \hat{X}'_{n,p} \quad (4)$$

In order to improve the correction of the estimated channel coefficients, we adopt moving average method in the frequency for noise reduction. They are given by:

$$\bar{H}_{n,p-qpsk} = \begin{cases} \frac{1}{21}(11-p)\hat{H}_{n,l} + \sum_{l=-10}^{10}\hat{H}_{n,l}, & p \leq 10 \\ \frac{1}{21}\sum_{l=-10}^{10}\hat{H}_{n,p+l}, & 10 < p < 1525 \\ \frac{1}{21}(p-1526)\hat{H}_{n,1535} + \sum_{l=p-10}^{1535}\hat{H}_{n,l}, & 1525 < p \end{cases} \quad (5)$$

$$\bar{H}_{n,p-4ASK} = \begin{cases} \sum_{m=0}^{p+3} W_7(m) |\hat{H}_{n,m}| / \sum_{m=0}^{p+3} W_7(m), & p \leq 3 \\ \sum_{m=p-3}^{p+3} W_7(m) |\hat{H}_{n,m}| / \sum_{m=p-3}^{p+3} W_7(m), & 3 < p \leq 1535 \\ \sum_{m=p-3}^{1535} W_7(m) |\hat{H}_{n,m}| / \sum_{m=p-3}^{1535} W_7(m), & 1535 < p \end{cases} \quad (6)$$

$$\text{Where } W_7(m) = 0.54 - 0.46 \cos(2\pi m/7) \quad (m \neq 7)$$

After noise reduction, we consider a piecewise-linear interpolation method for channel interpolation. So they are produced by:

$$\begin{aligned} \hat{H}_{n,k} &= \hat{H}_{n,pL+l} = \bar{H}_{n,p} + l/L(\bar{H}_{n,(p+l)} - \bar{H}_{n,p}) \\ 0 &\leq l \leq L, \quad k = 0, 1, \dots, 1536 \end{aligned} \quad (7)$$

Once $\hat{H}_{n,k}$ is obtained, the received hierarchical modulation signal $Y_{n,k}$ can be coherently detected by being equalized as follows

$$\hat{X}_{n,k} = Y_{n,k} / \hat{H}_{n,k} \quad (8)$$

The soft bit metric is given by applying the simplified LLR approximation functions presented in [11] as follows:

$$\begin{aligned} LLR(b_o^L) &= |\hat{H}_{n,k}|^2 \operatorname{Re}\{\hat{X}_{n,k}^L\} \\ LLR(b_o^L) &= |\hat{H}_{n,k}|^2 \operatorname{Im}\{\hat{X}_{n,k}^L\} \end{aligned} \quad (9)$$

Where $\hat{X}_{n,k}^L$ is given by

$$\hat{X}_{n,k}^L = \hat{X}_{n,k} - \hat{X}_{n,k}^H \quad (10)$$

In the LP symbol demodulation, the received complex symbols are converted into soft-bit information which is weighted by the channel state information (CSI) coefficients. Here $|\hat{H}_{n,k}|^2$ represents CSI.

3 Simulation Results

Computer simulation were performance by detecting they require SNR for a fixed Bit Error Rate (BER) of 10^{-4} in order to compare the performance of the proposed hierarchical LP-4ASK and LP-QPSK modulation for the AT-DMB system. The system parameters for T-DMB used in the simulation are given in Table 1. Simulation was performed in a transmission mode I environment of DAB [1]. Here we assume that guard intervals are longer than the maximum delay spread of the channel. The carry frequency is 200MHZ.The channel profile was COST207TU6 [12]. The channel parameters are shown in the Table 2.

The BER curves of soft-decision decoding of the LP signal are shown in Fig.5 and Fig.6.The SNR difference at the BER of 10-4 was 2 dB for the vehicle speed of 60 km/h. LP-4ASK needed about 4 dB higher SNR than the LP-4QPSK when a vehicle speed was 180 km/h. The results of Fig.5 and Fig.6 show that the hierarchical LP-QPSK modulation was better than the LP-4ASK modulation scheme for higher data rate transmission.

Table 1. Channel profile(cost207-tu6)

Path Number	Rayleigh Channel	
	Power(dB)	Delay(μ s)
1	-3.0	0.0
2	0.0	0.2
3	-2.0	0.5
4	-6.0	1.6
5	-8.0	2.3
6	-10.0	5.0

Table 2. Simulaiton prrameters

Parameters	Specification
FFT size	2048
Number of transmitted carriers	1536
Total symbol duration	1.246ms
Useful symbol duration	1ms
Guard interval duration	246 μ s
Bandwidth	1.536MHZ
Modulation	HP-DQPSK/LP-QPSK/4ASK
Time interleaving(depth)	384ms
Frequency interleaving(width)	1.536MHZ
Channel model	COST207 TU6
Convolutional code	Rc=1/2, constraint length=7

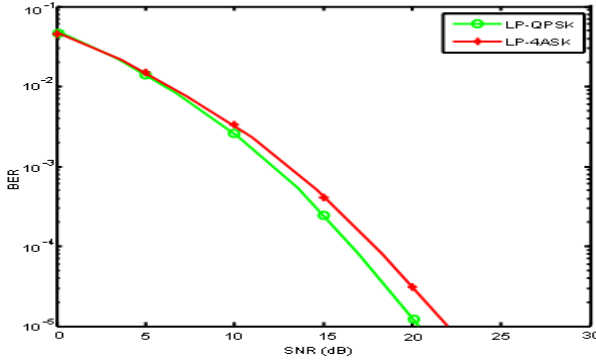


Fig. 5. The performance of the 4ASK ($\lambda_1 = 3$) and QPSK ($\lambda_2 = 2$) modulation ($v=60\text{km/h}$)

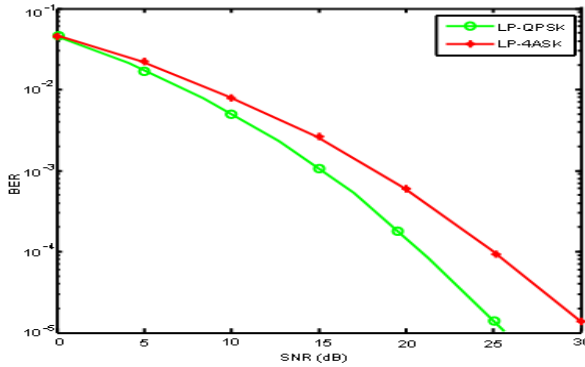


Fig. 6. The performance of the 4ASK ($\lambda_1 = 3$) and QPSK ($\lambda_2 = 2$) modulation ($v=180\text{km}$)

4 Conclusion

In a word, T-DMB with the higher data rate can meet our demands for the multi services. Hierarchical modulation is one practical implementation of superposition precoding, which can help achieve the maximum sum rate. It not only upgrades the data rate but also is performed with backward compatibility. In the paper, perfect channel estimation has been assumed in case of the LP symbol. Simulation results show that the performance of LP-QPSK is better than it of LP-4ASK. It provides a reference for our research about hierarchical modulation. Meanwhile, we can use the T-DMB with hierarchical modulation to the local networks for transmitting the local information for higher spectrum efficiency, as well as to cooperative communication system. However, the issue of the channel coding about LP data has not been addressed in this paper. Further studied are considered to find an optimal channel coding for hierarchical modulation.

Acknowledgments. This work is supported by Natural Science Foundation of Jiangsu Province (Grant No. BK2011475) and Foundation of Jiangsu University (Grant No. 10JDG115 and No. 11JDG034).

References

1. ETSI: Radio Broadcasting Systems, Digital Audio Broadcasting (DAB) to mobile, portable and fixed receivers. EN 300 401 (2001)
2. Han-gil, M.: Backward-compatible Terrestrial Digital Multimedia Broadcasting system for multichannel Audio service. *IEEE Trans. Consumer Electronics*, 1556–1561 (2010)
3. Jae Hang, L., Jong-Sool, Sang, W.: Development of Advance Terrestrial DMB system. *IEEE Trans. Consumer Electronics*, 28–35 (2010)
4. Chul Seung, K., Min, H., Ji won, J.: Hierarchical Transmission Algorithm in the Advanced T-DMB system. *IEEE* (2010)
5. Ma, L., Tang, C., Wang, C.: A Kind of Hierarchical Data Storage Management system Design for Wireless Sensor Network. In: *Proceeding of the 30th Chinese Control Conference*, pp. 5034–5038 (2011)
6. Dae-Ken, K., Wan-Jin, K., Ki-Hwan, S., Hyungsoo, L., Hyung-Nam, K.: A Higher Data Rate T-DMB Based on a Hierarchical A-DPSK Modulation. *IEEE Trans. Consumer Electronics*, 42–50 (2009)
7. Wan-Jin, K., Young-Jun, L., Hyung-Nam, K., Hyungsoo, L., JongSoo, L.: Coded decision-directed channel estimation for coherent detection in terrestrial DMB receivers. *IEEE Trans. Consumer Electronics*, 815–819 (2007)
8. Hsieh, M.-H., Wei, C.-H.: Channel estimation for OFDM systems based on comb-type pilot arrangement in frequency selective fading channels. *IEEE Trans. Consumer Electronics*, 217–225 (1998)
9. Sun, J., Song, S., Liu, Y.: Model checking hierarchical probabilistic systems. In: *ICFFM*, pp. 388–403 (2010)
10. Hewavithana, T.C., Brookes, M.: Soft decisions for DQPSK demodulation for the Viterbi decoding of the convolutional codes. In: *Proceedings of ICASSP 2003*, pp. 17–20 (2003)
11. Tosato, F., Bisaglia, P.: Simplified soft-output demapper for binary interleaved COFDM with application to HIPERLAN2, pp. 664–668. *IEEE* (2002)
12. COST 207 Report Digital Land Mobile Radio Communications, Commission of European Communities, Directorate General, Telecommunications, Information Industries and Innovation (1989)

Research on the Controlling Technique of DFBLD in the Spectrum Absorption Optical Fiber Gas Detecting System

Shutao Wang, Pengwei Zhang, and Xiaoqing Shao

Key Laboratory of Measurement Technology and Instrumentation of Hebei Province,
Yanshan University, Qinhuangdao, 066004, China
wangshutao@ysu.edu.cn, {zpw1004, sxqxf}@163.com

Abstract. Spectrum absorption optical fiber gas sensor has characteristics such as high measuring resolution, high precision and good gas differential ability when it is applied on gas measurement. As ideal illuminant and key component, the research of DFBLD seems especially important. Constant current driving and constant temperature controlling techniques are adopted on DFBLD to make the central wavelength of it accurately point at detected gas absorption. Theory analysis and experiment results indicate that the optical fiber gas sensor based on DFBLD has specified detecting precision and resolution.

Keywords: Optical fiber gas sensor, DFBLD, Constant current driving, Constant temperature controlling.

1 Introduction

To develop timely and accurate sensor system for prediction and control inflammable, explosive, toxic gas has become the serious problems to be solved in the current coal, petroleum, chemical industry, electric power, transportation, aerospace and environmental monitoring, etc. Many kinds of gas have been intrinsically absorbed in infrared wavelength region and have a weak wide-frequency absorption spectrum in near-infrared and visible wavelength regions. And the low loss window of quartz fiber, covers $0.8 \sim 1.7 \mu\text{m}$ range. In this band light emitting device and receiving device are very ideal photoelectric conversion devices. And optical fiber sensor gas can be used for environmental testing and process control, especially in continuous and online monitoring under adverse circumstances. It can play an important role in the fiber sensor based on the mechanism of spectrum absorption. It possesses lots of special characteristics such as low circuitry loss, great information transmission capacity, strong anti-jamming capability, corrosion resistance, insulation, fire and explosion prevention, realization of long distant real-time data gather and data management, and high gas selectivity, and is considered the replacement products for the traditional hot gas sensor. Fiber Bragg grating (FBG) sensor is one of the current heat research fields in the optical fiber sensor technique fields. This system, consisting of an illumination source DFBLD in $1.66\mu\text{m}$ waveband, through constant current source drive and constant temperature control technology to DFBLD,

overcomes the influence of temperature and other environmental factors, can greatly improve the detection system detection accuracy and resolution, reaches measurement engineering standards and requirements[1-2]

2 Testing Principle

When light pass through a medium, even if not occur reflection, refraction and diffraction phenomenon, its condition will also be changed. As the light frequency electromagnetic wave and composition of the medium atoms, molecules will exert an influence, and the light wave characteristic will change, producing absorption and scattering phenomenon. Based on the characteristics of gas absorption to light varying with concentration, absorption fiber-optic gas sensor is developed. Based on gas absorption characteristic it is known that when light passes the media some of it is absorbed, some is scattered and the remaining transfers directly. It can be concluded from law of Beer-Lambert that the intensity of light I passing the gas satisfies the relationship of

$$I = I_0 \exp(-\alpha_\lambda cL) \quad (1)$$

I is passing light intensity; I_0 is the incident monochromatic light intensity with wavelength, α_λ is the absorption coefficient of media with unit concentration and unit longitude, c is the concentration of the detected gas. L is the affecting longitude of detected gas and light. Transform (1), get

$$c = \frac{1}{\alpha_\lambda L} \ln \frac{I_0}{I} \quad (2)$$

The absorption spectrum line of gas to light is shown in figure 1. The gas spectrum line has a relative strong absorption peak nearby 1331.5nm. The absorption of peak value is 5.72dBm. At the direction of shortwave there also exist a set of weak absorption peak. But the light is almost not absorbed by the gas between 1318.26nm and 1320.18nm. So the wavelength of 1319.145nm and 1331.500nm are chosen for the measurement of gas.

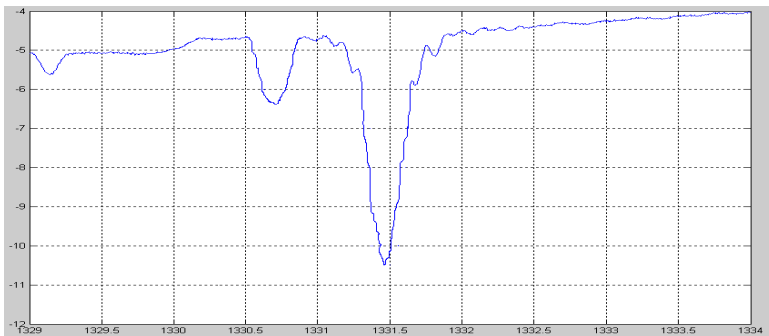


Fig. 1. The absorption spectrum of methane in the waveband of 1300nm

In fiber gas sensor design, the complex refractive index k_0 is used to describe the gas absorption. The incident plane wave is set for $E_0 \exp(i\omega t)$, the output light is

$$E = E_0 \exp(i(\omega t - \beta L)) \tag{3}$$

$$\beta = k_0 n = k_0 (n_r - ik) \tag{4}$$

$$k_0 = \frac{2\pi}{\lambda} \tag{5}$$

β is propagation Constants; n_r is the real part of gas refractive index; k is the Imaginary part of gas refractive index. The output light intensity is proportional to the square of the electric light strong intensity.

$$I = I_0 \exp(-2k_0 k L) \tag{6}$$

Compare (6) with (1), get

$$k = \frac{\alpha_\lambda c}{2k_0} = \frac{\alpha_\lambda c \lambda}{4\pi} \tag{7}$$

when there is no the sample gas, light intensity I_0 is measured by its detector in the air chamber. When the sample gas get a certain concentration c , the gas has a strong absorption to the light. Light intensity is I . The ratio of I_0 and I can characterize the gas concentration optical fiber gas detection[4].

3 Testing System Design

The system consists of the light source, the gas chamber, the fiber optic light path and the signal processing unit. Fiber spectrum absorption gas detection system based on frequency modulation DFBLD principle diagram is shown in Figure 2.

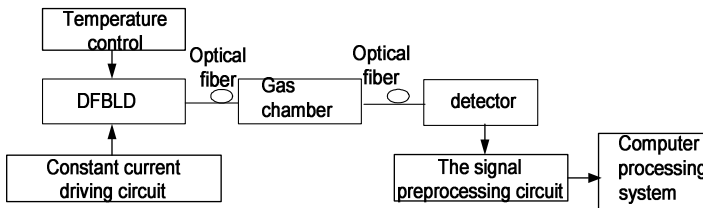


Fig. 2. Constitution block Diagram of optical fiber gas detecting system

The system is designed for testing the content of methane and acetylene gases. Methane has two strong absorption peaks respectively in 1.66 μ m band and 1.33 μ m band. The system based on 1.66 μ m DFBLD light (corresponding to the methane gas absorption peaks), using spectrum absorption testing technology, realize to detect the

concentration of the methane. Detection system use 1.65 KHZ incentive signal as the light source DFBLD constant drive current. The output light of the laser pass by 10 km long input optical fiber, get into gas sensor chamber assembled the beam of collimator on both ends [6]. When the beam passing by the measuring gas, enter into 10 km long output optical fiber through the beam collimator and is sent to InGaAsPIN photoelectric detector[7]. The output signal from photoelectric detector is sent into the lock-in amplifier, detecting modulation frequency and quadratic harmonic component. The signal passes signal processing circuit and then is processed by computer, displaying the content of the measuring gas. This system chooses DFBLD (Distributed Feedback Laser Diode) as stimulate light source. DFBLD light source has good characteristics such as the narrow spectral lines, the good performance wavelength stability, etc. In constant temperature conditions, it can steadily make stimulated light center wavelength accurately aim at measured gas absorption peak. DFBLD drive circuit with light power feedback has been designed to make that DFBLD output light power constant .Its principle is as shown in figure 3 below. It provides active control for driving power through the sampling current feedback. when DFBLD works, it will receive a small part of the internal PD optical power and translated it into monitoring current I_1 [8].

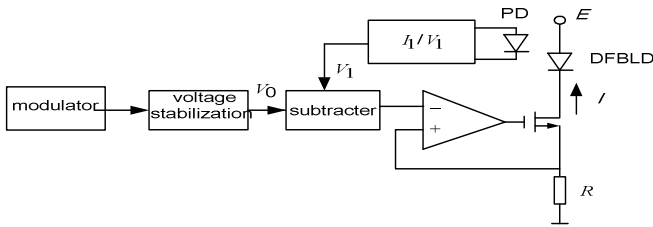


Fig. 3. Principle of light feedback and constant-current driving of DFBLD

In figure.3 , V_1 is modulation voltage of the optical current I_1 of PD conversion from the monitoring .Subtraction of the voltage and the modulation of drive circuit voltage form the DFBLD drive current end.

$$I = \frac{V_0 - V_1}{R} K \quad (8)$$

K is conversion coefficient. DFBLD driver current is adjusted by R . V_1 increase, and I decrease, making the DFBLD output power decrease. And vice is still. At the same time, the DFBLD drive current has relations with $(V_0 - V_1)$ in the drive circuit and R -resistance, noting to do with others factors, forming constant current drive with the ability of high anti-interference.

4 Experiment Results

Main performance of the system at 1665 center wavelength is tested for verifying the feasibility of DFBLD light source of methane gas detection. In the spectrum test

characteristics, the spectral instrument Canadian company EXFO with the resolution of 0.002 nm is used, and the EXFO company’s IQS-1100 optical power meter is adopted. DFBLD’s tail fiber is 2 m. The measured results: center wavelength is 1665.60 nm. 3 dB spectral bandwidth is 2.03 nm; The optical power output is 3.63 dBm; Short-term stability of Light output power is + 0.001 dB (15 minutes, constant temperature, continuous light output); Long-term stability of output light power is + 0.003 dB (8 hours, + 2°C, continuous light output). The relation of wavelength-temperature of the system using DFBLD is tested[9-10].

In the test, controllable constant water channel is adopted as temperature source ranging from 20 to 80°C, and the resolution for 1°C. Not in the DFBLD cooling control, in 20°C ~ 40°C temperature range, the measurement results of the output of the DFBLD wavelength along with the change of the environment temperature are shown in figure 4 .

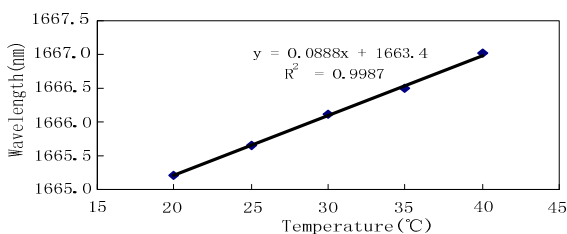


Fig. 4. The curve of output wavelength of the DFB-LD changing with temperature

Figure 4 shows that in the temperature range, with the increase of temperature, DFBLD output wavelength increase, which is approximately linear (linearity R2 = 0.9987). In 20°C ~ 40°C, the DFBLD peak wavelength range from 1665.21 to 1667.01 nm .The rate is about 0.09 nm / °C. The measurement result of DFBLD spectrum width is 2.03 nm.

Table 1. Experimental data of output wavelength of the DFB-LD changing with temperature

	D_1	D_2	D_3	D_4	D_5	average \bar{D}	Standard DeviationS
1	1665.270	1665.271	1665.269	1665.270	1665.270	1665.270	0.0007
2	1665.685	1665.686	1665.685	1665.685	1665.684	1665.685	0.0007
3	1666.101	1666.101	1666.100	1666.100	1666.008	1666.100	0.0012
4	1666.500	1666.500	1666.500	1666.501	1666.499	1666.500	0.0007
5	1667.010	1667.010	1667.011	1667.010	1667.009	1667.010	0.0007

Keep constant temperature cooling condition in 25°C. The stability of the DFBLD peak wavelength is tested to ensure the output of the DFBLD wavelength accurate locking in the methane gas absorption peaks. In the test, current I = 40 mA, frequency f = 1.65 KHz square wave signal is used by DFBLD driver signal.

