

David Taniar Osvaldo Gervasi
Beniamino Murgante Eric Pardede
Bernady O. Apduhan (Eds.)

LNCS 6018

Computational Science and Its Applications – ICCSA 2010

International Conference
Fukuoka, Japan, March 2010
Proceedings, Part III

3
Part III

 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 of Computer Science, Saarbruecken, Germany

David Taniar Osvaldo Gervasi
Beniamino Murgante Eric Pardede
Bernady O. Apduhan (Eds.)

Computational Science and Its Applications – ICCSA 2010

International Conference
Fukuoka, Japan, March 23-26, 2010
Proceedings, Part III

Volume Editors

David Taniar
Monash University, Clayton, VIC 3800, Australia
E-mail: david.taniar@infotech.monash.edu.au

Oswaldo Gervasi
University of Perugia, 06123 Perugia, Italy
E-mail: osvaldo@unipg.it

Beniamino Murgante
University of Basilicata, L.I.S.U.T. - D.A.P.I.T., 85100 Potenza, Italy
E-mail: beniamino.murgante@unibas.it

Eric Pardede
La Trobe University, Bundoora, VIC 3083, Australia
E-mail: e.pardede@latrobe.edu.au

Bernady O. Apduhan
Kyushu Sangyo University, Fukuoka 813-8503, Japan
E-mail: bob@is.kyusan-u.ac.jp

Library of Congress Control Number: 2010922807

CR Subject Classification (1998): C.2, H.4, F.2, H.3, C.2.4, F.1

LNCS Sublibrary: SL 1 – Theoretical Computer Science and General Issues

ISSN 0302-9743
ISBN-10 3-642-12178-0 Springer Berlin Heidelberg New York
ISBN-13 978-3-642-12178-4 Springer Berlin Heidelberg New York

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.

springer.com

© Springer-Verlag Berlin Heidelberg 2010
Printed in Germany

Typesetting: Camera-ready by author, data conversion by Scientific Publishing Services, Chennai, India
Printed on acid-free paper 06/3180

Preface

These multiple volumes (LNCS volumes 6016, 6017, 6018 and 6019) consist of the peer-reviewed papers from the 2010 International Conference on Computational Science and Its Applications (ICCSA2010) held in Fukuoka, Japan during March 23–26, 2010. ICCSA 2010 was a successful event in the International Conferences on Computational Science and Its Applications (ICCSA) conference series, previously held in Suwon, South Korea (2009), Perugia, Italy (2008), Kuala Lumpur, Malaysia (2007), Glasgow, UK (2006), Singapore (2005), Assisi, Italy (2004), Montreal, Canada (2003), and (as ICCS) Amsterdam, The Netherlands (2002) and San Francisco, USA (2001).

Computational science is a main pillar of most of the present research, industrial and commercial activities and plays a unique role in exploiting ICT innovative technologies. The ICCSA conference series has been providing a venue to researchers and industry practitioners to discuss new ideas, to share complex problems and their solutions, and to shape new trends in computational science.

ICCSA 2010 was celebrated at the host university, Kyushu Sangyo University, Fukuoka, Japan, as part of the university's 50th anniversary. We would like to thank Kyushu Sangyo University for hosting ICCSA this year, and for including this international event in their celebrations. Also for the first time this year, ICCSA organized poster sessions that present on-going projects on various aspects of computational sciences.

Apart from the general track, ICCSA 2010 also included 30 special sessions and workshops in various areas of computational sciences, ranging from computational science technologies, to specific areas of computational sciences, such as computer graphics and virtual reality. We would like to show our appreciation to the workshops and special sessions Chairs and Co-chairs.

The success of the ICCSA conference series, in general, and ICCSA 2010, in particular, was due to the support of many people: authors, presenters, participants, keynote speakers, session Chairs, Organizing Committee members, student volunteers, Program Committee members, Steering Committee members, and people in other various roles. We would like to thank them all. We would also like to thank Springer for their continuous support in publishing ICCSA conference proceedings.

March 2010

Oswaldo Gervasi
David Taniar

Organization

ICCSA 2010 was organized by the University of Perugia (Italy), Monash University (Australia), La Trobe University (Australia), University of Basilicata (Italy), and Kyushu Sangyo University (Japan)

Honorary General Chairs

Takashi Sago	Kyushu Sangyo University, Japan
Norio Shiratori	Tohoku University, Japan
Kenneth C.J. Tan	Qontix, UK

General Chairs

Bernady O. Apduhan	Kyushu Sangyo University, Japan
Oswaldo Gervasi	University of Perugia, Italy

Advisory Committee

Marina L. Gavrilova	University of Calgary, Canada
Andrès Iglesias	University of Cantabria, Spain
Tai-Hoon Kim	Hannam University, Korea
Antonio Laganà	University of Perugia, Italy
Katsuya Matsunaga	Kyushu Sangyo University, Japan
Beniamino Murgante	University of Basilicata, Italy
Kazuo Ushijima	Kyushu Sangyo University, Japan (ret.)