The experimental results show when environment temperature range 20 °C to 40 °C, the output of the DFBLD peak wavelength range from 1665.65 nm to 1665.71 nm as a rate of 0.003 nm / °C in 25 °C constant temperature cooling conditions. The change of environment temperature has a little effect on the peak wavelength change. The experimental results show that, using DFBLD for constant temperature cooling, its output wavelength basically accurately locks in the gas absorption peaks, which can guarantee the high measuring accuracy.

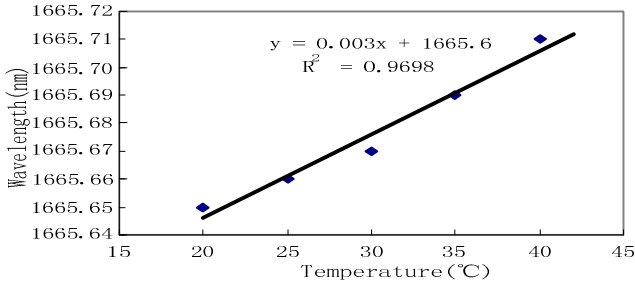


Fig. 5. The output wavelength of DFBLD changing with temperature by refrigerating at 25°C

5 Conclusion

Theoretical analysis and experimental results show that, the spectrum of light source based on DFBLD absorption optical fiber sensing system expounded by this paper satisfies actual requirements of the measuring gas. 1.66 μm DFBLD is adopted by the system as the light source, through control technology of the DFBLD constant current source and constant temperature, making the center of the DFBLD light wavelength accurate aim at methane gas absorption peaks.

Experiment of DFBLD light source test and sensitivity and resolution of system obtained a satisfactory result, proving that DFBLD light source control technology can restrain the influence of temperature and other environmental factors. Using long distance near infrared low loss quartz fiber optical as the stimulate light and the signal transmission medium and adopting multiplexing technique, can realize the continuous monitoring of environmental gas online and remote sensing. Further research on system feedback technology, the stability of light source to the influence of the system sensitivity, can make the system have a higher stability and reliability.

Acknowledgements. This work was financially supported by the Hebei Natural Science Foundation (F2010001313), the National Natural Science Foundation of China (60974115) and the Key Project of Chinese Ministry of Education (210025).

References

1. Othonos, A., Kalli, K.: *Fiber Bragg Gratings Fundamentals and Applications in Telecommunications and Sensing*. Artech House Press, Norwood (1999)
2. Weis, R.S., Kersey, A.D., Berkoff, T.A.: A Four-element Fiber Grating Sensor Array with Phase-sensitive Detection. *IEEE Photonics Technology Letters* 6(12), 1469–1472 (1994)

3. Jackson, D.A., Ribeiro, A.B.L., Reekie, L.: Simple Multi-plexing Scheme for a Fiber-optic Grating Sensor Network. *Optical Letter* (1993)
4. Ball, G.A., Morey, W.W., Cheo, P.K.: Fiber Laser Source/analyzer for Bragg Grating Sensor Array Interrogation. *Journal of Lightwave Technology* 12(4), 700–703 (1994)
5. Kersey, A.D., Berkoff, T.A., Morey, W.W.: Multiplexed Fiber Bragg Grating Strain-sensor System with a Fiber Fabry-Perot Wavelength Filter. *Optical Letter* 18(16), 1370–1372 (1993)
6. Gharavi, M., Buckley, S.G.: Diode laser absorption spectroscopy measurement of linestrengths and pressure broadening coefficients of the methane $2\nu_3$ band at elevated temperatures. *Journal of Molecular Spectroscopy* 229, 78–88 (2005)
7. Benounis, M., Aka-Ngnui, T.: NIR and optical fiber sensor for gases detection produced by transformation oil degradation. *Sensors and Actuators A Physical* 141, 76–83 (2008)
8. what's in Cyc [OL] (2007), <http://www.Cyc.CH4m>
9. Ding, K.-L., Li, S.-Y.: Carbon isotopic evidence for oxidation of CH₄ by TSR. *Utilization and Environmental Effects* 33(18), 1715–1725 (2011)
10. Zhou, J.: Infrared Gas Sensor, <http://www.intlsensor.com/p/df/infrared.Pdf>
11. Minming, T., Jieru, N., Cunliang, J.: Study on a fiber optic sensor for CH₄ in mine. *Geo-Marine Letters* 31(1), 37–49 (2011)
12. Peter, R., Andreas, S.: Fire detection in coal mines based on semiconductor gas sensors. *Sensor Review* 32(1), 47–58 (2012)

An Efficient and Provable Secure PAKE Scheme with Robust Anonymity

Cong Liu and Chuan-gui Ma

Zhengzhou Information Science and Technology Institute,
Zhengzhou, Henan Province, 450002, China
congliu85@163.com, chuanguima@sina.com

Abstract. In 2009, Sun et al. proposed an improved anonymous scheme on Juang et al.'s password-authenticated key exchange (PAKE) scheme using smart cards. Although, this scheme overcomes some weaknesses of Juang et al.'s scheme and achieves user anonymity, it is still lack of the untraceability property such that the third party over the communication channel can recognize whether or not the same user in different sessions. In this paper, we propose a new robust anonymous PAKE scheme using smart cards which not only can strengthen the security of Sun et al.'s scheme by addressing untraceability, but also can achieve greater communication efficiency. Moreover, we present a strict security proof for our scheme.

Keywords: untraceability, anonymity, provable secure, PAKE, smart card.

1 Introduction

In the present network, anonymity as an important and effective method of privacy protection arouses wide concern. For the design of anonymous PAKE schemes, the designers want to use some safe and efficient methods of mathematics to achieve anonymity, but they often neglect the randomness of the identity protection, this will result that the adversary can recognize the same user in different sessions, and then destroy user anonymity. Hence, there is a stronger property than anonymity called “untraceability” which must be needed in anonymous PAKE schemes.

In 2005, Fan et al. [1] proposed a robust remote authentication scheme with smart cards. Although it can resist the offline attack even if the smart card is compromised, but also has higher computation and communication cost and does not have the functions of session key agreement and anonymity. Juang et al. [2] proposed a PAKE scheme using smart cards in 2008. The merits of Juang et al.'s scheme are addressing the threat of the smart card loss and the use of the elliptic-curve algorithm for reducing the implementation costs. In 2009, Sun et al. [3] pointed out that Juang et al.'s scheme suffers from some weaknesses and proposed a more efficient improved scheme which fixes the weaknesses and reduces the storage and computation costs on the smart card. In 2010, Li et al.

[4] noticed the anonymity-related issues of Juang et al.’s scheme, and pointed out it does not satisfy untraceability property. The improved scheme is proposed by Li et al. which is similar in the process structure to Juang et al.’s, so it has lower efficiency than Sun et al.’s scheme. Moreover, some anonymous PAKE schemes only based on the password [5-7] are proposed but which fail to resist online dictionary attack. In this paper we point out Sun et al.’s scheme is also lack of untraceability, and propose an improved PAKE scheme which not only reserves the merits of Sun et al.’s, but also provides anonymity and untraceability. Due to using smart cards, it can reduce the probability of guessing user’s password and resist online dictionary attack. Comparing with the schemes of Sun et al.’s and Li et al.’s, our protocol only has two flows. Hence we believe our protocol is a safe, efficient and more suitable scheme for real-life applications. Furthermore, we strictly prove the security of our scheme in the random oracle model.

2 Review of Sun et al.’s Scheme

A. Parameter-Generation Phase

The server S chooses an elliptic curve E over a finite field F_p we suppose that the discrete logarithm problem is hard in $E(F_p)$. S also chooses a point P such that the subgroup G generated by P has a large order n . S publishes the parameters (p, E, P, G, n) .

B. Registration phase

We suppose the messages are transmitted in this phase via a secure channel.

1. U selects the sub-identifier ID_U and submits it to S .
2. Upon receiving messages from U , S checks if ID_U is valid. If yes, S selects the sub-identifier ID_S and generates $ID = ID_U || ID_S$. Then, S computes $V = h(ID || k) \oplus h(PW)$, $IM = E_k(ID || r)$, where PW is the selected by S , r is a random number and k is the master secret key only kept by S . S stores V and IM in smart card delivered to U .

C. Authentication Phase

1. U inputs his password PW^* . The smart card randomly selects an integer $r_C \in [1, n - 1]$ and computes $P_C = r_C \times P$.
2. Upon receiving the message $\{IM, P_C\}$, S computes $E_k^{-1}(IM) = ID || r$. Then S checks if ID is valid. If yes, S randomly selects an integer $r_S \in [1, n - 1]$, computes $P_S = r_S \times P$, $K_{SU} = h_1(h(ID^* || k) || r_S \times P_C)$ and $M_S = h_2(K_{SU} || P_C || P_S)$, and sends $\{M_S, P_S\}$ to U .
3. Upon receiving the message $\{M_S, P_S\}$, U computes $V' = V \oplus h(PW^*)$ and $K_{SU} = h_1(V' || r_C \times P_S)$, and compares M_S with $h_2(K_{SU} || P_C || P_S)$. If they are equivalent, the smart card computes $M_C = h_2(K_{SU} || P_S)$ and then sends $\{M_C\}$ to S .
4. Upon receiving the message $\{M_C\}$, S checks if M_C is equal to $h_2(K_{SU} || P_S)$. If yes, U and S successfully authenticate each other and establish the session key K_{SU} . Otherwise, S terminates this session.

D. Password-Change Phase

This phase is invoked whenever U wants to change his password with a new one PW' . U enters the old password PW' , and requests to change password. Next, U enters the new password PW' . Smart card computes $V^* = V \oplus h(PW) \oplus h(PW') = h(ID||k) \oplus h(PW')$, then replaces V with V^* .

3 The Weakness of Sun et al.'s Scheme

The schemes of Juang et al.'s and Sun et al.'s provide the initiator anonymity (i.e. user anonymity) that protects the user's identity, but they neglect a more ideal anonymity property—untraceability which means any third party over the communication channel cannot tell whether or not he has seen the same (unknown) smart card twice through the authentication sessions.

From the message flow of authentication phase in the step 1, U want to send the message $\{IM, P_C\}$ to S , one may notice that IM is fixed, since ID is determined by user's ID_U and server's ID_S , r provided by S in registration phase, and the long-term secret key k of the server. As long as U does not reregister to S , IM is always the same number. This will bring the weaknesses of internet security. On one hand, we supposed that the adversary obtains and stores IM in step 1 of authentication phase. If the adversary accidentally knows the real user's identity, he will recognize the user immediately by detecting the value of IM , once the user initiates a new authentication again. On the other hand, if the adversary do not know the real user's identity, he also can detect whether or not the users are the same one in different authentications through viewing IM which got from the message flow of authentication phase in the step 1. If the same one, he could pay attention to the user and get other useful information related to the identity of the user. The information may help the adversary to compromise the anonymity of Sun et al.'s.

Therefore this scheme would exigently need the untraceability property as the stronger property than anonymity. To make up the weakness of Sun et al.'s, we propose the new anonymous PAKE scheme using smart card in the following.

4 Our Proposed Scheme

The weakness of Sun et al.'s scheme root in the fact is that every message in the run of the scheme doesn't be randomized.

A. Parameter-Generation Phase

Our scheme still preserves the low costs of computation and storage based on ECC. In the phase, all parameters are generated, same as Sun et al.'s scheme. Therefore we omit the introduction.

B. Registration phase

The messages are also transmitted in this phase via a secure channel.

1. U selects the password PW , submits the identity $ID \in \{0, 1\}^{k_1}$ and PW to S for registration.
2. Upon receiving messages from U , S checks if ID is valid. If yes, S has a private secret key $a \in [1, n - 1]$, and then computes $A = a \times P$ and $V =$

$h_0((ID\|a) \times P) \oplus h_0(PW)$, and delivers a smart card to U , $h_0()$, $h_1()$, $h_2()$, A and V are stored in the smart card, where $h_i() : \{0, 1\}^* \rightarrow \{0, 1\}^{l_i}$. l_i, k_j are security parameters.

C. Authentication Phase

This phase is initiated whenever U wants log in to S . The steps of this phase are shown as follows.

1. U inputs his password PW^* and the identity ID^* . The smart card randomly chooses $r_U \in [1, n - 1]$, $n_U \in \{0, 1\}^{k_2}$. S computes $P_U = r_U \times P$, $D = r_U \times A$, $k = h_0(D\|A\|P_U)$, $V' = V \oplus h_0(PW^*)$, $R_U = \mathcal{E}_k(ID^*\|n_U)$, $M_U = h_1(D\|V'\|P_U\|n_U)$. U sends $\{M_U, R_U, P_U\}$ to S , where \mathcal{E} is an IND-CPA security symmetric encryption using key k .
2. Upon receiving the message $\{M_U, R_U, P_U\}$, In the same way, S computes $D^* = a \times P_U$, $k = h_0(D^*\|a \times P\|P_U)$, $\mathcal{E}_k^{-1}(R_U) = ID^*\|n_U$. Firstly, S checks if ID^* is valid. If yes, S computes $W = h_0((ID^*\|a) \times P)$. Then S checks whether $h_1(D^*\|W\|P_U\|n_U)$ equals M_U . If yes, S randomly selects an $r_S \in [1, n - 1]$, computes $P_S = r_S \times P$, $M_S = h_1(D^*\|n_U\|P_U\|P_S)$ and $K_{SU} = h_2(n_U\|P_U\|P_S\|r_S \times P_U)$. S delivers $\{M_S, P_S\}$ to U .
3. Upon receiving the message $\{M_S, P_S\}$, U checks whether $h_1(D\|n_U\|P_U\|P_S)$ equals M_S . If yes, U computes $K_{SU} = h_2(n_U\|P_U\|P_S\|r_C \times P_S)$. Otherwise, the smart card terminates this session.

D. Password-Change Phase

We reserve the merit of Sun et al.’s. in this phase such that U can change his password freely without any interaction with S so as to reduce the possibility of the insider attack. U enters the old password PW' , and requests to change password. Next, U enters the new password PW' . Smart card computes $V^* = V \oplus h_0(PW) \oplus h_0(PW') = h_0((ID\|a) \times P) \oplus h_0(PW')$, then replaces V with V^* .

5 Security of Our Scheme

5.1 Formal Security Model

In this subsection, we introduce the formal security model based on Bellare et al. [8] and Zhou et al. [9].

Participants and Initialization. we denote a user U and a server S that can participate in this PAKE scheme P . Each of them may have several instances called oracles involved in distinct, possibly concurrent, executions of P . We denote user instances and server instance by U^i and S^j (or by I^i when we consider any participant). S has a private key K_S . U has a password PW_U drawn from a small dictionary *Password* of size N according to the distribution D_{pw} . Additionally, when U enrolls S , S stores V_U in the smart card and deliver it to U , where V_U is an transformation of PW_U and K_S .

Queries. The adversary A interacts with the participants by making oracle queries. The oracle queries are explained in the following:

- *Excute*(U^i, S^j): This query models eavesdropping attacks of the adversary. It outputs the honest executions of P between U^i and S^j .

- *Reveal(I^i)*: This query models the misuse of the session key by instance I^i . The output of this query is the session key sk to A , if I^i actually “holds” sk .
- *Send(I^i, m)*: This query models active attacks. The output is the message generated by I^i in processing the message m according to P .
- *Corrupt(U, a)*: This query models corruption capabilities the adversary. If $a = 1$, output the password PW_U of U ; If $a = 2$, output the messages stored in the smart card.

5.2 Security Notions

Freshness. An instance is said to be **Fresh** in the current scheme execution if the instance has accepted and neither it nor the other instance with the same session tag have been asked for a Reveal-query.

Test-Query. The additional query $Test(I^i)$ models the semantic security of the session key. It can be asked by A only once if the instance I^i is **Fresh**. It outputs as following: one flips a (private) coin b and forwards sk (the value $Reveal(I^i)$ would output) if $b = 1$, or a random value if $b = 0$.

AKE Security. The semantic security of the session key is modeled by the game $Game(A, P)$ in which the adversary A makes a single Test-query $Test(I^i)$. Playing this game aims to guess the bit b involved in the Test-query, by outputting this guess b' . Event $Succ(A)$ occurs if $b = b'$. We define $Adv_P^{ake}(A) = 2Pro(b = b') - 1 = 2Pro(Succ(A)) - 1$. The scheme P is said to be **AKE-secure** if A 's advantage is negligible in the security parameter.

User Anonymity. To define anonymity, there are two queries are needed.

- *RevealID(I^i)*: This query models the misuse of user's identity by instance I^i . The output of this query is user's real identity of ID to A , if I^i actually “holds” ID .
- *TestAnon(I^i, ID_0, ID_1)*: This query models is to define anonymity of user identity. After querying the oracle, the transcript of I^i with identity ID_0 or ID_1 will be returned according to a predefined random bit c , If $c = j, j \in \{0, 1\}$, A would learn the transcript of I^i with identity ID_j .

For any adversary A , the semantic security of user anonymity is modeled by the game $GameAnon(A, P)$ in which the adversary A makes a single TestAnon-query. Playing this game aims to guess the bit c involved in the TestAnon-query, by outputting this guess c' . Event $Succ^{anon}(A)$ occurs if $c = c'$. We define $Adv_P^{anon}(A) = 2Pro(Succ^{anon}(A)) - 1$. The scheme P provides user anonymity if $Adv_P^{anon}(A)$ is negligible in the security parameter.

5.3 Security Proof

Theorem 1. Let G be a prepresent group, $Password$ be a finite dictionary of size N with uniform distribution and \mathcal{E} is an IND-CPA security encryption.

Let \mathcal{A} be an adversary against the AKE security of our scheme P with in a time bound t , with less q_s Send-queries and q_e Execution-queries, and asking q_h Hash-queries. Then we have

$$\begin{aligned}
 Adv_P^{ake} \leq & \frac{q_h^2}{2^{t_1}} + (q_s + q_e)^2 \left(\frac{1}{n} + \frac{1}{|\mathcal{C}|}\right) + 2(q_e + q_h) \left(\frac{1}{2^{k_1} \cdot 2^{k_2} \cdot (n-1)} + \frac{1}{2^{k_2} \cdot (n-1)}\right) \\
 & + \frac{1}{2^{k_2}} \times Succ_G^{dh}(t') + Adv_{\mathcal{A}, \mathcal{E}}^{cpa}(n') + 2q_s \left(\frac{1}{N} + Adv_{\mathcal{A}, \mathcal{E}}^{cpa}(n') + \frac{1}{2^{t_1-1}} + \right. \\
 & \left. 2q_h \left(\frac{1}{2^{k_1} \cdot 2^{k_2} \cdot (n-1)} + \frac{1}{2^{k_2} \cdot (n-1)}\right)\right)
 \end{aligned}$$

where $t' \leq t + (q_s + q_e + 1) \cdot \tau^G$, with τ^G denoting the exponentiation computational time in G .

Proof. We define a sequence of games starting at the real game G_0 and ending up the G_4 in which the adversary has no advantage. For each game G_i , We define an even $Succ_i$ occurs if $b = b'$, where b is the bit involved in the Test-queries, and b' is the output of the AKE-adversary. At the end of games, we compute the probability, $\Delta_i = |Pr[Succ_{i+1}] - Pr[Succ_i]|$ between in G_{i+1} and G_i for $0 \leq i \leq 4$.

Game G_0 . This is a real attack game in the random oracle model. We have $Adv_P^{ake} = 2Pr[Succ_0] - 1$. Therefore $Adv_P^{ake} = 2Pr[Succ_5] - 1 + 2(Pr[Succ_0] - Pr[Succ_5]) \leq 2Pr[Succ_5] - 1 + 2 \sum_{0 \leq i \leq 4} \Delta_i$.

Game G_1 . In this game, we simulate the random orales h_i (but also three additional random oracles $h'_i : \{0, 1\}^* \rightarrow \{0, 1\}^{l_i}$) as usual by maintaining a hash list $Ah(Ah')$. The *Excute, Reveal, Send, Corrupt* and *Test* oracles are also simulated as the real players do. We easily see that the game is indistinguishable from the real attack. Therefore, $\Delta_0 = 0$.

Game G_2 : In this game, all random oracles are simulated as in G_1 . we cancel the game in which the collisions occur in the transcript $\{\{M_U, R_U, P_U\}, \{M_S, P_S\}\}$. According to birthday paradox, the probability of collisions in the output of the random oracle is at most $\Delta_1 \leq \frac{q_h^2}{2^{t_1+1}} + (q_s + q_e)^2 \left(\frac{1}{2n} + \frac{1}{2|\mathcal{C}|}\right)$, where \mathcal{C} is the ciphertext message space.

Game G_3 : In this game, we consider the eavesdropping attacks generated via Execute query. We modify the Execute oracle so that the values of R_U, M_U, M_S, K_{SU} are selected uniformly at random. Due to the adversary \mathcal{A} can not know the server's private secret key a , ID , n_U and $r_U \times r_S \times P$, and then \mathcal{E} is an IND-CPA security encryption, \mathcal{A} does not distinguish four values from the random numbers. Hence G_3 and G_2 are indistinguishable unless the adversary can correctly guess the server's secret key, ID and n_U , solve the CDH problem and break the IND-CPA security. $\Delta_2 \leq q_e \times \left(\frac{1}{2^{k_1} \cdot 2^{k_2} \cdot (n-1)} + \frac{1}{2^{k_2} \cdot (n-1)} + \frac{1}{2^{k_2}} \times Succ_G^{dh}(t') + Adv_{\mathcal{A}, \mathcal{E}}^{cpa}(n')\right)$.

Game G_4 : We continue to weaken the ability of the passive attacks. Now, we compute the session key and the authentications using the private oracles h'_i , so that the values R_U, M_U, M_S and K_{SU} are completely independent from $h_i, D, r_u \times r_s \times P$. In the Execute queries, \mathcal{A} will get $R_U = \mathcal{E}_{k'}(ID^* || n_U)$, $M_U = h'_1(P_U)$, $M_S = h'_1(P_U || P_S)$, $K_{SU} = h'_2(P_U || P_S)$. The game G_4 and G_3

are indistinguishable unless break the IND-CPA security and some specific hash queries are asked by denoted the event $AskH_4$: \mathcal{A} queries the hash function h_1 on $D\|V'\|P_U\|n_U$, $D\|n_u\|P_U\|P_S$ and the hash function h_2 on $n_u\|P_U\|P_S\|r_u \times r_s \times P$. Therefore, $\Delta_3 \leq Pr[AskH] + q_h \times Adv_{\mathcal{A},\mathcal{E}}^{cpa}(n')$.

Next, we prove that if the event $AskH_4$ happened successfully, there is an algorithm which can break CDH-problem. Given two random instances $\{(U, V)(V, Q)\}$, we define an additional Game G'_4 . In Game G'_4 , we simulate the Execute queries by making $B = \alpha U, P_U = \beta V, P_S = \sigma Q$, where $\alpha, \beta, \sigma \in [1, n-1]$, and other oracles are same as Game G_4 . If the event $AskH_4$ happened, we can get the $D = CDH(\alpha U, \beta V) = CDH(U, V)^{\alpha\beta}, r_U \times r_S \times P = CDH(\beta V, \sigma Q) = CDH(V, Q)^{\beta\sigma}$ from Hash list \wedge_h so as to compute $CDH(U, V) = D^{-\alpha\beta}, CDH(V, Q) = (r_U \times r_S \times P)^{-\beta\sigma}$. Hence, the adversary \mathcal{A} can break CDH-problem. $Pr[AskH] \leq q_h \times (\frac{1}{2^{k_1 \cdot 2^{k_2} \cdot (n-1)}} + \frac{1}{2^{k_2 \cdot (n-1)}}) + \frac{1}{2^{k_2}} \times Succ_G^{cdh}(t')$,

Game G_5 : In this game, we consider the active attacks of the adversary via the Send query. Let the values of R_U, M_U and M_S be generated by adversary, This case is divided into three subcase:

Sbad1: R_U is generated by adversary as the input of the Send query, and is accepted. $Pr[Sbad1] \leq q_s \times Adv_{\mathcal{A},\mathcal{E}}^{cpa}(n')$.

Sbad2: The hash query has not been asked, but the authentications M_U and M_S are valid. $Pr[Sbad2] \leq 2q_s \times \frac{1}{2^t}$.

Sbad3: The adversary generates two valid authentications M_U and M_S by asking the hash oracle. $Pr[Sbad3] \leq q_s \times q_h \times (\frac{1}{2^{k_1 \cdot 2^{k_2} \cdot (n-1)}} + \frac{1}{2^{k_2 \cdot (n-1)}})$.

In addition, whatever the bit b involved in the Test-query, the answer is random, and independent for all the sessions, so we have $Pr[succ_5] = \frac{1}{2}$. Moreover, due to the $Corrupt(M, 1)$ query and $Corrupt(M, 2)$ query can be made at the same time. Hence, there is only one password for every transcript which can be tested by the adversary. In conclusion, we can get $\Delta_4 \leq q_s \times (\frac{1}{N} + Adv_{\mathcal{A},\mathcal{E}}^{cpa}(n') + \frac{1}{2^{t_1-1}} + q_h \times (\frac{1}{2^{k_1 \cdot 2^{k_2} \cdot (n-1)}} + \frac{1}{2^{k_2 \cdot (n-1)}}))$.

Theorem 2. Let G be a prepresent group. Our improved scheme can provide user anonymity in the random oracle model assuming \mathcal{E} is an IND-CPA security encryption in G .

Proof. In our improved scheme, the user's identity ID is protected in $R_U = \mathcal{E}_k(ID^*\|n_U)$. If the key $k = h_0(D\|A\|P_U)$ is safe, The adversary \mathcal{A} can make some $TestAnon$ queries about R_U . If the transcript of R_U with ID_0 or ID_1 is queried by \mathcal{A} , there must be $Adv_P^{anon}(\mathcal{A}) \leq Adv_{\mathcal{A},\mathcal{E}}^{cpa}(n')$, since \mathcal{E} is an IND-CPA security encryption. Next, we consider the security of k . Firstly, we compute $k = h'_0(D\|A\|P_U)$ using the private oracle h'_0 , where k is independent from D, A . Then, the probability that the adversary accurately guesses the value of k by Test query is $\frac{1}{2}$. Moreover the value of k using private h'_0 and the original k are indistinguishable unless the adversary uses $D\|A\|P_U$ to query the the original h_0 (denoted the event $AskH$). But $Pr[AskH] \leq q_h \times (Succ_G^{cdh}(t') + \frac{1}{(N-1)^2}) + \frac{q_s}{N-1}$. Hence k must be safe, the conclusion of theorem 2 is true.

Mean while, our improved scheme can provide user anonymity in the random oracle model. Furthermore, it is clear from our scheme that R_U is generated with

a random number n_U every time so that the adversary can not recognize the values of R_U in different authentication phase belong to the same user. Hence, our scheme also can satisfy the untraceability which Sun et al.'s scheme lacks.

6 Conclusion

In this paper, we point out the weakness of Sun et al.'s scheme, and propose an improved PAKE scheme with anonymity and untraceability using smart card for the network. In our improved scheme, there is only two flows in authentication phase, and no any interaction between U and S in password-change phase. In addition, we strictly prove the security and anonymity of our scheme. Hence, we think our proposed scheme is an efficient and provable secure PAKE scheme with robust anonymity.

Acknowledgment. This work was in part supported by Key Scientific and Technological Project of Zhengzhou City (No. 10PTGG341) and Key Scientific and Technological Project of Henan Province (No. 092101210502).

References

1. Fan, C., Chan, Y., Zhang, Z.: Robust remote authentication scheme with smart cards. *Comput. Secur.* 24(8), 619–628 (2005)
2. Juang, W.S., Chen, S.T., Liaw, H.T.: Robust and efficient password-authenticated key agreement using smart card. *IEEE Tran. Ind. Electron.* 55(6), 2551–2556 (2008)
3. Sun, D.-Z., Huai, J.-P., Sun, J.-Z., Li, J.-X., Zhang, J.-W., Feng, Z.-Y.: Improvements of Juang *et al.*'s Password-Authenticated Key Agreement Scheme Using Smart Cards. *IEEE Tran. Ind. Electron.* 56(6) (June 2009)
4. Li, X., Qiu, W., Zheng, D., Chen, K., Li, J.: Anonymity Enhancement on Robust and efficient Password-Authenticated Key Agreement Using Smart Cards. *IEEE Tran. Ind. Electron.* 57(2) (February 2009)
5. Shin, S., Kobara, K., Imai, H.: A Secure Threshold Anonymous Password-Authenticated Key Exchange Protocol. In: Miyaji, A., Kikuchi, H., Rannenberg, K. (eds.) *IWSEC 2007*. LNCS, vol. 4752, pp. 444–458. Springer, Heidelberg (2007)
6. Yang, J., Zhang, Z.: A New Anonymous Password-Based Authenticated Key Exchange Protocol. In: Chowdhury, D.R., Rijmen, V., Das, A. (eds.) *INDOCRYPT 2008*. LNCS, vol. 5365, pp. 200–212. Springer, Heidelberg (2008)
7. Shin, S., Kobara, K., Imai, H.: Very-Efficient Anonymous Password-Authenticated Key Exchange and Its Extensions. In: Bras-Amorós, M., Høholdt, T. (eds.) *AAECC 2009*. LNCS, vol. 5527, pp. 149–158. Springer, Heidelberg (2009)
8. Bellare, M., Pointcheval, D., Rogaway, P.: Authenticated Key Exchange Secure against Dictionary Attacks. In: Preneel, B. (ed.) *EUROCRYPT 2000*. LNCS, vol. 1807, pp. 139–155. Springer, Heidelberg (2000)
9. Zhou, T., Xiu, J.: Provable Secure Authentication Protocol with Anonymity for Roaming Services in Global Mobility Networks. *Computer Networks* 55(1), 205–213 (2011)

Study on QoS of Video Communication over VANET

Shouzhi Xu, Pengfei Guo, Bo Xu, and Huan Zhou

College of Computer and Information Technology, China Three Gorges University, Yichang, China

xsz@ctgu.edu.cn

Abstract. VANET (Vehicle Ad-Hoc Network) is an emerging hot technology, but it faces great challenge on quality of service issues since of limited transporting distance and high mobility. To evaluate the quality of video over VANET, a performance evaluation model and a simulation platform VANET-Evalvid is presented in this paper. The tool-set integrates myEvalvid, NS-2 and VanetMobiSim. The performance of different routing protocols under different network conditions is studied. Test results show that Position-based protocol is more suitable than Re-active protocol for video transmission over VANET, whereas Pro-active protocol plays bad performance. Further experiments discover that the quality of video transmitted in dense traffic scenario is better than it in sparse traffic scenario. Comparing with other research, the result testifies the correctness of our model and the efficiency of evaluation platform.

Keywords: VANET, QoS, Evaluation Model, Routing Protocol.

1 Introduction

VANET (Vehicular Ad-Hoc Network) is a kind of Mobile ad-hoc network to provide communications among vehicles and nearby fixed equipment, usually described as roadside equipment[1]. VANET can be used to solve not only traffic safety warning, but also traffic information inquiry, commercial advertisement and so on. Supporting video over VANET is an attractive feature for many VANET applications as follows:

(1) Driver assistance and safety applications. Adjacent vehicles can share traffic information video with each other, which are got from inter-vehicle equipment and roadside facilities. In case of a car accident in the distance, accident avoidance warnings could quickly notify drivers those conditions. Cooperative driving would allow vehicles to navigate without driver intervention by communicating with other vehicles about velocity, proximity, and other factors. Communication made to other vehicles prior to collision may allow the accident to be reconstructed more easily [2].

(2) Commercial and entertainment applications. Road side businesses, such as hotels and restaurants, can use content-rich video streams to broadcast advertisements to drivers on the road. Passengers in nearby cars can setup a video conversation by using the inter-vehicle streaming technology [3].

The main significance of VANET is providing safety and comfort for passengers, and giving assistance to drivers. However, in a VANET, quality of multimedia data

transmission is more difficult than other wireless networks, such as Ad Hoc and wireless sensor network. Many factors, such as high velocity, limited communication range, and dynamic network topology, may make data link disable[4]. The main goal of this paper is to analyze main the quality criteria among popular routing protocols.

Aiming at evaluating the quality of video transmitted over VANET, the rest of this paper is organized as: Section 2 introduces three typical QoS routing protocols, subsequently, two main evaluation models and a simulation platform is proposed for evaluating video transmission performance in Section 3, section 4 reports our experimental results and section 5 gives the concluding remarks finally.

2 Related Typical QoS Routing Protocols

There are three main classes of routing protocols in VANET: Proactive routing, Reactive routing, and Position-based routing [5]. Three typical routing protocols selected from each class respectively are analyzed, which include Destination Sequenced Distance Vector (DSDV)[6], ad-hoc on demand distance vector (AODV)[7], and Greedy Perimeter Stateless Routing (GPSR) [8] correspondingly.

(1) DSDV: DSDV is a table driven algorithm based on Bellmen-ford routing mechanism [6]. In this protocol, every mobile node maintains a routing table in which all of the possible destinations within the network and the number of hops to each destination are recorded. Each entry is assigned a sequence number by the destination node. The sequence numbers enable the mobile nodes to distinguish stale routes from new ones, thereby avoiding the formation of routing loops.

(2) AODV: AODV minimizes the number of required broadcast by creating routes on an on-demand basis to improve the DSDV. When a node wants to send a message to another node with an invalid route, it initiates a Path Discovery process to locate the destination node at first [7]. It broadcasts a route request (RREQ) packets to its neighbors, which then forward the request to their neighbors, and so on, until either the destination or an intermediate node with a newly route to the destination.

(3) GPSR: Each node in GPSR keeps states from immediate neighbors and uses only those states for data forwarding [8]. Forwarding nodes run greedy mode routing, which selects a node whose distance to a destination is shortest among all immediate neighbors and then drives the routing nodes to forward data to the destination node. If there is no neighbor whose distance to destination is greater than distance from forwarding node to destination, forwarding node runs perimeter mode routing.

At present, there are few works evaluating the quality of video over VANET in a realistic simulation environment [9]. Meng[10] proposes an architecture v3 to provide a live video streaming service through vehicle-to-vehicle (V2V) networks, but doesn't evaluate the quality of video transmitted using the video data. Fei [11] studies the performance of video streaming under different data forwarding and buffer management schemes, the dedicated short range communications (DSRC) are used in the simulation study, however, real vehicle mobility traces are not used in this study.

We introduce a simulation tool to study the quality of video transmitted under different traffic conditions and the above three typical routing protocols, which integrated Evalvid [12], NS-2 [13] and VanetMobiSim [14].

3 Evaluation Model and Test Bed

There are two approaches to support video transmitted over VANET: V2V approach and V2I (Vehicle to Infrastructure) approach. In this paper, we are interested in investigating the problem of supporting video transmitted over V2V, because it doesn't need any roadside infrastructure, it's easier to deploy. In this section, we mainly present an evaluation model, and design a simulating platform to evaluate performance of quality of multimedia under different routing protocols.

Many literatures have done some research on video transmission over VANET ([9]), authors in [8] propose a framework that facilitates the transmission of video over multi-hop wireless networks, and compare different routing protocols under different conditions. Since our goal is to evaluate the quality of video transmitted over VANET under different traffic conditions and different routing protocols, we mainly consider following characters:

- (1) Wireless channel quality can be easily affected by many factors, including street construction, road conditions, vehicle type and so on;
- (2) High dynamics of mobile nodes, which may incur frequent link disconnection and even network partition;
- (3) Application background is totally different from general Ad hoc network, while VANET is specially designed to achieve multi-hop communication between vehicle-to-vehicle and vehicle-to-Infrastructure on road.

So, we have to set a suitable evaluation model and then evaluate the quality of video transmitted in a realistic simulation environment[15]. In order to better match the reality, we use real video data and realistic vehicle mobility traces to evaluate the performance of video transmitted over VANET.

3.1 Evaluation Models

In this paper, we use two important performance parameters to evaluate the quality of video transmitted over VANET: Frame Loss Rate and PSNR.

(1) Frame Loss Ratio Model

The frame loss rate is a fundamental criterion of performance evaluation. In video streaming, a single video frame is decomposed into many smaller packets and sent into the network, so decoded video quality at the receiver is affected by two factors: encoder compression performance and distortion due to the packet loss or late arrivals. Based on the model in [13], the video distortion can be modeled as:

$$D_{dec} = D_{enc} + D_{loss} \quad (1)$$

The encoder distortion can be modeled by:

$$D_{enc} = D_0 + \theta / (R - R_0) \quad (2)$$

Where R is the rate of the video stream, and the parameters, D_0 and R_0 are estimated from empirical rate-distortion curves via regression techniques.

In this paper, we mainly care about the packets lost rate, represented by criterion D_{loss} . If the percentage of lost packets exceeds the bound of error correction, the receiver can not playback this frame. Similarly, if the packet arrival time is later than the playback deadline of the corresponding frame, it will also be dropped by the decoder in receiver cache. Therefore D_{loss} can be modeled by:

$$D_{loss} = P_{loss} + P_{delay} \tag{3}$$

In our study, we consider the Frame Loss Rate as the performance parameter based on (1). We use a video file "foreman.yuv" which has 400 frames, which contains 45 I frames, 89 P frames and 266 B frames. For each frame type, we measure the lost frames separately, and then we calculate the percentage of the overall frame loss. Overall Frame Loss Rate % = (Lost I frames + Lost P frames + Lost B frames) * 100 / Total Number of Frames.

(2) PSNR Model

Peak Signal-to-Noise Ratio (PSNR) is an important criterion to measure the error between the reconstructed image and the original one frame by frame. The PSNR has become the most widespread objective metric used to assess the application-level QoS of video transmissions [12].

For frame n, its PSNR between the source image S and destination image D is defined by PSNR model as:

$$PSNR(n)_{dB} = 20 \lg \left(\frac{V_{peak}}{\sqrt{\frac{1}{N_{col} N_{row}} \sum_{i=0}^{N_{col}} \sum_{j=0}^{N_{row}} [Y_S(n, i, j) - Y_D(n, i, j)]^2}} \right) \tag{4}$$

Where Y denotes the luminance component, $V_{peak} = 2^k - 1$ and k = number of bits per pixel (luminance component). $Y_S(n, i, j)$ and $Y_D(n, i, j)$ are the values of the luminance component of the nth frame at pixel for the source and destination images respectively. N_{col} and N_{row} are the dimensions of the frame.

3.2 Test Bed

To evaluate the performance of video transmitted over VANET, we propose a simulation tool-set called VANET-EvalVid. This tool-set integrates myEvalvid[12], and VanetMobiSim. MyEvalvid is a tool-set for evaluating the quality of video transmitted over a real or simulated communication network, which combines Evalvid and NS-2. VanetMobiSim is a tool for generating realistic vehicle mobility trace file. With the integration, researchers can easily evaluate the quality of video transmitted over VANET and analyze designed mechanisms, such as network protocols or QoS control schemes in a realistic simulation environment.

The framework of VANET-EvalVid is shown in fig. 1.

The VANET-EvalVid framework includes 6 main components as follows:

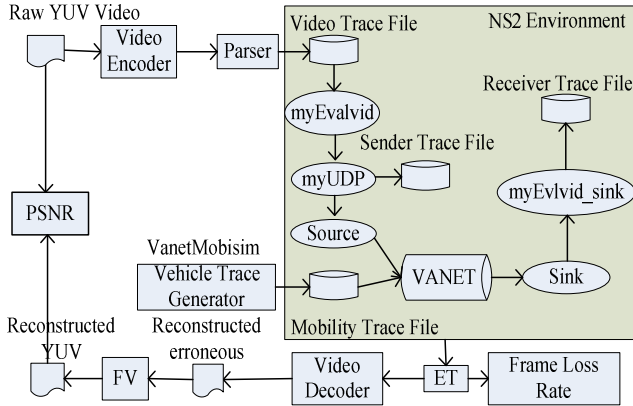


Fig. 1. Framework of VANET-Evalvid

Video Source: The video source is responsible for generating video streams which video format can be either in YUV CIF (352 x 288) or YUV QCIF (176*144).

Video Encoder and Video Decoder: These coders are used to convert the video file from YUV format to MPEG4 format at the sender side and transfer it back to YUV format at the receiver side.

Fix Video (FV): Since video frames are compared frame by frame, the total number of video frames at the receiver side must match that of the original video at the sender side. If there are missing frames that video code cannot be decoded, the FV is responsible for fixing them by substituting the last successfully decoded frame for each lost frame.

Evaluate Trace (ET): The ET component is responsible for the evaluation task. The evaluation begins at the sender side when the video transmission finishes. Information for delivering video packets is transported back to the sender side. The ET generates a report of video frame/packet delay, frame/packet loss, and frame/packet jitter by comparing the trace files, including the original encoded video file, the video trace file, the sender trace file and the receiver trace file. In addition, the ET also creates a reconstructed video file, corresponding to the possible corrupted video frame at the receiver side.

VanetMobiSim: This component is responsible for generating vehicle mobility trace file used as the mobility trace file in NS-2. It is designed for generating vehicle mobility trace file, which supports both macro-mobility and micro-mobility representation to defined mobility models. The macro-mobility models is mainly used to define the characteristic of the roads, as lines, speed limit, traffic signs, etc.. The micro-mobility models include Intelligent Driver Model with Intersection Management (IDM-IM) and Intelligent Driver Model with Lane Changes (IDM-LC). IDM-IM is used to define the behavior when the driver is at the intersection, while IDM-LC is used to define the behavior when the driver is changing lane.

The communication interfaces: The interfaces between myEvalvid and NS-2 are given by MyTrafficTrace MyUDP and MyUDPSink in our test bed.

MyTrafficTrace is response to read the video trace file and extract the video frame type and the video frame size. Furthermore, MyTrafficTrace fragments the video frames into smaller video packets and sends these packets to the lower layer at the appropriate time according to the time settings in the simulation script file.

MyUDP enhances the original UDP component in NS-2. This interface allows users to designate the output file name of the sender trace file; moreover, it also records the information of transmitted video packets, like the packet id, the timestamp, and the payload size.

MyUDPSink is a receiving agent for video packets sent by MyUDP. When receiving a video packet, it records information from the transmitted video packet, such as the packet id, the timestamp, and the payload size.

4 Simulation and Analysis

4.1 Simulation Scenario Setup

We setup a scenario of transmitting real video data over a 2km road with a bi-directional multiple lanes, which is generated by VanetMobiSim. Vehicles on the road can react according to the vehicles ahead, like changing lanes. Two different traffic environments is constructed[16], which are sparse (150 cars distribute on the road) and dense (500 cars distribute on the road). Then we use the simulation tool-set introduced in section 3 to evaluate the quality of video transmitted under three different routing protocols. Main simulation parameters list in Table 1.

Table 1. Simulation Parameters

Parameter	Discription
Simulator	NS-2.28
Routing Protocols	DSDV, AODV, GPSR
Simulation Time	100s
Video File	Foreman.yuv
Communication Range	300m
Bandwidth	6Mbps
MAC Protocol Packet Size	IEEE 802.11DCF 1024 Bytes

The video file is foreman.yuv containing 400 frames, which average PSNR is 34.89. According to the suggestion of DSRC [17], we set the bandwidth 6Mbps and the communication range 300m. In our experiments, we mainly compare three different routing protocols introduced in section 2: DSDV, AODV and GPSR.