Program Committee Chairs

Oswaldo Gervasi	University of Perugia, Italy
David Taniar	Monash University, Australia
Eric Pardede (Vice-Chair)	LaTrobe University, Australia

Workshop and Session Organizing Chairs

Beniamino Murgante	University of Basilicata, Italy
Eric Pardede	LaTrobe University, Australia

Publicity Chairs

Jemal Abawajy	Deakin University, Australia
Koji Okamura	Kyushu Sangyo University, Japan
Yao Feng-Hui	Tennessee State University, USA
Andrew Flahive	DSTO, Australia

International Liaison Chairs

Hiroaki Kikuchi	Tokay University, Japan
Agustinus Borgy Waluyo	Institute for InfoComm Research, Singapore
Takashi Naka	Kyushu Sangyo University, Japan

Tutorial Chair

Andrès Iglesias	University of Cantabria, Spain
-----------------	--------------------------------

Awards Chairs

Akiyo Miyazaki	Kyushu Sangyo University, Japan
Wenny Rahayu	LaTrobe University, Australia

Workshop Organizers

Application of ICT in Healthcare (AICTH 2010)

Salim Zabir	France Telecom /Orange Labs Japan
Jemal Abawajy	Deakin University, Australia

Approaches or Methods of Security Engineering (AMSE 2010)

Tai-hoon Kim	Hannam University, Korea
--------------	--------------------------

Advances in Web-Based Learning (AWBL 2010)

Mustafa Murat Inceoglu	Ege University (Turkey)
------------------------	-------------------------

Brain Informatics and Its Applications (BIA 2010)

Heui Seok Lim	Korea University, Korea
Kichun Nam	Korea University, Korea

Computer Algebra Systems and Applications (CASA 2010)

Andrès Iglesias	University of Cantabria, Spain
Akemi Galvez	University of Cantabria, Spain

Computational Geometry and Applications (CGA 2010)

Marina L. Gavrilova	University of Calgary, Canada
---------------------	-------------------------------

Computer Graphics and Virtual Reality (CGVR 2010)

Oswaldo Gervasi	University of Perugia, Italy
Andrès Iglesias	University of Cantabria, Spain

Chemistry and Materials Sciences and Technologies (CMST 2010)

Antonio Laganà University of Perugia, Italy

Future Information System Technologies and Applications (FISTA 2010)

Bernady O. Apduhan Kyushu Sangyo University, Japan
 Jianhua Ma Hosei University, Japan
 Qun Jin Waseda University, Japan

Geographical Analysis, Urban Modeling, Spatial Statistics (GEOG-AN-MOD 2010)

Stefania Bertazzon University of Calgary, Canada
 Giuseppe Borruso University of Trieste, Italy
 Beniamino Murgante University of Basilicata, Italy

Graph Mining and Its Applications (GMIA 2010)

Honghua Dai Deakin University, Australia
 James Liu Hong Kong Polytechnic University, Hong Kong
 Min Yao Zhejiang University, China
 Zhihai Wang Beijing JiaoTong University, China

High-Performance Computing and Information Visualization (HPCIV 2010)

Frank Dévai London South Bank University, UK
 David Protheroe London South Bank University, UK

International Workshop on Biomathematics, Bioinformatics and Biostatistics (IBBB 2010)

Unal Ufuktepe Izmir University of Economics, Turkey
 Andres Iglesias University of Cantabria, Spain

International Workshop on Collective Evolutionary Systems (IWCES 2010)

Alfredo Milani University of Perugia, Italy
 Clement Leung Hong Kong Baptist University, Hong Kong

International Workshop on Human and Information Space Symbiosis (WHISS 2010)

Takuo Suganuma Tohoku University, Japan
 Gen Kitagata Tohoku University, Japan

Wireless and Ad-Hoc Networking (WADNet 2010)

Jongchan Lee
Sangjoon Park

Kunsan National University, Korea
Kunsan National University, Korea

WEB 2.0 and Social Networks (Web2.0 2010)

Vidyasagar Potdar

Curtin University of Technology, Australia

Workshop on Internet Communication Security (WICS 2010)