Video data is transmitted over a sparse traffic environment and a dense traffic environment separately to test performance in different scenarios. From the VANET-Evalvid model, the simulation test works as:

- Setting up a certain traffic scenario to simulate transport status;
- Selecting different embeded standard routing protocols as network routing;
- Transmitting video data over the slected routing protocol;

Recording network data and making a comparison analysis.

In each test, by changing the speed of vehicles, the quality of video transmitted under different routing protocols can be measured. The video starts at 0.0 second and the simulation lasts 100 seconds.

4.2 Simulation Results and Analysis

(1) In sparse traffic scenario

In the sparse scenario, 150 cars distribute on a 2km road. Fig.2 shows the performance comparison of frame lost rate and average PSNR among different protocols. From fig.2(a), the Frame Loss Rate of DSDV increase greatly with the increase in vehicle speed, while GPSR protocol plays best. The network is more likely to disconnect with the increase in vehicle speed, so the Frame Lost Ratio of different protocols is the biggest when the speed of cars is at 30m/s.

Fig.2(b) shows the average PSNR in different protocols are all greater than 32dB. It indicates that the quality of the received video is at the excellent level. In addition, the Frame Loss Rate of DSDV protocol is so large that the video can not be reconstructed, so we can't get the simulation results of the average PSNR of DSDV protocol when the speed of cars is 20m/s and 30m/s.

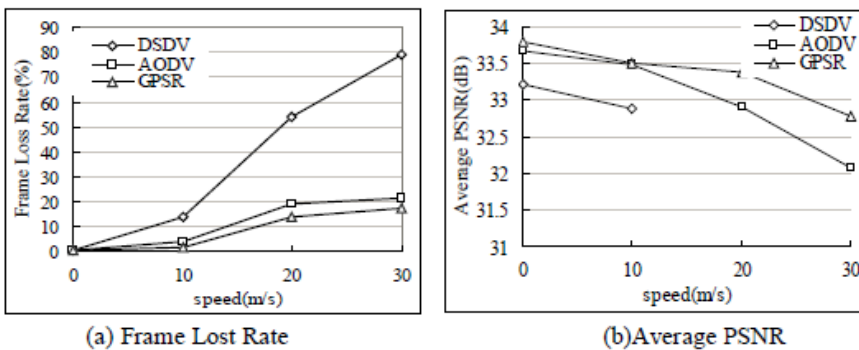


Fig. 2. Performance comparison in sparse traffic scenario

(2) In dense traffic scenario

In this case, 500 cars distribute on a 2km road. Fig. 3 shows the result of performance comparison among different protocols. From fig. 3(a), DSDV protocol plays worst, and GPSR protocol plays best too, with the increase in vehicle speed. The simulation result shows that the Position-based protocol can handle the video transmission over high mobility scenario well. The main reason is that the position-based protocol does not rely on the maintained neighboring information which is likely inaccurate in high mobility scenario. From fig. 3(b), another conclusion is that the frame loss rate is obviously less, but the average PSNR is also larger in the dense scenario. The main reason of it is that the network is highly connected in the dense scenario, thus the quality of video transmitted in dense traffic scenario is better than that in sparse traffic scenario. This result is quite similar with other research, and it testifies our model correct and efficiency of evaluation.

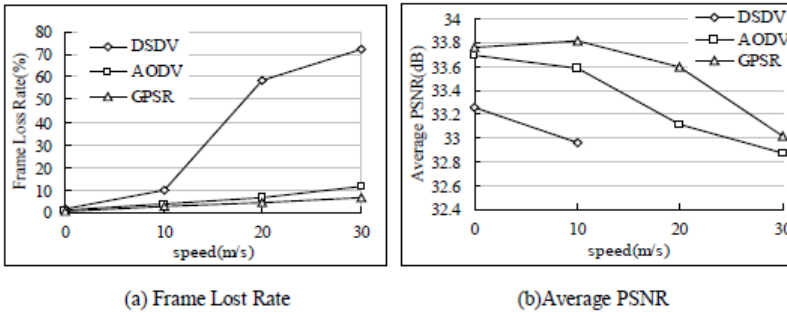


Fig. 3. Performance comparison in dense traffic scenario

From the above analysis, protocols of three classes play quite different in VANET. Proactive protocols such as DSDV have considerable difficulties in maintaining valid routes, and lose many packets easily. With increasing mobility, it strives to continuously maintain routes to every node, which results in increasing network load. The rapidly changing routes through the fast vehicle nodes add inter-group traffic fastly. So such protocol cannot adapt well to such fast route changes. In reactive routing protocols like AODV, uncontrolled flooding generates redundant transmissions, which may cause broadcast storm problem, and the network suffers from the increasing administrative load as the number of vehicle nodes increases. Whereas Position-based protocols like GPSR play well in VANET in highway, the reason is the forwarding node always tries to find a best node closed to pre-calculated routing node, so it simplifies the check if the destination node is in its neighborhood.

5 Conclusion and Future Work

In this paper, we presented two main evaluation models of QoS performance or routing protocols over VANET and designed a simulation platform of VANET. With the integration of simulation platform, researchers can easily evaluate the quality of video transmitted over VANET and analyze their designed mechanisms. We created two different scenarios to test the performance of multimedia data transmission under different routing protocols. From the result analysis, Position-based protocol is more suitable for video transmission over VANET than Re-active protocol, whereas Proactive protocol is not suitable. It also illustrates scalability of different protocols. In the future work, we'll compare different Position-based protocols to choose a better routing protocol for video transmission over VANET. The result testifies the correctness of our model correct and the efficiency of evaluation platform.

Acknowledgments. This work was partially supported by National Natural Science Foundation of China (61174177), Yichang STF(A2011-302-13) and CTGU 2011CX051.

References

1. Su, X.: A comparative survey of routing protocol for vehicular sensor networks. In: 2010 IEEE International Conference on Wireless Communications Networking and Information Security (WCNIS), June 25-27, pp. 311–316 (2010)
2. Fracchia, R., Meo, M.: Alert service in VANET: Analysis and design. In: 2006 4th International Symposium on Modeling and Optimization in Mobile, Ad Hoc and Wireless Networks, April 03-06, pp. 1–8 (2006)
3. Toor, Y., Muhlethaler, P., Laouti, A.: Vehicle Ad Hoc networks: applications and related technical issues. *IEEE Communications Surveys & Tutorials* 10(3), 74–88 (2008)
4. Mittal, M.: A Study of Live Video Streaming over Highway Vehicular Ad hoc Networks. *International Journal of Computer Applications* 1(21), 86–90 (2010)
5. Tomer, P., Chandra, M.: An application of routing protocols for Vehicular Ad-hoc Networks. In: 2010 International Conference on Networking and Information Technology (ICNIT), June 11-12, pp. 157–160 (2010)
6. Perkins, C., Bhagwat, P.: Highly Dynamic Destination Sequenced Distance-vector Routing (DSDV) for Mobile Computers. *ACM SIGCOMM Computer Communication Review* 24(4), 234–244 (1994)
7. Perkins, C., Royer, E.: Ad-hoc on-demand distance vector routing. In: Proc. of IEEE WMCSA 1999, New Orleans, Louisiana, pp. 90–100 (February 1999)
8. Karp, B., Kung, H.T.: GPSR: Greedy perimeter stateless routing for wireless network. In: Proc. of MobiCom 2000, Boston, MA, USA, pp. 243–254 (August 2000)
9. Mao, S., et al.: Video Transport over Ad Hoc Networks: Multistream Coding with Multipath Transport. *IEEE JSAC* 21(10), 1721–1737 (2003)
10. Guo, M., Ammar, M.H., Zegura, E.W.: V3: A vehicle-to-vehicle live video streaming architecture. In: Proc. of the 3rd Int'l Conf. on Pervasive Computing and Communications, pp. 171–180 (March 2005)
11. Xie, F., Hua, K.A., Wang, W., Ho, Y.H.: Performance study of live video streaming over highway vehicular Ad hoc networks. In: Proc. of the 1st IEEE International Symposium on Wireless Vehicular Communications (WiVeC), Baltimore, MD, pp. 2121–2125 (September 2007)
12. Klaue, J., Rathke, B., Wolisz, A.: EvalVid - A Framework for Video Transmission and Quality Evaluation. In: Proc. of the 13th Int.l Conference on Modeling Techniques and Tools for Computer Performance Evaluation, Illinois, USA, pp. 255–272 (September 2003)
13. The Network Simulator – NS-2, <http://www.isi.edu/nsnam/ns/>
14. VanetMobiSim, <http://vanet.eurecom.fr>
15. Ke, C.-H., Lin, C.-H., Shieh, C.-K., Hwang, W.-S.: A Novel Realistic Simulation Tool for Video Transmission over Wireless Network. In: Proc. of SUTC, Taichung, Taiwan, June 5-7 (2006)
16. Wang, W., Xie, F., Chatterjee, M.: Small-Scale and Large-Scale Routing in Vehicular Ad Hoc Networks. *IEEE Transactions on Vehicular Technology* 58(9), 5200–5213 (2009)
17. DSRC: Dedicated Short Range Communications (DSRC) Home [EB/OL], <http://www.leearmstrong.com/Dsrc/DSRCHomeset.html>

Phase Noise Estimation and Mitigation for Cognitive OFDM Systems

Yuan Jing, Haoyu Li, Xiaofeng Yang, Li Ma, Ji Ma, and Bin Niu

College of Information, Liaoning University, Shenyang 110036, Liaoning, P. R. China
eejingyuan@yahoo.com.cn

Abstract. In this paper, a novel maximum likelihood (ML) method is proposed to estimate and mitigate phase noise (PN) for cognitive orthogonal frequency division multiplexing (OFDM) systems. In the proposed method, PN estimation is formulated as a unitary-constrained optimization problem based on the ML criterion. Using the obtained estimate of the PN vector, both common phase error and inter-carrier interference are effectively mitigated and reduced for symbol error rate (SER) performance improvement. Simulation results show that the proposed method can mitigate PN effectively, and obtain better SER performance for cognitive OFDM systems compared with conventional methods.

Keywords: Orthogonal frequency division multiplexing (OFDM), phase noise, maximum likelihood (ML), inter-carrier interference (ICI).

1 Introduction

Orthogonal frequency division multiplexing (OFDM) is an attractive multi-carrier modulation technique. Since its high spectrum efficiency, OFDM has been widely adopted both in many existed wireless communication standards and the novel cognitive radio system [1-3]. Unfortunately, OFDM is highly sensitive to the phase noise (PN) mainly introduced by the mismatch between the phase of carrier signal and the phase of the local oscillator. The effect of PN on OFDM cognitive systems includes an inter-carrier interference (ICI) term and a common phase error (CPE) term, which generally leads to a rotation of all constellation points and damages the orthogonality between sub-carriers. Therefore, the symbol error rate (SER) performance of cognitive OFDM systems may significantly degrade under the PN.

Many methods [4-9] have been presented to estimate or mitigate PN for OFDM systems. Some methods such as [4] are mainly focused on estimating and correcting the CPE term, and the ICI term is proposed as an additive noise in these methods. However, they perform well only for small PN levels. In [5], an efficient least-squares (LS) based phase-noise suppression (PNS) algorithm is proposed. Since the minimum mean-square error (MMSE) equalization technique is used to reduce the effect of ICI term, the obtained SER performance in [5] is better than that in [4]. Recently, in [6], a two-stage PN compensation method which can effectively mitigate the ICI is also proposed. In the first stage, channel coefficients are estimated by using block-type

pilot symbols; Then, the data symbols and PN are jointly estimated in the second stage. In [7], two novel PN estimation and mitigation methods have been proposed for OFDM systems. These two PN estimation methods are based on ML estimation criterion and linear MMSE (LMMSE) technique, respectively. Both these two methods perform well with the known channel knowledge. In addition, the statistics of PN and background noise are also required in the deviation of the LMMSE-based method in [7]. In addition, [8-9] also present the new results of PN mitigation.

In this paper, a novel ML-based PN estimator is proposed for cognitive OFDM systems. Through utilizing the unitary character of PN vector in OFDM system model, the PN estimation is formulated as a unitary-constrained optimization problem based on the ML criterion. Then the unitary-constrained optimization technique in [10] is utilized in this paper to estimate the PN vector. Simulation results show that, the proposed method can effectively mitigate the PN and achieve a good SER performance for cognitive OFDM system.

2 System Model

In an OFDM system, source data are mapped into the complex signal constellation points at the transmitter. The N constellation points $X_m, (m=0, \dots, N-1)$ are grouped and modulated onto N parallel sub-carriers by an inverse discrete Fourier transform (IDFT): $x_n = (1/\sqrt{N}) \sum_{m=0}^{N-1} X_m \exp(j2\pi nm/N)$, where $n=0, \dots, N-1$. Before transmission, a cyclic prefix (CP) is inserted between successive OFDM symbols to mitigate the inter-symbol interference.

Assuming perfect frequency and timing synchronization at the receiver, after removing the CP and performing OFDM demodulation, the received $(n+1)$ th subcarrier signal y_n can be written as [6,7]

$$y_n = x_n H_n I_0 + \sum_{k=0, k \neq n}^{N-1} x_k H_k I_{k-n} + z_n \quad (1)$$

where H_n is the channel frequency response on the $(n+1)$ th subcarrier, z_n denotes the additive white Gaussian noise (AWGN) with variance σ^2 . I_0 is the common phase error (CPE), and the second term on the right side of (1) denotes the inter-carrier interference (ICI). Let $\phi_n, (n=0, \dots, N-1)$ denote the phase noise (PN) in current OFDM symbol, we can obtain that $I_k = (1/N) \sum_{n=0}^{N-1} \exp(j2\pi kn/N + j\phi_n)$.

Let $W[n, m] = (1/\sqrt{N}) \exp(j2\pi mn/N)$ denote the $(n+1, m+1)$ th element of the $N \times N$ matrix W . The vector form of (1) can be written as

$$Y = WPW^H XW_L h + Z \quad (2)$$

where $Y = [y_0, \dots, y_{N-1}]^T$, $Z = [z_0, \dots, z_{N-1}]^T$, and $(\cdot)^T$ denotes transpose. h is the $L \times 1$ channel impulse response (CIR) vector, where L is the channel length. W_L is the

$N \times L$ sub-matrix of W . Diagonal matrices P and X are defined as $P = \text{diag}\{\exp(j\phi_0), \dots, \exp(j\phi_{N-1})\}$ and $X = \text{diag}\{X_0, \dots, X_{N-1}\}$, respectively.

3 Constrained ML Estimation

As shown in (2), if both P and h are known, data matrix X can be easily estimated. Due to the nonlinearity of the exponential term in P , we try to estimate the PN vector $\Phi = [\exp(-j\phi_0), \dots, \exp(-j\phi_{N-1})]^T$ instead of $\phi_n, (n = 0, \dots, N-1)$ in this section.

Assuming X is known, from (2), the probability density function of Y , given Φ and h , may be written as

$$p(Y | \Phi, h) = \frac{1}{(\pi\sigma^2)^N} \exp\left(-\frac{\|Y - WPW^H XW_L h\|^2}{\sigma^2}\right) \tag{3}$$

where $P = \text{diag}\{\exp(j\phi_0), \dots, \exp(j\phi_{N-1})\}$. Since $\Phi = [\exp(-j\phi_0), \dots, \exp(-j\phi_{N-1})]^T$, we can obtain

$$\Phi^H \Phi = \text{tr}\{P^H P\} = N \tag{4}$$

where $\text{tr}\{\cdot\}$ is the trace operator.

Therefore, from (3), the joint ML estimation problem of Φ and h can be given as

$$[\Phi, h] = \arg \min_{\Phi, h} J(\Phi, h) \quad \text{Subject to } \Phi^H \Phi = N \tag{5}$$

where $J(\Phi, h) = \|Y - WPW_H XW_L h\|^2$.

Given the PN vector Φ , the ML estimate of h is given by

$$\hat{h} = (W_L^H X^H XW_L)^{-1} \bar{W}(\Phi) Y \tag{6}$$

where $\bar{W}(\Phi) = W_L^H X^H WP^H W^H$. Substituting (6) to (5), the ML estimate of Φ is

$$\hat{\Phi} = \arg \min_{\Phi} J(\Phi) \quad \text{Subject to } \Phi^H \Phi = N \tag{7}$$

Where

$$J(\Phi) = \|Y - WPW^H XW_L \hat{h}\|^2 = \|WPP^H W^H Y - \bar{W}^H(\Phi)(W_L^H X^H XW_L)^{-1} \bar{W}(\Phi) Y\|^2$$

Let $S = I_N - W^H XW_L (W_L^H X^H XW_L)^{-1} W_L^H X^H W$, and $N \times 1$ vector $y = W^H Y$ denote the time-domain received signal vector, then (7) may be rewritten as

$$\hat{\Phi} = \arg \min_{\Phi} \|WPS P^H W^H Y\|^2 \tag{8}$$

where $A = D_y^H S^H S D_y = D_y^H S D_y$, $D_y = \text{diag}\{y\}$.

As shown in (8), the estimation of Φ can be formulated as the following constrained optimization problem

$$\hat{\Phi} = \arg \min_{\Phi} J(\Phi) = \Phi^H \mathbf{A} \Phi, \quad \text{subject to } \Phi^H \Phi = N \quad (9)$$

Let $\Theta = \Phi / \sqrt{N}$, (9) is equivalent to the following unitary-constrained optimization problem

$$\hat{\Theta} = \arg \min_{\Theta} J(\Theta) = \frac{1}{2} \Theta^H \bar{\mathbf{A}} \Theta, \quad \text{subject to } \Theta^H \Theta = 1 \quad (10)$$

where $\bar{\mathbf{A}} = 2NA$. Then the PN vector estimate $\hat{\Phi} = \sqrt{N}\hat{\Theta}$.

The unitary constrained optimization problem like (10) has been deeply investigated in [8]. Through reformulating the constrained optimization problem as an unconstrained one on the manifold, we can use the iterative algorithm (such as steepest-descent like method) to converge to the local minimum of the cost function $J(\Theta)$. Then, the value $\hat{\Theta}$ corresponding to the minimum of $J(\Theta)$ is the desired estimate of the Θ .

In (8), since $A = D_y^H S^H S D_y = A^H$, the matrix $\bar{\mathbf{A}}$ in the cost function $J(\Theta)$ is an Hermitian matrix. Therefore, when $J(\Theta)$ achieves its minimum with the constrain $\Theta^H \Theta = 1$, Θ is actually the eigenvector corresponding to the smallest eigenvalue of $\bar{\mathbf{A}}$ [8]. From [8], we can obtain the following algorithm to solve the constrained optimization problem (10).

- 1) Choose an initial $\hat{\Theta} = 1 / \sqrt{N} \hat{\Phi}$.
- 2) Set $\gamma = \hat{\Theta}^H \bar{\mathbf{A}} \hat{\Theta}$, $a_0 = \hat{\Theta}^H \bar{\mathbf{A}}^3 \hat{\Theta} - \gamma^2$, $a_1 = \hat{\Theta}^H \bar{\mathbf{A}}^3 \hat{\Theta} - 3a_0\gamma - \gamma^3$ and $a_2 = -a_0$, then compute the convergence step μ , which is the positive root of the following equation:

$$a_0^2 \mu^2 + a_1 \mu + a_2 = 0 \quad (11)$$

- 3) Compute the deepest decent direction $d = \gamma \hat{\Theta} - \bar{\mathbf{A}} \hat{\Theta}$. If $d^H d > \varepsilon$ where ε is a small value, then stop iteration; else, update $\hat{\Theta} = \hat{\Theta} + \mu d$ and $\hat{\Theta} = \hat{\Theta} / \sqrt{\hat{\Theta}^H \hat{\Theta}}$, then go to step 2.

- 4) Obtain the PN vector estimate $\hat{\Phi} = \sqrt{N} \hat{\Theta}$.

It should be noted that the channel knowledge is not required in our PN estimation method. It also means that, unlike some conventional methods in [4], the estimation performance of the proposed method does not depend on the channel estimation errors.

4 Equalization and Decoding

At the receiver, the transmitted data symbols are distorted by CPE and ICI due to the PN. To effectively recover the data symbols, equalization is performed to mitigate PN

by pre-multiplying the frequency-domain received signal Y in (2) with $W\hat{P}^H W^H$ where $\hat{P}^H = \text{diag}\{\hat{\Phi}\}$. Then the PN corrected received signal is given by

$$Y = W\hat{P}^H W^H Y = W\hat{P}^H y \quad (12)$$

After equalization, the estimate \hat{X}_n of data symbol on the n th subcarrier can be obtained by the hard decision: $\hat{X}_n = \hat{y}_n / H_n$, where \hat{y}_n and H_n are the $(n+1)$ th element of \hat{Y} and the channel frequency response H , respectively.

Note that, like the conventional ML method in [6], the data symbol is required and used to calculate the matrix \bar{A} in (10) in the proposed PN estimation method. As mentioned in [6], however, we can obtain an initial estimate of data symbol by using the PNS method in [5]. Then the obtained initial data estimate may be used in the proposed method to estimate the PN, thus improving the system performance. Moreover, to speed up the convergence of the proposed method, the estimated CPE using the conventional method may also be used as the initial value $\hat{\Phi}_i^* = \hat{I}_0 / \sqrt{\hat{I}_0^* \hat{I}_0}$, where \hat{I}_0 is the CPE estimate.

It should be also noted that, since it requires the inversion of $W_L^H X^H X W_L$ when calculating the matrix A in our method, the computational complexity of our method is higher than the ML method in [6] whose computational complexity is perfectly reduced, even though the performance of our method is better than ML method.