José Maria Sierra Camara

University of Madrid, Spain

Wireless Multimedia Sensor Networks (WMSN 2010)

Vidyasagar Potdar
Yan Yang

Curtin University of Technology, Australia
Seikei University, Japan

Program Committee

Kenneth Adamson

Ulster University, UK

Margarita Albertí Wirsing

Universitat de Barcelona, Spain

Richard Barrett

Oak Ridge National Laboratory, USA

Stefania Bertazzon

University of Calgary, Canada

Michela Bertolotto

University College Dublin, Ireland

Sandro Bimonte

CEMAGREF, TSCF, France

Rod Blais

University of Calgary, Canada

Ivan Bleic

University of Sassari, Italy

Giuseppe Borruso

Università degli Studi di Trieste, Italy

Martin Buecker

Aachen University, Germany

Alfredo Buttari

CNRS-IRIT, France

Carlo Cattani

University of Salerno, Italy

Alexander Chemeris

National Technical University of Ukraine
“KPI”, Ukraine

Chen-Mou Cheng

National Taiwan University, Taiwan

Min Young Chung

Sungkyunkwan University, Korea

Rosa Coluzzi

National Research Council, Italy

Stefano Cozzini

National Research Council, Italy

José A. Cardoso e Cunha

Univ. Nova de Lisboa, Portugal

Gianluca Cuomo

University of Basilicata, Italy

Alfredo Cuzzocrea

University of Calabria, Italy

Ovidiu Daescu

University of Texas at Dallas, USA

Maria Danese

University of Basilicata, Italy

Pravesh Debba

CSIR, South Africa

Oscar Delgado-Mohatar

University Carlos III of Madrid, Spain

Roberto De Lotto

University of Pavia, Italy

Jean-Cristophe Desplat	Irish Centre for High-End Computing, Ireland
Frank Dévai	London South Bank University, UK
Rodolphe Devillers	Memorial University of Newfoundland, Canada
Pasquale Di Donato	Sapienza University of Rome, Italy
Carla Dal Sasso Freitas	UFRGS, Brazil
Francesco Gabellone	National Research Council, Italy
Akemi Galvez	University of Cantabria, Spain
Marina Gavrilova	University of Calgary, Canada
Nicoletta Gazzea	ICRAM, Italy
Jerome Gensel	LSR-IMAG, France
Andrzej M. Goscinski	Deakin University, Australia
Alex Hagen-Zanker	Cambridge University, UK
Muki Haklay	University College London, UK
Hisamoto Hiyoshi	Gunma University, Japan
Choong Seon Hong	Kyung Hee University, Korea
Fermin Huarte	University of Barcelona, Spain
Andrès Iglesias	University of Cantabria, Spain
Antonio Laganà	University of Perugia, Italy
Mustafa Murat	Inceoglu Ege University, Turkey
Ken-ichi Ishida	Kyushu Sangyo University, Japan
Antonio Izquierdo	Universidad Carlos III de Madrid, Spain
Daesik Jang	Kunsan University, Korea
Peter Jimack	University of Leeds, UK
Korhan Karabulut	Yasar University, Turkey
Farid Karimipour	Vienna University of Technology, Austria
Baris Kazar	Oracle Corp., USA
Dong Seong Kim	Duke University, USA
Pan Koo Kim	Chosun University, Korea
Ivana Kolingerova	University of West Bohemia, Czech Republic
Dieter Kranzmueller	Ludwig Maximilians University and Leibniz Supercomputing Centre Munich, Germany
Domenico Labbate	University of Basilicata, Italy
Rosa Lasaponara	National Research Council, Italy
Maurizio Lazzari	National Research Council, Italy
Xuan Hung Le	University of South Florida, USA
Sangyoun Lee	Yonsei University, Korea
Bogdan Lesyng	Warsaw University, Poland
Clement Leung	Hong Kong Baptist University, Hong Kong
Chendong Li	University of Connecticut, USA
Laurence Liew	Platform Computing, Singapore
Xin Liu	University of Calgary, Canada
Cherry Liu Fang	U.S. DOE Ames Laboratory, USA
Savino Longo	University of Bari, Italy
Tinghuai Ma	NanJing University of Information Science and Technology, China
Antonino Marvuglia	University College Cork, Ireland