5 Simulation Results

In this section, a single-antenna OFDM system with 64 subcarriers is constructed in our simulation. To evaluate the PN mitigation performance of the proposed method, channel coding is not used in our simulations. Since the M-QAM modulation signals are more sensitive to the PN than M-PSK modulation signals, in our simulations, 16 QAM modulation is employed in the constructed OFDM system. In addition, each OFDM symbol is transmitted through a 6-taps Rayleigh fading frequency-selective channel which is assumed to be constant within each OFDM frame (16 OFDM symbols) but varying frame by frame. We also assume that the length of cyclic prefix is greater than the channel length. The range of the considered signal-noise-ratio (SNR) is between 0 dB and 35 dB. In our simulations, the PN is generated by using the Wiener PN model, which has been described in [4].

For performance comparison, CPE estimation algorithm in [4] and ML estimator in [6] are also realized. The symbol error rate (SER) performance comparisons of these three methods are shown in Fig. 1. And the PN variance equals to 10^{-2} in this simulation. When there is no PN correction, as shown in Fig. 1, the receiver can not effectively detect the transmitted data symbols. For conventional CPEC method, due to neglecting the ICI, there is an error floor when the SNR is relative high. Since both the CPE and the ICI are considered in ML estimator in [6], the SER performance is better than that of the conventional CPEC method. Furthermore, since the unitary constrain is considered, we can also observe that the proposed constrained ML method is best for all the SNRs.

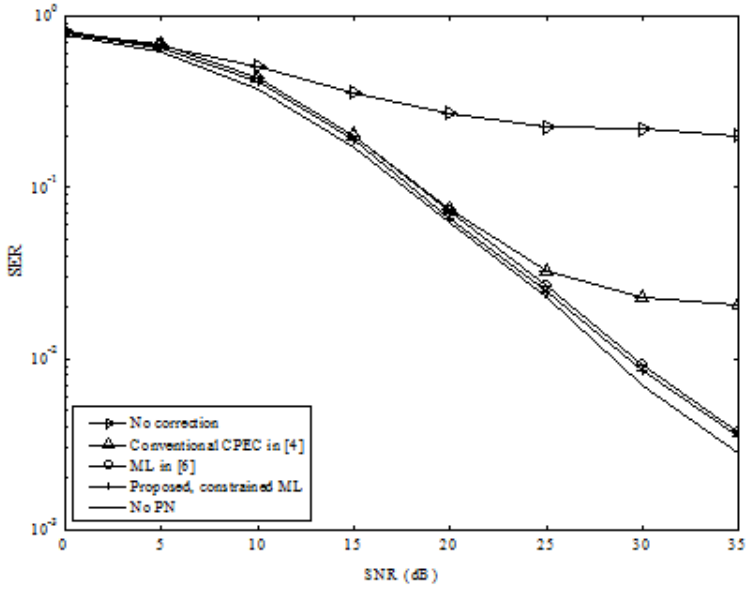


Fig. 1. Symbol error rate performance comparison between the proposed method and the conventional methods

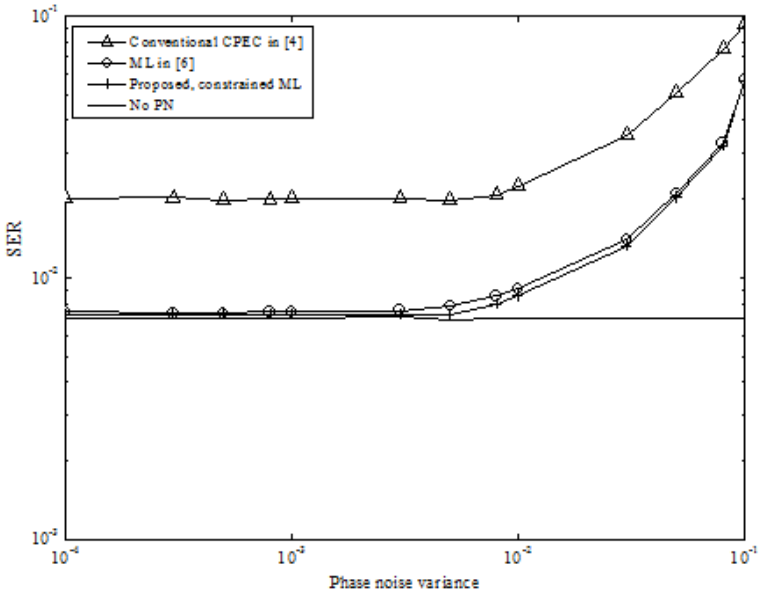


Fig. 2. Symbol error rate performance versus phase noise variance

In Fig. 2, we also exhibit the performances of these three methods versus different PN variances, when the SNR is 30 dB. It is shown that, like in [6], the performance of

the conventional CPE method is worst in the proposal estimators. In addition, the proposed constrained ML method outperforms the conventional ML method in [6] when the PN variance is smaller than 10^{-1} .

6 Conclusions

In this paper, a novel maximum likelihood (ML) phase noise (PN) estimator is proposed for cognitive OFDM systems. The PN estimation problem is first formulated as a ML estimation problem with unitary constraint. Accordingly, PN may be estimated and mitigated by solving this constrained ML optimization problem. Simulation results show that, the proposed method can effectively mitigate the PN for the cognitive OFDM system, thus obtaining a good SER performance.

Acknowledgments. The work of this paper is supported by the National Natural Science Foundation of China (No. 61101115).

References

1. Haykin, S.: Cognitive Radio: Brain-empowered Wireless Communications. *IEEE J. Sel. Areas Commun.* 23, 201–220 (2005)
2. Sun, D., Zheng, B.: A Novel Sub-carrier and Power Joint Allocation Algorithm for Multi-user Cognitive OFDM. In: 10th IEEE International Conference on Signal Processing, pp. 1458–1462. IEEE Press, Nanjing (2010)
3. Zhang, H., Le Ruyet, D., Terre, D.: Spectral Efficiency Analysis in OFDM and OFDM/OQAM Based Cognitive Radio Networks. In: 69th IEEE International Conference on Vehicular Technology, pp. 1–5. IEEE Press, Barcelona (2009)
4. Robertson, P., Kaiser, S.: Analysis of the Effects of Phase Noise in Orthogonal Frequency Division Multiplexing (OFDM) Systems. In: 1995 IEEE International Conference on Communications, pp. 1652–1657. IEEE Press, Seattle (1995)
5. Wu, S., Bar-Ness, Y.: A Phase Noise Suppression Algorithm for OFDM Based WLANs. *IEEE Commun. Lett.* 6, 535–537 (2002)
6. Wu, S., Liu, P., Bar-Ness, Y.: Phase Noise Estimation and Mitigation for OFDM Systems. *IEEE Trans. Wireless Commun.* 5, 3616–3625 (2006)
7. Zou, Q.Y., Tarighat, A., Sayed, A.H.: Compensation of Phase Noise in OFDM Wireless Systems. *IEEE Trans. Signal Process.* 55, 5407–5424 (2007)
8. Riihonen, T., Tchamov, N., Werner, S., Valkama, M., Wichman, R.: Characterization of OFDM Radio Link Under PLL-Based Oscillator Phase Noise and Multipath Fading Channel. *IEEE Trans. Commun.* 99, 1–8 (2012)
9. Lee, M.K., Lim, S.C., Yang, K.: Blind Compensation for Phase Noise in OFDM Systems over Constant Modulus Modulation. *IEEE Trans. Commun.* 60, 620–625 (2012)
10. Manton, J.H.: Optimization Algorithms Exploiting Unitary Constrains. *IEEE Trans. Signal Process.* 50, 635–650 (2002)

An Effective Multipath Approach to Reducing Congestion in WSNs

Laomo Zhang¹, Ying Ma¹, and Guodong Wang²

¹School of Software, Henan Institute of Engineering, 450053, Zhengzhou, China

²Graduate University, Chinese Academy of Sciences, 100049, Beijing, China
{cyechna, myfirst}@163.com, wgdaaa@126.com

Abstract. In recent years, many researchers have devoted themselves to investigate the issues of network congestion in the WSNs. In order to resolve the congestion in the WSN quickly, we propose an effective multipath approach to alleviate network congestion. The proposed approach establishes multiple paths before congestion occurs so that traffic can be rerouted to another path when congestion is detected. Our simulation results show that the approach can effectively alleviate congestion and provide better protection ability with less extra nodes and links than other braided multipath approaches.

Keywords: WSN, Congestion, Multipath, Base Station.

1 Introduction

The Wireless Sensor Network (WSN) [1-3] is an emerging technology that can be used as the platform for the versatile applications. In [4], Kang et al. proposed a Topology-Aware Resource Adaptation (TARA) approach, which can build a new topology with enough capacity to handle the increased traffic by making use of the capacity analysis model.

In this paper, we proposed an effective multipath approach to alleviate congestion in a WSN. The effective multipath approach can establish multiple paths for an active route before congestion occurs and would detour the data packets on the original path to another path when congestion is detected. Thus, constructing an effective multipath before congestion development can quickly react to congestion. Moreover, the effective multipath approach can support higher protection ability for the original path. It can find the alternate path quickly and isolate the failure effectively when congestion occurs. The rest of this paper is organized as follows. Section 2 briefly introduces the related work. Section 3 presents the proposed effective multipath approach. Thereafter, Section 4 discusses the simulation study. Finally, we conclude this paper in Section 5.

2 Related Work

In this section, we briefly describe the constraints of WSNs. Then, we present the TARA approach and the multipath approaches.

2.1 Constraints of WSNs

The WSN has some strict constraints on traffic pattern, energy, bandwidth, buffer size, memory, processing capability, etc. We briefly describe these constraints as follows. (1)Unbalanced Traffic and Data Redundancy. (2)Energy Limitation. (3)Bandwidth Limitation. (4)Buffer Size Limitation.[5]

2.2 TARA Approach

TARA approach is based on a resource control scheme, which uses a capacity analysis model to determine the needed topology to increase capacity and accommodate high incoming traffic during the emergent state. Moreover, it can detour the path to isolate the failure when congestion in intersection hotspot is detected.[6]

2.3 Multipath Approaches

The TARA approach establishes the detour path and delivers the data via detour path toward the sink when congestion is detected. However, during the establishment of the detour path, the congestion could result in the packet loss and lengthen transmission delay. Thus, to build alternate paths before congestion occurs is a feasible approach to avoid packet loss.[7]Some multipath approaches were proposed in the literature. We briefly discuss these approaches below.

(1) Disjoint Multipath (DM) Approach: In [8], a Split Multipath Routing (SMR) protocol is proposed. Although the disjoint multipath approach has some attractive resilience properties [9], it can be energy inefficient because the alternate node-disjoint path is longer than the primary path and, thus, consumes significantly much energy.(2) Braided Multipath (BM) Approach: In [9], a braided multipath routing is proposed. (3) Disjoint and Braided Multipath (DBM) Approach: Hoang et al. proposed a disjoint and braided multipath (DBM) approach in [10].

3 Effective Multipath Approach

3.1 Topology of the Effective Multipath Approach

The proposed Effective Multipath approach is illustrated in Fig.1. An extra link is established on the alternate path between node a_1 and node a_3 . Note that node n_2 and node a_1 in the topology still introduce the key node problem. In other words, if the two nodes fail simultaneously, the approach fails to resolve node failure due to congestion. However, it can provide more paths from source to sink to improve the tolerance to node failure and can improve the protection ability on primary path with less cost. Moreover, if the hop count of the primary path increases, it can support more alternate paths from source to sink.

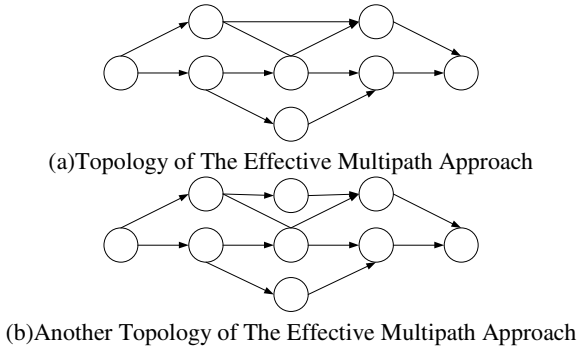


Fig. 1. The Proposed Effective Multipath Approach

3.2 Analysis of the Effective Multipath Approach

In order to measure and compare the performance of the BM, DBM, and approaches we define three performance metrics as follows.

(1) Average Cost to Protect a Node/Link: The average cost to protect a node or a link for the effective multipath approach is defined as follow.

$$C = \frac{n_e + l_e}{\sum_{i=1}^a n_i^p + l_i^p} \tag{1}$$

Where n_e and l_e are the numbers of extra nodes and extra links, respectively, and a is the number of alternate paths and n_i^p and l_i^p are the numbers of nodes and links protected by the i -th alternate path. With an example shown in Fig. 1(a), the number of extra nodes and links is 10 and n_i^p and l_i^p for the five alternate paths are 1+2, 1+2, 2+4, 1+2 and 3+4, respectively. Thus, the cost for this case is 0.4545. Obviously, the smaller cost is preferred since the resource of a WSN is scarce.

(2) Average Protection Ability (AP) on Primary Path: First, we assign a number to each node and link on the primary path as shown in Fig. 7. Moreover, the source and sink are assumed to be very robust so that they never fail. Let denote the total number of nodes and links on the primary path from source to sink. And, n_i^a represents the total number of alternate paths if the i -th node on the primary path fails. Similarly, l_i^a represents the total number of alternate paths if the i -th link on the primary path fails. Thus, we define the average protection (AP) ability acquired by each node and link on the primary path as follows.

$$AP = \frac{\sum_{i=2, i+2}^{k-1} n_i^a + \sum_{i=2, i+2}^k l_i^a}{k} \tag{2}$$

The numbers of alternate paths for the link or node on the primary path shown in Fig. 2 are 1, 1, 3, 2, 3, 2, and 1. Thus, the average protection acquired by each node and link is 1.8571. Obviously, a better topology has the higher AP value since more alternate paths are ready for protecting each link or node on the primary path.

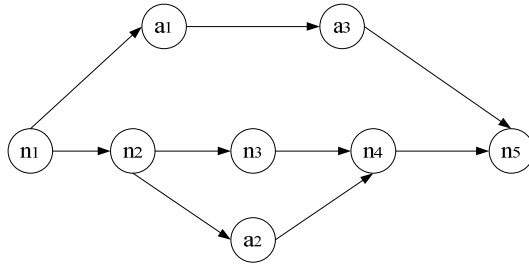


Fig. 2. Average protection ability

(3) Ratio of Protection to Cost (RPC): The goal of braided multipath approach is to construct a robust topology which can provide better average protection ability but make use of less resources (i.e., lower cost) to isolate failure and alleviate congestion. Thus, we define the ratio of protection to cost as follows.

$$RPC = \frac{AP}{C} \tag{3}$$

With the example shown in Fig. 4, the RPC is 4.086. Obviously, the higher the RPC value, the better the network topology.

Table 1. Comparison of different approaches

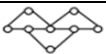



Approaches	Cost	AP	RPC
A 	0.6	1.2857	2.1428
B 	0.4545	1.8571	4.0860
C 	0.5455	1.8571	3.4044
D 	1	2.1429	2.1429

Table 1 shows the comparison of different multipath approaches. The proposed approach (B and C) can achieve the better RPC values than BM and DBM approaches. Moreover, if the total number of nodes and links on the primary path from source to sink is increased, our analysis results also show that the proposed approach has higher RPC value than the other approaches. To further validate the analysis results, we conducted simulation study.

4 Simulation Study

We conducted simulations experiments by making use of MATLAB to compare the performance of the different multipath approaches. The performance is measured in terms of C , AP , and RPC . And, we considered the factors, which affect the performance, as follows: node failure probability on primarily path and the path length between the sources and sink. In this simulation study, we denote the braided multipath as approach A, the proposed enhanced braided multipath as approach B, the enhanced braided multipath with the medium node as approach C, and the disjoint and braided multipath as approach D. We address the simulation results as follows.

(1) Cost: In this simulation experiment, the hop count of the path from source and sink is 4. Fig.3 shows the cost, in terms of extra nodes and links, of the different multipath approaches. Fig.4 shows how the path length affects the cost. The proposed approach B and C have better performance than approach A and D. Thus, the proposed approach needs less extra nodes and links to construct the braided multiple paths than the BM and DBM approach.

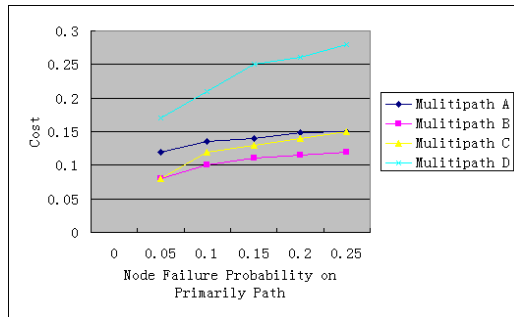


Fig. 3. Cost vs. node failure probability

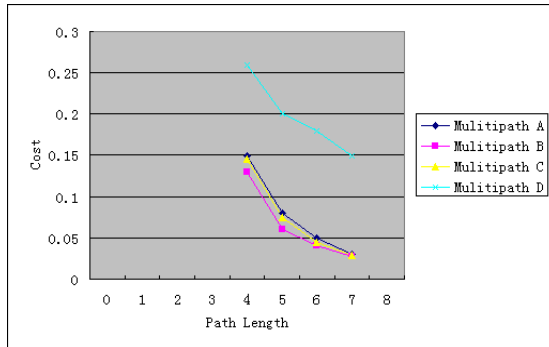


Fig. 4. Cost vs. path length

(2)Average Protection Ability on Primary Path: The alternate path that can protect longer segment of the primary path provides more protection on the primary path. And, the AP value represents the protection ability to isolate failure and alleviate congestion. Fig.5 shows the average protection abilities for the multipath approaches. The DBM approach has the highest protection ability because it uses more extra nodes and links to protect the primary path. Fig.6 illustrates the relationship between the path length and protection ability.

(3)Ratio of Protection to Cost: Fig.7 shows the simulation results of the RPC values for different multipath approaches. Approach B and C can provide higher

RPC values than approach A and D do. Although approach D has the highest protection ability as shown in Fig.8, it requires more extra nodes and links, so that the average protection ability contributed by each node and link for approach D is lower. From the above simulation results, the proposed approaches can achieve better protection ability per extra node and link than the other approaches.

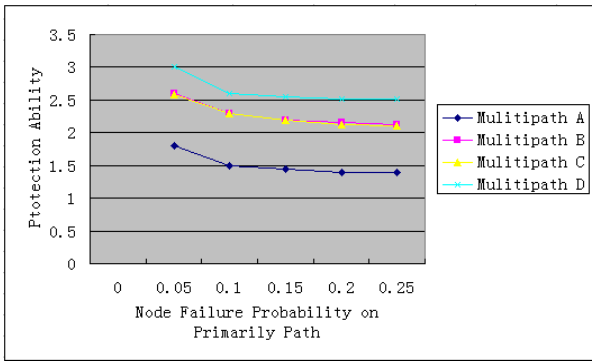


Fig. 5. AP vs. node failure probability

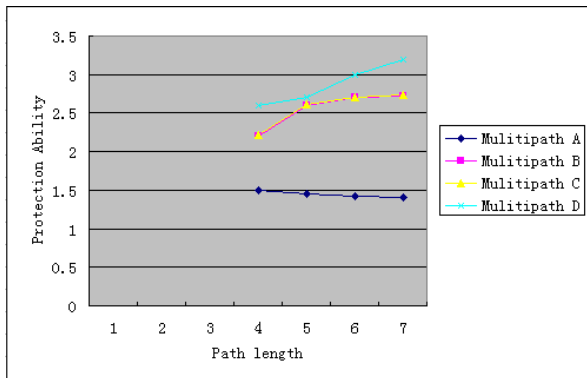


Fig. 6. AP vs. path length

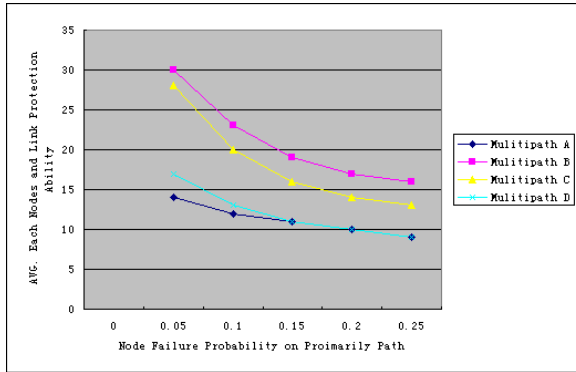


Fig. 7. RPC vs. node failure probability

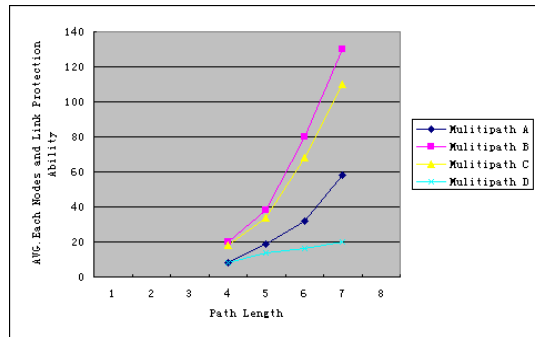


Fig. 8. RPC vs. path length

5 Conclusions

Wireless sensor networks are the promising platform on which versatile applications can be deployed. In this paper, we propose an effective multipath approach to resolve the congestion issue. The Effective Multipath approach not only can react to congestion faster but also can isolate the node failure. Moreover, we conducted simulation experiments by making use of MATLAB. The simulation results show that the performance of The Effective Multipath approach is better than those of the other multipath approaches.