Michael Mascagni	Florida State University, USA
Nikolai Medvedev	Institute of Chemical Kinetics and Combustion SB RAS, Russia
Nirvana Meratnia	University of Twente, The Netherlands
Alfredo Milani	University of Perugia, Italy
Sanjay Misra	Atilim University, Turkey
Asish Mukhopadhyay	University of Windsor, Canada
Beniamino Murgante	University of Basilicata, Italy
Takashi Naka	Kyushu Sangyo University, Japan
Jiri Nedoma	Academy of Sciences of the Czech Republic, Czech Republic
Laszlo Neumann	University of Girona, Spain
Belen Palop	Universidad de Valladolid, Spain
Dimos N. Pantazis	Technological Educational Institution of Athens, Greece
Luca Paolino	Università di Salerno, Italy
Marcin Paprzycki	Polish Academy of Sciences, Poland
Gyung-Leen Park	Cheju National University, Korea
Kwangjin Park	Wonkwang University, Korea
Paola Perchinunno	University of Bari, Italy
Carlo Petrongolo	University of Siena, Italy
Antonino Polimeno	University of Padova, Italy
Jacynthe Pouliot	Université Laval, France
David C. Prospero	Florida Atlantic University, USA
Dave Protheroe	London South Bank University, UK
Richard Ramaroso	Harvard University, USA
Jerzy Respondek	Silesian University of Technology, Poland
Alexey Rodionov	Institute of Computational Mathematics and Mathematical Geophysics, Russia
Jon Rokne	University of Calgary, Canada
Octavio Roncero	CSIC, Spain
Maytham Safar	Kuwait University, Kuwait
Haiduke Sarafian	The Pennsylvania State University, USA
Bianca Schön	University College Dublin, Ireland
Qi Shi	Liverpool John Moores University, UK
Dale Shires	U.S. Army Research Laboratory, USA
Olga Sourina	Nanyang Technological University, Singapore
Henning Sten	Copenhagen Institute of Technology, Denmark
Kokichi Sugihara	Meiji University, Japan
Francesco Tarantelli	University of Perugia, Italy
Jesús Téllez	Universidad Carlos III de Madrid, Spain
Parimala Thulasiraman	University of Manitoba, Canada
Giuseppe A. Trunfio	University of Sassari, Italy
Mario Valle	Swiss National Supercomputing Centre, Switzerland

Pablo Vanegas	Katholieke Universiteit Leuven, Belgium
Piero Giorgio Verdini	INFN Pisa and CERN, Italy
Andrea Vittadini	University of Padova, Italy
Koichi Wada	University of Tsukuba, Japan
Krzysztof Walkowiak	Wroclaw University of Technology, Poland
Jerzy Wasniewski	Technical University of Denmark, Denmark
Robert Weibel	University of Zurich, Switzerland
Roland Wismüller	Universität Siegen, Germany
Markus Wolff	University of Potsdam, Germany
Kwai Wong	University of Tennessee, USA
Mudasser Wyne	National University, USA
Chung-Huang Yang	National Kaohsiung Normal University, Taiwan
Albert Y. Zomaya	University of Sydney, Australia

Sponsoring Organizations

ICCSA 2010 would not have been possible without the tremendous support of many organizations and institutions, for which all organizers and participants of ICCSA 2010 express their sincere gratitude:

University of Perugia, Italy
Kyushu Sangyo University, Japan
Monash University, Australia
La Trobe University, Australia
University of Basilicata, Italia
Information Processing Society of Japan (IPSJ) - Kyushu Chapter
and with IPSJ SIG-DPS

Table of Contents – Part III

Workshop on Mobile Communications (MC 2010)

A Control Loop Reduction Scheme for Wireless Process Control on Traffic Light Networks	1
<i>Junghoon Lee, Gyung-Leen Park, In-Hye Shin, Choel Min Kim, and Sang-Wook Kim</i>	
Performance Measurement of a Dual-Channel Intersection Switch on the Vehicular Network	11
<i>Junghoon Lee, Gyung-Leen Park, In-Hye Shin, Ji-Ae Kang, Min-Jae Kang, and Ho-Young Kwak</i>	
A Rapid Code Acquisition Scheme for Optical CDMA Systems	21
<i>Youngyoon Lee, Dahae Chong, Chonghan Song, Youngpo Lee, Seung Goo Kang, and Seokho Yoon</i>	
Optimal and Suboptimal Synchronization Schemes for Ultra-Wideband Systems	31
<i>Dahae Chong, Chonghan Song, Youngpo Lee, Myungsoo Lee, Junhwan Kim, and Seokho Yoon</i>	
Partial Information Relaying with Multi-Layered Superposition Coding	42
<i>Jingyu Kim and Dong In Kim</i>	
Performance Measurement of the Hybrid Prefetch Scheme on Vehicular Telematics Networks	52
<i>Junghoon Lee, Gyung-Leen Park, Youngshin Hong, In-Hye Shin, and Sang Joon Lee</i>	
Power Control for Soft Fractional Frequency Reuse in OFDMA System	63
<i>Young Min Kwon, Ok Kyung Lee, Ju Yong Lee, and Min Young Chung</i>	
Authentication – Based Medium Access Control to Prevent Protocol Jamming: A-MAC	72
<i>Jaemin Jeung, Seungmyeong Jeong, and Jaesung Lim</i>	
Femtocell Deployment to Minimize Performance Degradation in Mobile WiMAX Systems	85
<i>Chang Seup Kim, Bum-Gon Choi, Ju Yong Lee, Tae-Jin Lee, Hyunseung Choo, and Min Young Chung</i>	