References

1. Karl, H., Willig, A.: A Short Survey of Wireless Sensor Networks, Technical Report TKN-03-018, Telecommunication Networks Group, Technical University Berlin (October 2003)
2. Chen, D., Varshney, P.K.: QoS Support in Wireless Sensor Networks: A Survey. In: The Proc. of the In. Conf. on Wireless Networks 2004 (ICWN 2004), Las Vegas, Nevada, USA, June 21-24 (2004)

3. Akyildiz, I.F., Su, W., et al.: A Survey on Sensor Networks. *IEEE Communication Mag.* 40(8), 102–114 (2002)
4. Kang, J., Zhang, Y., et al.: TARA: Topology-Aware Resource Adaptation to Alleviate Congestion in Sensor Networks. *IEEE Trans. on Parallel and Distributed Systems* 18(7), 919–931 (2007)
5. He, T., Ren, F., Lin, C., Das, S.: Alleviating Congestion Using Traffic-Aware Dynamic Routing in Wireless Sensor Networks. In: 5th Annual IEEE CSC on Sensor, Mesh and Ad Hoc Communications and Networks, SECON 2008, June 16-20, vol. 2, pp. 233–241 (2008)
6. Wu, C., Shao, F., Wang, L.: MAC layer intelligent split-stream scheme for alleviating congestion in Ad hoc networks. In: ICWITS 2010, August 28, vol. 1, pp. 1–4 (2010)
7. Hou, Y., Leung, K.K., Misra, A.: Enhancing congestion control with adaptive per-node airtime allocation for wireless sensor networks. In: 2009 IEEE 20th IS on Personal, Indoor and Mobile Radio Communications, September 13-16, vol. 2, pp. 67–71 (2009)
8. Chauhan, N., Awasthi, L.K., Chand, N.: Global cooperative caching for Wireless Sensor Networks. In: 2011 World Congress on Information and Communication Technologies (WICT), December 11-14, pp. 235–239 (2011)
9. Akshay, N., Kumar, M.P., Harish, B., Dhanorkar, S.: An efficient approach for sensor deployments in wireless sensor network. In: Emerging Trends in Robotics and Communication Technologies (INTERACT), December 3-5, vol. 10, pp. 350–355 (2010)
10. Hoang, N.M., Son, V.N.: Disjoint and Braided Multipath Routing for Wireless Sensor Networks. In: Int. Symposium on Electrical and Electronics Engineering, October 11-12 (2005)

Immunity-Based Gravitational Search Algorithm

Yu Zhang, Yana Li, Feng Xia, and Ziqiang Luo

College of Information Science and Technology, Hainan Normal University, 571158 Haikou, China

bullzhangyu@yahoo.com.cn

Abstract. GSA (Gravitational Search Algorithm) is inspired by the Newton's law of universal gravitation and considered as a promising evolutionary algorithm, which has the advantages of easy implementation, fast convergence, and low computational cost. However, GSA has the disadvantages that its convergence speed slows down in the later search stage and it is easy to fall into local optimum solution. We proposed a novel immunity-based Gravitational Search Algorithm (IGSA) that is inspired by the biological immune system and the traditional gravitational search algorithm. The comparison experiments of GSA, IGOA and PSO (Particle Swarm Optimization) on 5 benchmark functions are carried out. The proposed algorithm shows competitive results with improved diversity and convergence speed.

Keywords: Gravitational Optimization Algorithm, Artificial Immune System, Evolution, Vaccine.

1 Introduction

The traditional gravitational search algorithm (GSA) firstly proposed by Rashedi et al. [1-4] in 2009 is a new intelligent heuristic optimization algorithm, which is based on the metaphor of gravitational interaction between masses. Like Particle Swarm Optimization (PSO) [5,6], GSA is an agent-based iterative optimization algorithm. In GSA, agents are considered as objects and their performance is measured by their masses. All these objects attract each other by a gravity force, and this force causes a movement of all objects globally towards the objects with heavier masses. The heavy masses correspond to good solutions of the problem. Specifically speaking, each mass (agent) in GSA has four specifications: its position, its inertial mass, its active gravitational mass, and its passive gravitational mass. The position of the mass corresponds to a solution of the problem, and its gravitational and inertial masses are determined using a fitness function. With time goes by, masses navigated by properly adjusting the gravitational and inertia masses are attracted by the heaviest mass that presents an optimum solution in the search space.

GSA has the advantages of easy implementation, fast convergence and low computational cost. However, GSA driven by the gravity law is easy to fall into local optimum solution. Moreover, the convergence speed slows down in the later search stage, and the solution precision is not good. Therefore, we introduce other operators

from biological immune system that can raise agent population diversity and solution accuracy to improve GSA performance.

Biological immune system (BIS) plays an important role in the defense of foreign invasion with a variety of antibodies, and thus keeps body healthy. BIS can quickly search for the best matching antibody when invaded by a known antigen. When an unknown antigen enters the body, BIS can also generate the optimum antibody through the dynamic adaptive learning. This shows that BIS has internal mechanisms of immune memory and antibody diversity. Inspired by the immune mechanisms, we proposed an Immunity-based Gravitational Search Algorithm (IGSA) based on GSA. In IGSA, we use memory antibody as vaccine to improve the convergence speed, and use antibody diversity to raise the diversity of agents. As a result, IGSA can avoid falling into local optimum solution and premature degradation, and therefore improves the global search capability and solution accuracy.

In the remaining sections of the paper, we first provide a brief review of GSA in Section 2. In Section 3, we introduce our IGSA in detail. Experiment and results are discussed in Section 4. Finally, a conclusion is stated in Section 5.

2 Traditional Gravitation Optimization Algorithm

2.1 Newtonian Laws of Gravitation and Motion

GSA (Gravitational Search Algorithm) could be considered as a closed system composed of a variety of masses obeying the Newtonian laws of gravitation and motion.

The Law of Gravity states that every particle of matter in the universe attracts every other particle with a force that is directly proportional to the product of the masses of the particles and inversely proportional to the square of the distance between them. Mathematically, this law can be translated into the equation shown below:

$$F = G \frac{M_1 M_2}{R^2} \quad (1)$$

Where F = the force of gravity, G = the gravitational constant, which adds the proper level of proportionality to the equation, M_1 & M_2 = the masses of the two particles, R = the straight-line distance between the two particles.

The law of Motion declares that the acceleration of an object produced by a total applied force is directly related to the magnitude of the force, the same direction as the force, and inversely related to the mass of the object. Mathematically, this law can be presented as following equation:

$$\sum F = Ma \quad (2)$$

The equations of motion are:

$$V = V_0 + at \quad (3)$$

$$S = V_0t + \frac{1}{2}at^2 \tag{4}$$

From equation (1), equation (2) and equation (3), we get

$$S = V_0t + \frac{1}{2} \frac{GM_2}{R^2}t^2 \tag{5}$$

2.2 The Principle of GSA

In GSA, the isolated system with N agents (masses), the position of the i -th agent is defined by:

$$X_i = (x_i^1, \dots, x_i^d, \dots, x_i^n), i = 1, 2, \dots, N \tag{6}$$

where x_i^d presents the position of i -th agent in the d -th dimension and n is the space dimension.

At a specific time t , the force acting on mass i from mass j is defined as following:

$$F_{ij}^d = G(t) \frac{M_i(t) \times M_j(t)}{R_{ij}(t) + \mathcal{E}} (x_j^d(t) - x_i^d(t)) \tag{7}$$

where M_j is the active gravitational mass related to agent j , M_i is the passive gravitational mass related to agent i , $G(t)$ is gravitational constant at time t , \mathcal{E} is a small constant, and $R_{ij}(t)$ is the Euclidian distance between two agents i and j .

The total force that acts on agent i in a dimension d is a randomly weighted sum of d -th component of the forces exerted from K agents:

$$F_i^d(t) = \sum_{j=1, j \neq i}^K \text{rand}_j F_j^d(t) \tag{8}$$

where rand_j is a random number in the interval $[0,1]$ and K is the set of first K agents with the best fitness value and biggest mass. K is a function of time, initialized to K_0 at the beginning and decreasing with time.

From the law of motion, the acceleration of the agent i at time t , and in direction d , $a_i^d(t)$, is given as follows:

$$a_i^d(t) = \frac{F_i^d(t)}{M_i(t)} \tag{9}$$

where $M_i(t)$ is the inertial mass of i -th agent.

The next velocity and position of an agent could be calculated as follows:

$$v_i^d(t+1) = \text{rand}_i \times v_i^d(t) + a_i^d(t) \tag{10}$$

$$x_i^d(t+1) = x_i^d(t) + v_i^d(t+1) \tag{11}$$

where $rand_i$ is a uniform random variable in the interval $[0, 1]$. This random number gives a randomized characteristic to the search.

The values of masses are calculated using the map of fitness. The gravitational and inertial masses are updated by the following equations:

$$m_i(t) = \frac{fit_i(t) - worst(t)}{best(t) - worst(t)} \quad (12)$$

$$M_i(t) = \frac{m_i(t)}{\sum_{j=1}^N m_j(t)} \quad (13)$$

where $fit_i(t)$ represent the fitness value of the agent i at time t , and, $worst(t)$ and $best(t)$ are defined as follows (for a minimization problem):

$$worst(t) = \max_{j \in N} fit_j(t) \quad (14)$$

$$best(t) = \min_{j \in N} fit_j(t) \quad (15)$$

3 Immunity-Based Gravitation Search Algorithm

From the principle of traditional gravitational search algorithm, we can reach these conclusions that GSA is efficient to reach the optimum value space, but inefficient in attaining the optimum value. Moreover, GSA has poor performance at the later search stage. These phenomena are due to the lack of agents' diversity in GSA. So, if we can enhance the diversity of agents, then the performance of GSA can be improved.

In artificial immune system, people generally consider an antibody as a candidate solution for the problem, antigen as the problem to be solved. The affinity between antibodies and antigens is calculated to evaluate the level how the antibody closes to the optimal solution to the problem. Some of prior knowledge and the characteristics of the problem be solved can be considered as vaccine [9,10]. Therefore, we proposed a novel Immunity-based Gravitation Search Algorithm (IGSA) that includes the mechanisms of antibodies diversity and immunity memory based on the GSA. The main idea of the proposed IGSA is that the characteristic of antibody diversity is to improve the solution space, and the characteristic of immunity memory is to enhance the solution quality. By taking these immune characteristics, IGSA will help to speed the convergence of evolutionary algorithms and improve the optimization capability.

The general flow chart of IGSA is as follows.

Step 1: Initialize parameters, including the gravitational constant G_0 , α , masses number N , and the maximum number of $Max_iteration$.

Step 2: Randomly generate the initial population A_k consisting of N masses with an initial velocity $V_0 = 0$.

Step 3: Extract the information to the problem be solved as vaccine V_a .

Step 4: Calculate the fitness of each mass in the current population A_k , and save the mass that has the best fitness as an immune memory mass M_j . Thereafter, determine

whether the end conditions are met; if the termination conditions are met, the algorithm stops and returns the result, otherwise, continue.

Step 5: Vaccinate the mass in the initial population A_k with a certain probability to form new mass population B_k .

Step 6: Substitute a part of masses with poor fitness in the population B_k for an immune memory mass, thus to generate new mass population C_k .

Step 7: Select part of masses in the population C_k based on mass concentration and randomly generate an additional part of masses to form the next generation A_{k+1} .

Step 8: Update the gravitational constant G , $best$, $worst$, and $M_i, i = 1, 2, \dots, N$.

Step 9: Calculate the suffered gravity for each mass.

Step 10: Calculate the acceleration and velocity for each mass.

Step 11: Update the location of each mass in the population A_{k+1} .

Step 12: Go to *Step 4*.

4 Experiments and Results

To objectively evaluate the performance of the proposed algorithm IGSA, comparison experiments with the traditional gravity algorithm GSA and particle swarm optimization PSO are carried out, and the parameters of the three algorithms are the same as much as possible. Numerical function optimization problems are selected as benchmark problems to test the performance of the three algorithms.

4.1 Benchmark Functions

The benchmark functions used in our experiments are listed in Table 1 from literature [2,11-14]. In these tables, n is the dimension of function. The minimum value of the functions in Table 1 is zero.

Table 1. Benchmark functions

Benchmark function	Optimum value
$F_1(X) = \sum_{i=1}^n x_i^2, x \in [-100, 100]^n$	0
$F_2(X) = \sum_{i=1}^n x_i + \prod_{i=1}^n x_i , x \in [-10, 10]^n$	0
$F_3(X) = \sum_{i=1}^n \left(\sum_{i=1}^n x_i \right)^2, x \in [-100, 100]^n$	0
$F_4(X) = \max\{ x_i , 1 \leq i \leq n\}, x \in [-100, 100]^n$	0
$F_5(X) = \sum_{i=1}^{n-1} [100(x_{i+1} - x_i^2)^2 + (x_i - 1)^2], x \in [-30, 30]^n$	0
$F_6(X) = \sum_{i=1}^n ([x_i + 0.5])^2, x \in [-100, 100]^n$	0
$F_7(X) = \sum_{i=1}^n ix_i^4 + random[0, 1], x \in [-1.28, 1.28]^n$	0
$F_8(X) = \sum_{i=1}^n [x_i^2 - 10 \cos(2\pi x_i) + 10], x \in [-5.12, 5.12]^n$	0

4.2 Experimental Results

We compare IGSA with GSA as well as PSO on the benchmark functions above where the population size $N=50$, dimension size $n=30$, and iteration number is 1000. The results are averaged over 50 runs and the average best solution, average mean fitness function, and median of the best solution in the last iteration are listed in Table 2.

Table 2. Experiment result of benchmark functions

		PSO	GSA	IGSA
F ₁	Average best	1.8×10^{-3}	7.3×10^{-11}	2.6×10^{-17}
	Median best	1.2×10^{-3}	7.1×10^{-11}	2.2×10^{-17}
	Average mean fitness	5.0×10^{-3}	2.1×10^{-10}	7.6×10^{-17}
F ₂	Average best	2.0	4.03×10^{-5}	2.23×10^{-8}
	Median best	1.9×10^{-3}	4.07×10^{-5}	2.47×10^{-8}
	Average mean fitness	2.0	6.9×10^{-5}	4.5×10^{-8}
F ₃	Average best	$4.1 \times 10^{+3}$	$0.16 \times 10^{+3}$	$0.08 \times 10^{+3}$
	Median best	$2.2 \times 10^{+3}$	$0.15 \times 10^{+3}$	$0.11 \times 10^{+3}$
	Average mean fitness	$2.9 \times 10^{+3}$	$0.16 \times 10^{+3}$	$0.10 \times 10^{+3}$
F ₄	Average best	8.1	3.7×10^{-6}	2.63×10^{-9}
	Median best	7.4	3.7×10^{-6}	2.7×10^{-9}
	Average mean fitness	23.6	8.5×10^{-6}	4.2×10^{-9}
F ₅	Average best	$3.6 \times 10^{+4}$	25.16	25.83
	Median best	$1.7 \times 10^{+3}$	25.18	25.95
	Average mean fitness	$3.7 \times 10^{+4}$	25.16	25.89
F ₆	Average best	1.0×10^{-3}	8.3×10^{-11}	2.08×10^{-17}
	Median best	6.6×10^{-3}	7.7×10^{-11}	2.28×10^{-17}
	Average mean fitness	0.02	2.6×10^{-10}	2.56×10^{-17}
F ₇	Average best	0.04	0.018	0.011
	Median best	0.04	0.015	0.012
	Average mean fitness	1.04	0.533	0.053
F ₈	Average best	55.1	15.32	12.93
	Median best	55.6	14.42	13.25
	Average mean fitness	72.8	15.32	14.23

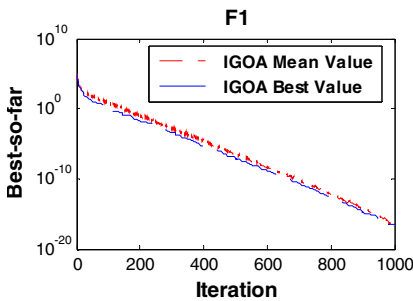


Fig. 1. The evolutionary curve of IGSA for minimization of F1

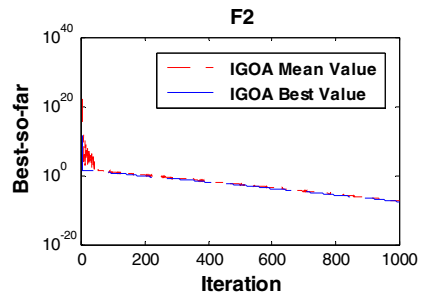


Fig. 2. The evolutionary curve of IGSA for minimization of F2

To better understand the search process and see how IGSA reaches the best solution, we show the evolutionary curve of IGSA for benchmark functions $F1$ to $F8$ below.

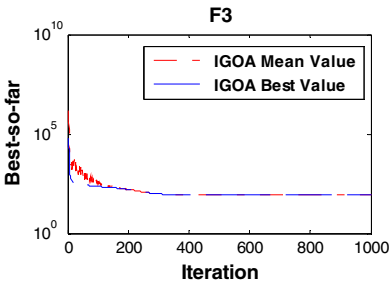


Fig. 3. The evolutionary curve of IGSA for minimization of $F3$

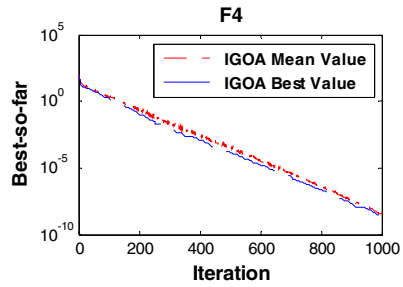


Fig. 4. The evolutionary curve of IGSA for minimization of $F4$

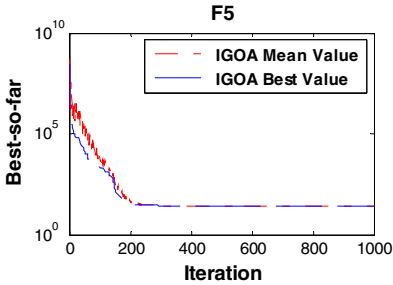


Fig. 5. The evolutionary curve of IGSA for minimization of $F5$

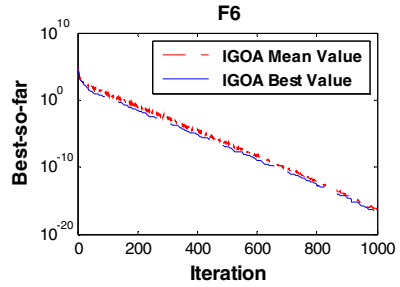


Fig. 6. The evolutionary curve of IGSA for minimization of $F6$

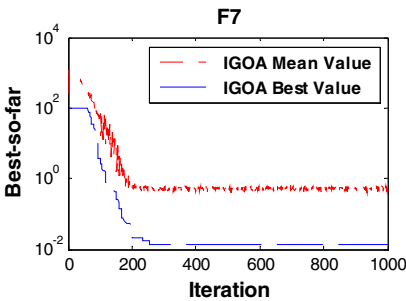


Fig. 7. The evolutionary curve of IGOA for minimization of $F7$

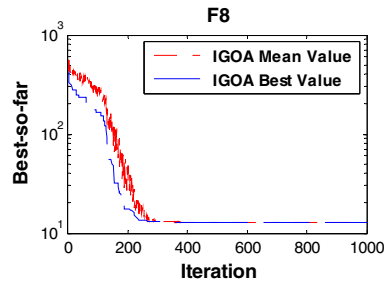


Fig. 8. The evolutionary curve of IGOA for minimization of $F8$

From the experimental results listed in Table 2, we can clearly see that IGSA provides better results than GSA and PSO on all benchmark functions except for $F5$. From the figure 1 to 8, we can see that IGSA reaches the best solution with 200 to 400 iterations.

The better performance of IGSA benefits from the following features: (1) elitist antibody gene segments are retained, and inherited to offspring through vaccination, which will enhance local search capability and improve the convergence speed; (2) antibody population diversity is achieved by mutation that avoid the degenerative phenomenon and to enhance the global optimization capability.

5 Conclusion

When people solve the optimization problems with high dimension, they often resort to evolutionary algorithms. Over the past decades, there are many evolutionary algorithms that inspired by the behaviors of natural phenomena. We draw the mechanisms of immune system and gravity law, and propose a novel immunity-based gravitation search algorithm (IGSA). The proposed algorithm, which is based on the gravitational search algorithm (GSA), includes the mechanisms of antibodies diversity and immunity memory. The mechanism of antibodies diversity is to enhance the population diversity, and the mechanism of immunity memory is to save some optimum value of current population. With these features, the proposed IGSA can help speed the convergence of evolutionary algorithms and improve the optimization capability. The comparison experiments of IGSA, GSA and PSO on some benchmark functions are carried out. The proposed algorithm shows competitive results with improved diversity and convergence.

Acknowledgement. This work was supported by the Laboratory Project of Hainan Normal University (kfsy11048), Hainan Natural Science Foundation (610220), Project of Zhejiang Key Laboratory of Information Security (2010ZISKL007), and the Science and Technology Project of Hainan Normal University (00203020214).

References

1. Rashedi, E.: Gravitational Search Algorithm. MS Thesis, Shahid Bahonar University of Kerman, Iran (2007)
2. Rashedi, E., Nezamabadi-pour, H., Saryazdi, S.: GSA: A Gravitational Search Algorithm. *Information Sciences* 179(13), 2232–2248 (2009)
3. Rashedi, E., Nezamabadi-pour, H., Saryazdi, S.: BGSA: binary gravitational search algorithm. *Natural Computing* (December 2009), <http://dx.doi.org/10.1007/s11047-009-9175-3>
4. Rashedi, E., Nezamabadi-pour, H., Saryazdi, S., et al.: Allocation of Static Var Compensator Using Gravitational Search Algorithm. In: *First Joint Congress on Fuzzy and Intelligent Systems*. Ferdowsi University of Mashhad, Iran (2007)
5. Zhan, Z., Zhang, J., Li, Y., et al.: Adaptive Particle Swarm Optimization. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics* 39(6), 1362–1381 (2009)
6. Chen, D., Zhao, C.: Particle swarm optimization with adaptive population size and its application. *Applied Soft Computing* 9(1), 39–48 (2009)

A Clustering Method Based on Time Heat Map in Mobile Social Network

Wang Ye, Wang Jian, and Yuan Jian

Department of Electronic Engineering, Tsinghua University, Beijing 100084, China
iequinox@126.com

Abstract. Mobile Social Network Services (MSNS) have collected massive amount of users' daily positioning information, which could be used for data mining to learn people's habits and behaviors. This paper proposes a novel clustering method to group nodes according to timestamp information, by analyzing the Time Heat Map (THM), i.e. the activity level distribution of a node during different time intervals. We have employed large amounts of anonymized positioning records coming from a real MSNS, which has extinguished this paper from other researches that use volunteers' daily GPS data. Experiment results have shown that this method not only reveals some interesting features of human activities in real world, but also can reflect clusters' geographical "interest fingerprints" affectively.

Keywords: mobile social computing, MSNS, clustering, time heat map.

1 Introduction

MSNS (Mobile Social Network Service) such as Foursquare [1] has allowed people to share one's current location, which gives researchers an opportunity to discover the true meaning behind raw data. Recent researches on mobile social networks vary from user context conception [2], location sharing [3, 4], network metrics [5], to designs of geo-social applications [6, 7]. These articles focus mainly on the architect and application of mobile end and hardly did any data analysis. [8] did an outstanding work on identifying user groups based on user co-presence through mobile devices. However, their data comes from a rather small sample in a campus, and the scalability of their method for large scale data in real mobile social applications is not validated.

Besides these new researches on mobile social computing, there are also articles related to traditional geographic data mining. In [9, 10, 11, 12], detailed mathematical models and methods are given to dig into user GPS trajectories in order to learn user behavior patterns. However, these methods all rely on continuous or intensive GPS trajectories, which are not available in real MSNS applications, due to GPS's high energy consumption and users' privacy concern. Moreover, there are some unique facts about social network data: 1. a large part of the nodes are inactive ones with few records; 2. the original raw data is high-dimensional and tempo-spatially sparse. Thus not all the classical methods are suitable for MSNS data.

There are two ways to solve these problems: one is through data cleaning and pretreatment; and the other one is to aggregate the original data into fewer

dimensions. Our method is to cluster the nodes, based on the timestamp tag of each positioning record. By aggregating the records according to several selected time interval, we will be able to obtain a time heap map, which should depict this node's average activity level in different intervals. And this process can reduce the dimension from tens of thousands to a remarkably small level. And this is a creative work in utilizing not only geographical data but also temporal information. Paper [13] proposed some paradigms and methodologies in spatio-temporal data mining, with lots of marvelous understanding about the essence of time and space. But no methods or examples are given.

The rest of this paper will be organized as follows: Section 2 demonstrates the data prototype and the process of our clustering algorithm; Section 3 displays the clustering results and evaluation with interest fingerprints; Section 4 will give a summarization of our conclusion.

2 Data Collection and Clustering Algorithm

2.1 Anonymized Raw Data

Our data comes from a MSNS service in China and has been anonymized for analyzing. During a period of 11 months, more than 600 thousand positioning records of about 17,000 users have been collected. Each record consists of an anonymous user id, a location id and the timestamp of this record. Besides this, the locations have been tagged with a category field, which roughly indicates the type of a place.

2.2 Definition of Time Heat Map (THM)

A Time Heat Map (THM) is a normalized vector obtained by dividing the timestamp field into several periodical intervals (say, 24 hours) and sum up all the records accordingly. It can indicate the activity level distribution. When there are M time intervals, the THM vector v has M components.

Assuming the total number of records for a given node i in time interval t_j is n_{ij} , the exact definition of time heat map for this node is as follows:

$$v_{ij} \triangleq \frac{n_{ij}}{\sum_{j=1}^M n_{ij}}, \text{ where } 1 \leq j \leq M. \quad (1)$$

Note that this is a normalized definition, and thus v can be also regarded as a probability vector. This paper uses "hour-of-day" THM, by dividing each day into 24 hour intervals and aggregate records according to the hour information.

2.3 K-Means Clustering Method

K-Means clustering method [14] is a classical clustering method in Data Mining. Its basic operation flow could be represented with Fig. 1. This paper uses the K-Means Clustering Program in Weka Platform [15] to complete this part. And Euclidian distance is selected as the distance function.

1. At the beginning of the algorithm, choose K random points as the centers of K clusters;
2. For each point to be sorted, calculate its distances to the K clusters, and sort it to the nearest cluster;
3. After all the points are sorted, there will be K clusters. Calculate the average value of each cluster and use it as the new center.
4. Go to step 2 again, and iterate until the result converges.

Fig. 1. The process of simple K-Means Clustering method

One of the short comings of simple K-Means is that the cluster number K needs to be determined manually. Different initial setting will result in different cluster groups. But in our case, the cluster number K has a clear physical meaning, so we do not expect the number to be too large or too small. A rough range of K could be specified, say, from 3 to 9. But a standard process is still needed.

On the other hand, we would like to obtain K groups with approximately the same size. Thus our aim is to maximize the uniformity of the clusters.

Let U_K denote the uniformity of clusters when cluster number equals to K , and p_k denote the proportion of the k -th cluster in total. U_K is defined as follows:

$$U_K \triangleq -\frac{\sum_{k=1}^K p_k \ln p_k}{\ln K}, \text{ where } 1 \leq k \leq K. \quad (2)$$

For a zero-one distribution, U_K equals to 0, and for a uniform distribution, $U_K = 1$. Combined with the rough range of K mentioned above, now we can determine the value of parameter K according to U_K :

$$K = \arg \max_{3 \leq k \leq 9} U_k. \quad (3)$$

2.4 Data Pretreatment

For real-world MSNS records, there are inevitably large numbers of inactive nodes, which could severely undermine the effectiveness of clustering. We will filter out nodes with less than n_{thresh} records, where n_{thresh} stays to be determined.

Assume that there are already L sampling records for a given node, with the corresponding THM vector v_L , and the real THM vector is v_{real} . The record number in time interval j ($1 \leq j \leq M$) is $n_{L,j}$, and the M components of vector v_L is represented by $v_{L,j}$. A new-arriving sampling record will cause the THM vector change from v_L to v_{L+1} , and this new record will appear in time interval j with the probability:

$$p_j = v_{\text{real},j} \quad (4)$$

The consequential relative change of vector is:

$$\frac{|v_L - v_{L+1}|}{|v_L|} = \frac{\sqrt{\sum_{j=1}^M (v_{L,j} - v_{L+1,j})^2}}{|v_L|} \tag{5}$$

Apply (1) to the equation above, and we will get:

$$\frac{|v_L - v_{L+1}|}{|v_L|} = \frac{\sqrt{\sum_{j=1}^M p_j \left(\left(\frac{n_{L,j}}{L} - \frac{n_{L,j} + 1}{L+1} \right)^2 + \sum_{1 \leq k \leq M, k \neq j} \left(\frac{n_{L,k}}{L} - \frac{n_{L,k}}{L+1} \right)^2 \right)}}{|v_L|} \tag{6}$$

When L is large enough, we can assume that v_L is already close to v_{real} , thus $p_j \approx \frac{n_{L,j}}{L}$, and $L+1 \approx L$. And finally we will get:

$$\frac{|v_L - v_{L+1}|}{|v_L|} = \frac{\sqrt{1 - |v_L|^2}}{|v_L|} \cdot \frac{1}{L} \tag{7}$$

According to (1), it is easy to deduce that $\frac{1}{\sqrt{M}} \leq |v_L| \leq 1$, where M is the number of time intervals, so the factor $\frac{\sqrt{1 - |v_L|^2}}{|v_L|}$ in (7) has a limited range $[0, \sqrt{M-1}]$. In our case, $M=24$, then we can get a criteria for determining n_{thresh} :

$$\frac{|v_L - v_{L+1}|}{|v_L|} \leq \frac{\sqrt{23}}{n_{thresh}} \tag{8}$$

3 Clustering Results and Evaluation

3.1 User Clusters

According to (8), we decide to set $n_{thresh}=50$, which could ensure that the relative change in (8) will not exceed 10% even in the worst case. We got 1822 user nodes eventually. According to (3), the calculated K equals to 5 where U_K reaches its maximum. The clustering result is shown as in Fig. 2, where there are 5 different THM patterns. Here are some interesting findings:

Three minor peaks in three meal-time period: This could be observed in almost all THM patterns, around 8~9, 12~13 and 17~18 o'clock respectively;

The reflection of “nine to five” working schedule in human activities;

The patterns can quite differ from each other: (c), (d) and (f) have one sharp peak, while (e) has two moderate peaks. Group (b) is so ordinary that its pattern resembles the overall mean vector in (a).

The majority lives a “mediocre” life, and is active in the daytime. Group (b) and (e) are the two biggest clusters, and they have flatter peaks and account for 55% altogether.

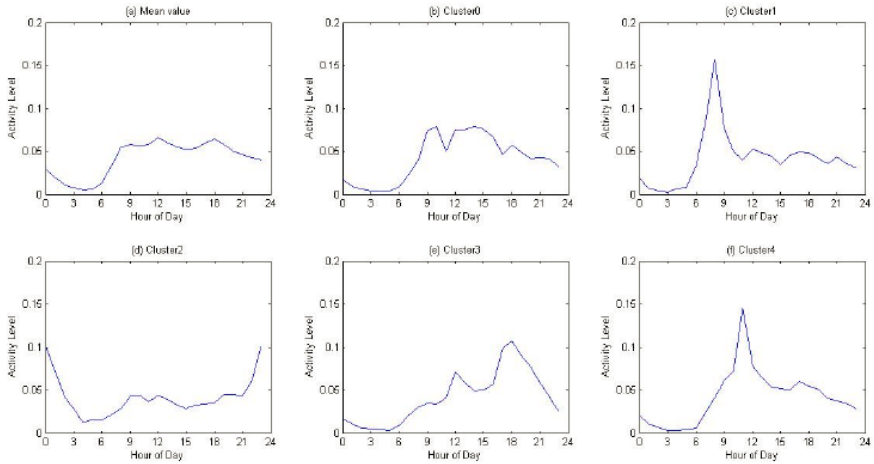


Fig. 2. THM vectors for different user clusters: (a) overall; (b) Cluster 0: the “Ordinary” group; (c) Cluster 1: the “Early Birds”; (d) Cluster 2: the “Night Owls”; (e) Cluster 3: the “Evening” group; (f) Cluster 4: the “Morning-Midday” group

3.2 Location Clusters and “Interest Fingerprints” Evaluation

For location nodes, $n_{\text{thresh}}=50$ would result in few available data. Thus we select location nodes with the following rule: in each location category, select the first 50 popular location nodes that have less than 10 records. Finally we got a sample set of 546 nodes. From equation (3) we calculated that the ideal value of K is 4.

There are several observations:

Also has three “meal-time” minor peaks and “nine-to-five” pattern; but this is only clearly observed in the overall THM in Fig. 4(a);

The THM in each cluster has flatter peaks than in Fig. 3; the boundaries seem to be more vague; and the local minor peaks seem to be less regular;

The biggest cluster (37%) is prosperous in the afternoon, while the smallest group (12%) operates more at night. But unlike human users, these “Night” locations’ heat drops before midnight.

We use the term “Interest Fingerprint” to denote the distribution of location category in a given cluster. The category’s meanings are shown in Table 1. For each cluster, calculate the normalized distribution of location categories, as in Table 2. Then we plot Fig. 4, the interest fingerprints for each cluster according to Table 2.

Fig. 4 below has clearly illustrated that the clusters based on temporal information can also effectively reflect geographical preferences. For example, we can know that cluster 0 is more connected to schools, outdoor places and transportation, and they are usually active in the afternoon around 3 o’clock. Cluster 1 is active at night, and the

corresponding locations are more likely to be entertainment or life related ones, while impossible to be malls and shops, which really makes sense. Cluster 2 shows strong favor to markets, malls and transportation, and it is not difficult to understand why this group is usually active from late in the afternoon to evening. Finally, cluster 3 has extremely intimate relation with offices and working places, which explains why this cluster is the “Morning Group”.

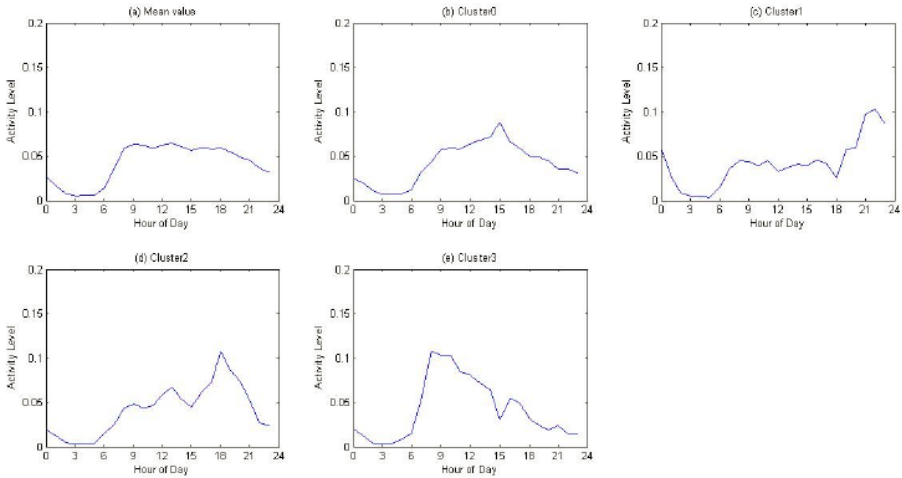


Fig. 3. THM vectors for different location clusters; (a) overall; (b) Cluster 0: the “Afternoon” group; (c) Cluster 1: the “Night” group; (d) Cluster 2: the “Evening” group; (e) Cluster 3: the “Morning” group;

Table 1. The meanings of different location category

ID	CATEGORY	ID	CATEGORY
0	Markets	6	Schools & colleges
1	Restaurants	7	Communities
2	Malls & shops	8	Companies & offices
3	Cinemas & theatres	9	Outdoor places
4	Life facilities	10	Subways & airports
5	Entertainments		

Table 2. The normalized distribution of location categories in each cluster (in percentage %)

Clusters	Location Category Distribution (%)										
	0	1	2	3	4	5	6	7	8	9	10
Cluster 0	5.4	9.4	9.9	10.9	5	9.4	12.4	9.4	3.5	11.4	13.4
Cluster 1	11.8	10.3	0	8.8	26.5	14.7	4.4	8.8	8.8	4.4	1.5
Cluster 2	14.7	10	17.3	7.3	8.7	9.3	4	8	3.3	5.3	12
Cluster 3	4	7.1	3.2	8.7	7.1	5.6	12.7	10.3	25.4	12.7	3.2

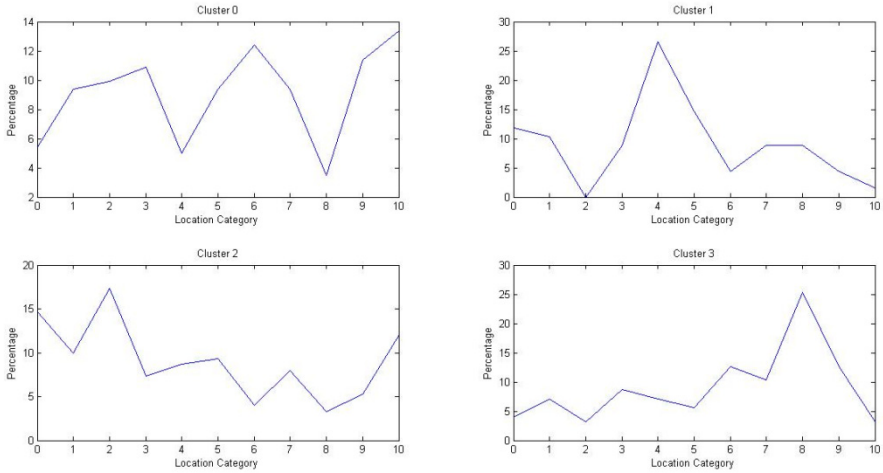


Fig. 4. Different “Interest Fingerprints” in four location clusters (Cluster 0~3)

4 Conclusions

Our clustering method has successfully clustered nodes into groups, and each of the groups shows a strong preference to a certain active time period. We have discovered the three minor peaks for “meal time” and a “nine-to-five” band-pass pattern in the THMs, which implies that this can highly reflect human behaviors in reality. Moreover, after plotting the distribution of location category for clustering results, different groups have shown rather different “interest fingerprints”, which can validate the effectiveness of our method.

References

1. <https://foursquare.com/>
2. Ofstad, A., Nicholas, E., Szcodronski, R., Choudhury, R.R.: AAMPL: accelerometer augmented mobile phone localization. In: MELT (2008)
3. Gaonkar, S., et al.: Micro-Blog: Sharing and Querying Content Through Mobile Phones and Social Participation. In: Proc. of MobiSys 2008, Breckenridge, CO, USA (2008)
4. Mody, R.N., Willis, K.S., Kerstein, R.: WiMo: Location-Based Emotion Tagging. In: Proceedings of the 8th International Conference on Mobile and Ubiquitous Multimedia, Cambridge (November 2009)
5. Scellato, S., Mascolo, C., Musolesi, M., Latora, V.: Distance Matters: Geo-social Metrics for Online Social Networks. In: WOSN 2010 (2010)
6. Pietiläinen, A.K., Oliver, E., Lebrun, J., Varghese, G., Diot, C.: MobiClique: middleware for mobile social networking. In: WOSN 2009 (August 2009)
7. Zhang, L., Ding, X., Wan, Z., Gu, M., Li, X.-Y.: Wiface: a secure geosocial networking system using wifi-based multi-hop manet. In: Proceedings of the 1st ACM Workshop on Mobile Cloud Computing & Services, pp. 1–8 (2010)

8. Gupta, A., Paul, S., Jones, Q., Borcea, C.: Automatic identification of informal social groups and places for geo-social recommendations. *International Journal of Mobile Network Design and Innovation* 2(3/4), 159–171 (2007)
9. Hariharan, R., Toyama, K.: Project Lachesis: Parsing and Modeling Location Histories. In: Egenhofer, M., Freksa, C., Miller, H.J. (eds.) *GIScience 2004*. LNCS, vol. 3234, pp. 106–124. Springer, Heidelberg (2004)
10. Zheng, Y., et al.: Learning transportation modes from raw GPS data for geographic applications on the Web. In: *Proceedings of WWW 2008*, pp. 247–256. ACM Press (2008)
11. Zheng, Y., Zhang, L., Xie, X., Ma, W.-Y.: Mining interesting locations and travel sequences from GPS trajectories. In: *WWW 2009: Proc. of the 18th Intl. World Wide Web Conference*, pp. 791–800 (April 2009)
12. Patterson, D.J., Liao, L., Fox, D., Kautz, H.: Inferring High-Level Behavior from Low-Level Sensors. In: Dey, A.K., Schmidt, A., McCarthy, J.F. (eds.) *UbiComp 2003*. LNCS, vol. 2864, pp. 73–89. Springer, Heidelberg (2003)
13. Roddick, J.F., Lees, B.G.: Paradigms for spatial and spatio-temporal data mining. In: Miller, H.G., Han, J. (eds.) *Geographic Data Mining and Knowledge Discovery*. Taylor & Francis, London (2001)
14. Lloyd, S.P.: Least Squares Quantization in PCM. *IEEE Transactions on Information Theory* IT-28, 129–137 (1982)
15. <http://www.cs.waikato.ac.nz/ml/weka/>

TNC-eSA: An Enhanced Security Access Solution to Office Networks

Jun Ma¹, Yuan-bo Guo¹, and Jinbo Xiong²

¹Zhengzhou Information Science and Technology Institute Shangchengdong Road, Guancheng District, Zhengzhou 450004 ZhengZhou, China

²Faculty of Software Fujian Normal University Minhou Section, 350108 Fuzhou, China
sijunhan@gmail.com, yuanbo_g@hotmail.com, jinbo810@163.com

Abstract. With the computer technology booming and the Internet applications gradually spreading all over, the office network speeds up its pace to informationization, automation and networking. However, the complexity of its users, the increasingly expanding of the fields it involves in, the disparity of its management, and the variety of the access to the network make the security of the current office network severe challenging. At this background, the trusted network connection bases on the endpoint's security and trust, proposes the whole architecture of the trusted network, emerges as a new approach for the problem of the office network's secure access. In this paper, TNC-eSA, an enhanced security solution to access office networks prototype is designed and implemented. TNC-eSA not only provides the features of TNC, but also achieves stronger security and higher performance by introducing extends 802.1X and dynamic extraction of endpoint characters. Experiment demonstrated the advantages of the TNC-eSA solution.

Keywords: office network, TNC, 802.1X.

1 Introduction

With the increasing popularization of computer network technology, it is an urgent demand to develop office automation through network and applies it to high-security level institutions, including trades, governments and militaries. However, in the actual process of office network operating, due to a wider coverage area, complicated staff structure, different levels of management, and the development of wireless mobile office network based WLAN, administrator of office network can not guarantee the safe transmission and management of documents on the network. Furthermore, hack attacks and vulnerability attacks have made great damage to office network. Therefore, it is necessary to present security solution to office networks.

Access authentication is the dominant security procedure to office network. There has already been some network connect and access control mechanisms, such as 802.1X, VPN, as well as encryption techniques which includes Extensible Authentication Protocol-Message Digest V5 (MD5), Transport Layer Security (TLS), and Secure Sockets Layer (SSL).However, security breaches are still causing

significant losses and damages. A security survey for the first half of 2010 states [1] that, malware infection and hack attacks to office networks have already run up to 30%, and demonstrate the raising tendency because of untrusted and insecure endpoint users.

Some papers [2] [3] present their views to improve security and efficiency of accessing office network, whereas available endpoint security has not been considered any more. This condition possibly causes that unavailable user with available credential still could access office network. The Trusted Network Connection (TNC) bases on the endpoint security and trust, proposes the whole architecture of the trusted network, and emerges as a new approach for the problem of the office networks' secure access.

An enhanced Access Authentication System in TNC architecture (TNC-eSA) is designed and implemented in this paper. Aimed to the vulnerability of the password authentication pattern the traditional endpoints use, the theory of data transition is adopted over the trusted endpoints and the information of the endpoint system, the hardware information, and the users' identification is extracted as the multiple authentication factors to access authenticate over terminal and to enhance the trust and integrity of the terminal authentication information.

The rest of the paper is organized as follows: Section 2 presents overview of TNC, and describes the process of integrity checking in TNC architecture. In Section 3, we propose our solution and design details TNC-eSA. Some main implement are illustrated and analyzed in section 4, and finally we give our conclusion in section 5.

2 Overview of TNC

In this section, we describe the architecture of TNC, and then give a brief analysis of integrity checking in TNC.

2.1 The Architecture of TNC

Trusted Network Connect (TNC) is an open standard network access control architecture, which is defined and promoted by Trusted Computing Group (TCG). It is a network access control standard with a goal of multi-vendor endpoint policy enforcement, which is compatible with existing network connect and authentication technologies, provides security access specification and standard authentication architecture to protect network's security and trusted. The TNC architecture focuses on interoperability of network access control solutions and on the use of trusted computing as the basis for enhancing security of those solutions. Integrity measurements are used as evidence of the security posture of the endpoint so access control solutions can evaluate the endpoint's suitability for being given access to the network. The TNC Architecture is shown in Figure 1[4].

The architecture incorporates several roles, functions, and interfaces which are discussed below. (1) five roles: the Access Requestor (AR), the Policy Enforcement Point (PEP), the Policy Decision Point (PDP), the Metadata Access Point (MAP), and the MAP Client (MAPC); Among of these roles, PEP, MAP and MAPC are optional

roles.(2)three horizontal shaded layers: Integrity Measurement Layer, Integrity Evaluation Layer, Network Access Layer;(3)and more interfaces: Integrity Measurement Collector Interface (IF-IMC), Integrity Measurement Verifier Interface (IF-IMV), Network Authorization Transport Protocol (IF-T),etc.

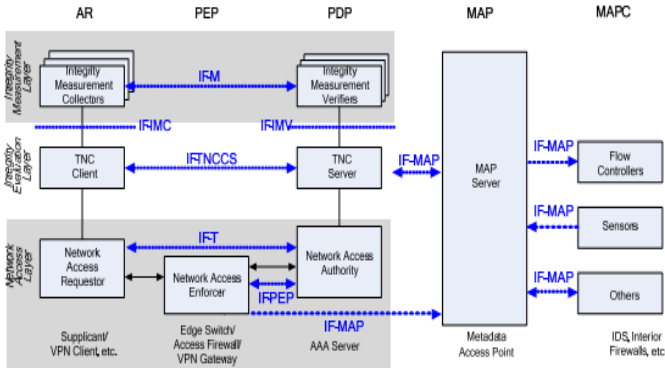


Fig. 1. The TNC Architecture

2.2 Brief Description of Integrity Checking Flow

Compared with now available access authentication mechanism, integrity checking of client and server are provided in TNC architecture. Client acts as AR, Server acts as PDP must check integrity of their own, and make integrity checking handshake before AR accesses network.

The description of integrity checking flow is as follows: Prior to beginning a network connection, TNCC of AR initializes the IMC of AR, which includes defining the necessary connection IDs and IMC IDs, and ensuring that the TNCC has a valid connection state with the IMC. TNCC checks the integrity of the IMCs. Similarly, the TNCS of PDP checks the integrity of the IMVs of PDP. The components of AR and PDP begin the exchange of messages pertaining to the integrity check for an integrity check handshake. After completing integrity check handshake between AR and PDP, they begin to network access process with authentication protocols.

3 Design of TNC-eSA

In this section, we propose the TNC-eSA in TNC architecture. Abide by integrity checking of TNC, we show basic design idea, and describe access flow based on extend 802.1X.

3.1 The Basic Design Idea of TNC-eSA

Transparency Endpoint Integrity Checking To User

In TNC-eSA, we design a client agent program, named wire and wireless AAA Supplicant (W2AAASupplicant) that accomplishes access authentication of AR. In the

access process, the agent firstly checks local position integrity, including the integrity of hardware platform and software platform. All components whichever integrity checking is failed, the agent can not connect the protected office network subsequently. For example, Series number of hard disk or OS is inconsistent with registered one formerly, it is reduces that the agent can not initial to connect protected office network.

Transparency Access Mode to User

Since the present office network tends to be mobile, W2AAASupplicant is designed which can be used in both the wired network and the wireless network, makes the endpoint users access safely anytime and anywhere without caring about the access way. After checking integrity, the agent diagnoses automatically network interface, and then begins to connect protected office network.

Extensible 802.1X Protocol

The 802.1X standard provides a framework for port based access control (PBAC) that is in accordance with TNC standard. In TNC-eSA, integrity measurement and reporting is provided to an 802.1X deployment by additional data regarding the integrity status of the Agent. After integrity checking handshake between W2AAASupplicant (AR) and (Radius Server) PDP was finished, they can perform mutual authentication.

Supporting Different Authentication Protocol

The security requests and levels of different users are different in the office networks. IN TNC-eSA, we propose a low security-level but simple access scheme based on EAP-MD5 and a high security-level with certificate access scheme based on EAP-TLS that unify the access authentication of different users.

3.2 Authentication Procedure of TNC-eSA

Figure 2 shows a typical TNC-eSA authentication procedure. We now consider the following four activities that describe how TNC-eSA works.

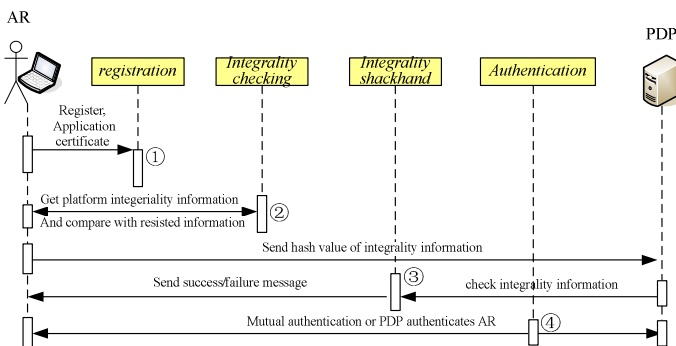


Fig. 2. TNC-eSA authentication procedure

Registration of Endpoint Platform

In order to integrity checking and authentication, endpoint needs to install W2AAASupplicant and runs it to registration platform information via GUI interface. The agent access a MAP database which store basic information of every user's platform, to submit available information and to get public key certificate (if using EAP-TLS, a certification is indispensable).

Integrity Checking of Endpoint

Both user platform and AAA server platform are required integrity checking respectively. Take AR for an instance, The agent of AR, W2AAASupplicant, firstly gets hardware information such as MAC address, CPU SN, HD SN etc; Software information such as OS SN; application program version etc, and certification information (if using EAP-TLS protocol) or user name/password information (if using EAP-MD5 protocol). Thereafter, the agent calculates the hash value of this information, and compares the consequence with registered value from MAP database. If only checking passed, the agent sends out connect request message.

Procedure of Integrity Checking Handshake between AR and PDP

In this procedure, the first step that AR sends connect request and PDP sends integrity challenge message to AR via PEP for beginning of integrity checking. The next step that AR sends hash value of integrity to PDP, PDP gets hash value of the AR rejected from MAP, and check whether the value is same with the value from AR. If the checking passes, the process of authentication is started.

Choice Authentication Protocol

According to different security level, AR may specify one way certification, EAP-MD5, or mutual authentication, EAP-TLS that need public key certificate in the TNC-eSA. If authentication successes, this procedure of authentication is end and AR can securely use office network service.

4 An Implement of TNC-eSA

In this section, we develop main modules of the TNC-eSA. Due to freeradius[5] which is adopted in the TNC-eSA has been supported by TNC standard, the main work of this paper is engaged in the implement of W2AAASupplicant which is a extend program of 802.1X open source software Xsupplicant [6]. Consider easy use for endpoint users, we implement W2AAASupplicant in Windows environment.

4.1 Acquirement of Integrity Information

We use a program specification of Windows Management Instrumentation (WMI) to acquire system information of software platform and hardware platform. Using WMI Query Language (WQL), the core code for getting information of AR platform is as follows:

Get Bios

```

ManagementClass MyClass = new ManagementClass
("Win32_BIOS");
ManagementObjectCollection MyCollection = MyClass.
GetInstances();
foreach(ManagementObject mo in MyCollection){
    MBIOSID = mo.Properties ["SerialNumber "
].Value.ToString();}

```

Get NIC

```

ManagementObjectSearcher query = new
ManagementObjectSearcher ("SELECT * FROM
Win32_NetworkAdapterConfiguration WHERE IPEnabled =
'TRUE'");
ManagementObjectCollection queryCollection =
query.Get();
foreach (ManagementObject mo in queryCollection){
NetworkCardDescription =
mo.Properties["Description"].Value.ToString();
    NetworkCardMACAddress = mo.Properties ["MAC Address "]
.Value.ToString();
}

```

Get OS SN

```

OperatingSystem os = Environment.OSVersion; string
usersn;
usersn = System.Environment.SerialNum;

```

4.2 The Choice of Authentication Protocol

According to context of proceeding, W2AAASupplicant needs to choice dynamically method of authentication to transfer message to freeradius. The code of our implement is as follows:

Eap-Md5

```

#include "ProPageMD5.h"
if (pPropertySheet->md5_enable) {
...
strcpy(tempmd5name,pPropertySheet-
>m_md5_username);strcpy(tempmd5pass,pPropertySheet-
>m_md5_password);}

```

Eap-tls

```

#include "ProPageTLS.h"
if (pPropertySheet->tls_enable) {
...
strcpy(temptlsname,pPropertySheet-> m_tls_username);
strcpy(temptlspass,pPropertySheet->m_tls_password);
strcpy(temptlscert,pPropertySheet->m_tls_certificate);
}

```

4.3 The Choice of NIC

The choice of NIC is supported by W2AAASupplicant, regardless of wire or wireless network device. The code reality is as follows:

```

while ((eapol_authenticate() == 0)&&( close_program
==false)) { if(running == TRUE){
    auth1 = get_current_state();
    if(strcmp(auth1,auth2) != 0){
        m_STA.SetWindowText(auth1);}
    auth2=auth1; }
    if( auth2 == "authentication success\n(you
passed!)") { m_logoff.EnableWindow(true);}
    if( auth2 == "authentication failure\n(you
failure!)"){ if (enable_menu == 1) {
        menu = GetMenu();
        sub = menu->GetSubMenu(0);
        sub->EnableMenuItem (ID_MENUIITEMSetting,
MF_BYCOMMAND | MF_ENABLED);
        sub->EnableMenuItem (ID_MENUIITEMload, MF_BYCOMMAND |
MF_ENABLED );
        enable_menu = 0; }}}

```

5 Conclusion

For the security challenge of current office network, TNC architecture was introduced in this paper. Within the idea of trusted office network presented, this paper uses and extends 802.1X via the trusted terminal ideology and implements the trusted authentication system of the secure access without changing present access method in the office network environment. A TNC-eSA prototype is under development. We have implemented the part of our solution based on the open source project. More works will appear in our future research.

Acknowledgments. This paper is supported by the scientific innovation talents Foundation of Henan, China under Grant No. 104100510025; Education department financing projects of Fujian, China under Grant No. JA09046.5

References

1. CNCERT/CC2010 first half year internet security survey,
<http://www.cert.org.cn/UserFiles/File/2010%20first%20half.pdf>
2. Chang, S., Yang, C.: Security of Office Management Information System Analysis. In: 2010 Second UTA International Conference on Geoscience and Remote Sensing, pp. 179–181 (2010)
3. Yin, Z., Zhang, L.: Study on Security Strategy of Wireless Mobile Office System. In: 2009 First International Workshop on Education Technology and Computer Science, pp. 495–498 (2009)
4. TCG Specification Trusted Network Connect-TNC Architecture for Interoperability Revision 1.4 (2009)
5. Freeradius, <http://wiki.freeradius.org/Main-Page>
6. Openlx open source project, <http://openlx.sourceforge.net/>

Author Index

- Abdullah, Zailani 51, 500, 592
Ahmad, Noraziah 508
- Bao, Lili 307
Barjini, Hassan 299
Bu, Youjun 684
- Cai, Ying 560
Cao, Hanyue 175
Cao, Xuyang 315
Cao, Zhengjun 175
Cavalcante, Tarcísio Pequeno 648
Chen, Bo 468
Chen, Guohua 323
Chen, Jie 236
Chen, Ruei-Chang 136
Chen, Shuai-Shuai 492
Chen, Shun 569
Chen, Wenying 331
Chen, Xiang-yu 261
Chen, Xiao 104
Chen, Xiaofeng 405
Chen, Xiao-Yue 569
Chen, Ying 128
Cheng, An-Chun 569
Cui, Gang 429
- Deng, Jian 468
Deng, Rui 608
Deris, Mustafa Mat 51, 397, 500,
508, 592
Dong, Linlin 707
Dou, Xue 203
Du, Guanglong 692
Du, Jiao 167
Du, Juan 277
Du, Qiao-qiao 261
Du, Yusong 37
- Fauzi, Ainul Azila Che 584
Feng, Ganzhong 343
Feng, Jing 624
Fujii, Mizuki 632
- Gan, Mingxin 203, 600
Gao, Ge 261
Gao, Jun 677
Gao, Shunde 315, 484
Guan, Wei 144
Gui, Zhanji 112, 120
Guo, Haimin 608
Guo, Huifang 15
Guo, Pengfei 730
Guo, Ping 365, 492
Guo, Wensheng 616, 640
Guo, Yuan-bo 770
- Han, Bo 45
Han, Xiaoting 23
He, Kaida 516
He, Qiudong 516
He, Shi wei 75
He, Ying 365, 492
He, Zhiming 343
Herawan, Tutut 51, 381, 397, 500, 508,
584, 592
Hou, Jun 88
Hou, Lei 284
Hou, Zhenjie 152, 454
Hu, Hao 437
Hu, Jianbin 229
Hu, Jinghong 82
Hu, Lin 189
Hu, Na 437
Huang, Junsheng 152
Huang, Peijie 413
Huang, Wei 104
Huang, Xipei 421
Huang, Yue 600
Huang, Zhiqin 350
Hui, Xingjie 545
Hui, Yuan 307
Hung, William N.N. 616, 640
- Ibrahim, Hamidah 299
- Ji, Junzhong 445
Jia, Ren-Yong 569
Jiang, Dailin 413

- Jiang, Haixin 269
 Jiang, Jiacheng 413
 Jiang, Kunpeng 15, 684
 Jiang, Rui 203, 600
 Jiang, Xiaoming 707
 Jiao, Lang 445
 Jin, Bin 429
 Jin, Zong-da 261
 Jing, Lin 421
 Jing, Yuan 739

 Kapelko, Rafał 477
 Khan, Nawsher 381, 397
 Khor, Kok-Chin 576
 Kuo, Yeong-Chau 136

 Lai, Wen-Hsing 98
 Lan, Julong 15, 684
 Le, Qianqi 616, 640
 Lee, Shih-Fong 136
 Lei, Xu 315
 Li, Bin 692
 Li, Dong 159
 Li, Hanling 284
 Li, Haoyu 739
 Li, Lijun 516
 Li, Ning 560
 Li, Qianmu 88
 Li, Wei 261
 Li, Weisheng 167
 Li, Xiaohao 700
 Li, Xin 31
 Li, Yana 754
 Li, Yijun 539
 Li, Yuanxiang 531
 Lin, Dezhi 284
 Lin, Piyuan 413
 Lin, Yuanshan 484
 Liu, Caiming 229
 Liu, Chengxia 560
 Liu, Chunnian 445
 Liu, Cong 722
 Liu, Fulai 196
 Liu, Gang 531
 Liu, Jianbo 331, 337
 Liu, Jun 531
 Liu, Ke Hui 75
 Liu, Li 196
 Liu, Pingping 104
 Liu, ShuaiShi 159

 Liu, Xiaotian 677
 Liu, Yanli 246
 Liu, Ying 253
 Liu, Zhijun 445
 Liu, Zhujin 365
 Liu, Zongtian 1, 253
 Lu, Yueming 221
 Luo, Ricai 211
 Luo, Zhijun 323
 Luo, Ziqiang 754
 Lv, Tianyang 59
 Lv, Yuanhai 144

 Ma, Chuan-gui 722
 Ma, Chun-guang 181
 Ma, Ji 739
 Ma, Jun 770
 Ma, Li 739
 Ma, Lingjiao 104
 Ma, Ying nan 75, 746
 Mamat, Rabiei 508
 Men, Yafeng 331
 Menezes, Andrea Carvalho 648
 Meng, Xiandong 343

 Nishizaki, Shin-ya 461, 632
 Niu, Bin 739
 Niu, Li 23
 Niu, Xiaowei 343
 Noraziah, A. 51, 381, 397, 500, 584, 592

 Othman, Mohamed 299

 Pan, Jianhui 669
 Peng, Lingxi 229
 Pinheiro, Mirian Caliope Dantas 648
 Pinheiro, Placido Rogerio 648

 Qi, Li 236
 Qi, Yong 88
 Qiu, Lin 284
 Qu, ZeHui 358

 Ren, Junling 67

 Sam, Yok-Cheng 576
 Sasajima, Takayuki 461
 Shan, Jianfang 1
 ShangGuan, Tingjie 624
 Shao, Xiaoping 715

- Shen, Dingcai 389
 Shen, Jing 37
 Shen, Ping 656
 Shi, Lan 181
 Shi, Qiugan 88
 Song, Rui 75
 Sun, Lixin 707
 Sun, Ping 75
 Sun, Yu 531
 Sun, Zhimeng 291

 Tan, Hua 211
 Tang, Weiwen 229
 Tao, Wen 337
 Tian, Yantao 159
 Ting, Choo-Yee 576
 Tong, Wei 269

 Udzir, Nur Izura 299

 Wan, Chuan 159
 Wan, Luoia 524
 Wang, Ding 181
 Wang, Guodong 746
 Wang, Jian 762
 Wang, Jinkuan 196
 Wang, Kaihua 112, 120
 Wang, Liang 429
 Wang, Ming-Shu 569
 Wang, Naisheng 663
 Wang, Qing 45
 Wang, Shutao 715
 Wang, Siou-Lin 98
 Wang, Su 373
 Wang, Weijia 284
 Wang, Xin 315, 484
 Wang, Xiukun 484
 Wang, Yanxia 545
 Wang, Ye 762
 Wang, Yourong 246
 Wang, Yu-heng 181
 Wang, Yun 350
 Wu, Di 484
 Wu, Yaping 677
 Wu, Yin 421
 Wu, Zhenfu 236

 Xia, Feng 754
 Xia, Xingyou 405
 Xia, Xuewen 389

 Xiao, Lin 437
 Xiong, Jinbo 770
 Xu, Bo 730
 Xu, Shouzhi 730
 Xu, Xiaojing 350
 Xu, Xing 437
 Xue, Youcai 669

 Yan, Dongmei 196
 Yan, Yan 120
 Yang, Guowu 616, 640
 Yang, Li-Sha 569
 Yang, Xiaofeng 739
 Yang, Yue-fang 31
 Ye, Ming 358
 Yin, Ying 59
 Ying, Weiqin 437
 You, Shibing 307, 656
 Yu, Fajiang 45
 Yu, Ruiyun 405
 Yu, Xue 307
 Yuan, Jian 762

 Zeng, Jiazhi 468
 Zeng, Jinquan 229
 Zhang, Aidong 445
 Zhang, Bin 59
 Zhang, Ce 429
 Zhang, Fang 692
 Zhang, Hong 88
 Zhang, Hui 421
 Zhang, Jian 246
 Zhang, Jianhua 152, 454
 Zhang, Laomo 746
 Zhang, Li 67
 Zhang, Liang 608
 Zhang, Ming 284
 Zhang, Pengfei 524
 Zhang, Pengwei 715
 Zhang, Shaobai 128
 Zhang, Wenxiang 112
 Zhang, Xianda 45
 Zhang, Xizhe 59
 Zhang, Yanyan 331, 337
 Zhang, Yu 754
 Zhang, Zhenwei 700
 Zhao, Guolong 59
 Zhao, Hongwei 104
 Zhao, Wenyu 677
 Zhao, Yanfeng 236

- Zhen, Li 553
Zheng, Jun 373
Zheng, Kai 373
Zhong, Qichun 82
Zhong, Yiwen 421
Zhou, Huan 730
Zhou, Lifang 167
Zhou, Yang 315
Zhu, De-Kang 569
Zhu, Lin 365
Zhu, Min 373
Zhu, Na 707
Zhu, Qin 337
Zhu, Shengping 15
Zhu, Wei 75
Zhu, Zhibin 323
Zin, Noriyani Mohd 584
Zou, Chao 221
Zou, Weixia 692