A Novel Frequency Planning for Femtocells in OFDMA-Based Cellular Networks Using Fractional Frequency Reuse	96
<i>Chang-Yeong Oh, Min Young Chung, Hyunseung Choo, and Tae-Jin Lee</i>	
Association Analysis of Location Tracking Data for Various Telematics Services	107
<i>In-Hye Shin and Gyung-Leen Park</i>	
An Efficient ICI Cancellation Method for Cooperative STBC-OFDM Systems	118
<i>Kyunghoon Won, Jun-Hee Jang, Se-bin Im, and Hyung-Jin Choi</i>	
Low-Cost Two-Hop Anchor Node-Based Distributed Range-Free Localization in Wireless Sensor Networks	129
<i>Taeyoung Kim, Minhan Shon, Wook Choi, MoonBae Song, and Hyunseung Choo</i>	
SecDEACH: Secure and Resilient Dynamic Clustering Protocol Preserving Data Privacy in WSNs	142
<i>Young-Ju Han, Min-Woo Park, and Tai-Myoung Chung</i>	
Avoidance of Co-channel Interference Using Switched Parasitic Array Antenna in Femtocell Networks	158
<i>Yeonjune Jeong, Hyunduk Kim, Byung-Sung Kim, and Hyunseung Choo</i>	
User Policy Based Transmission Control Method in Cognitive Wireless Network	168
<i>Noriki Uchida, Yoshitaka Shibata, and Kazuo Takahata</i>	
Workshop on Rough and Soft Sets Theories and Applications (RSSA 2010)	
Development of a Hybrid Case-Based Reasoning for Bankruptcy Prediction	178
<i>Rong-Ho Lin and Chun-Ling Chuang</i>	
The Time Series Image Analysis of the HeLa Cell Using Viscous Fluid Registration	189
<i>Soichiro Tokuhisa and Kunihiko Kaneko</i>	
Matrices Representation of Multi Soft-Sets and Its Application	201
<i>Tutut Herawan, Mustafa Mat Deris, and Jemal H. Abawajy</i>	
Clustering Analysis of Water Quality for Canals in Bangkok, Thailand	215
<i>Sirilak Areerachakul and Siripun Sanguansintukul</i>	

Workshop on Wireless and Ad Hoc Networking (WADNet 2010)

A Review of Routing Protocols for UWB MANETs	228
<i>Yahia Hasan Jazyah and Martin Hope</i>	
An Efficient and Reliable Routing Protocol for Wireless Mesh Networks	246
<i>Jaydip Sen</i>	
A Context-Aware Service Model Based on Workflows for u-Agriculture	258
<i>Yongyun Cho, Jongbae Moon, and Hyun Yoe</i>	
A History-Based Scheduler for Dynamic Load Balancing on Distributed VOD Server Environments	269
<i>Jongbae Moon, Hyun-joo Moon, and Yongyun Cho</i>	
A Secure Routing Protocol for Wireless Sensor Networks	277
<i>Jaydip Sen and Arijit Ukil</i>	
Efficient Pairwise Key Establishment Scheme Based on Random Pre-distribution Keys in WSN	291
<i>Hao Wang, Jian Yang, Ping Wang, and Pu Tu</i>	
Agent Based Approach of Routing Protocol Minimizing the Number of Hops and Maintaining Connectivity of Mobile Terminals Which Move One Area to the Other	305
<i>Kohei Arai and Lipur Sugiyanta</i>	
Ensuring Basic Security and Preventing Replay Attack in a Query Processing Application Domain in WSN	321
<i>Amrita Ghosal, Subir Halder, Sanjib Sur, Avishek Dan, and Sipra DasBit</i>	

Workshop on Wireless Multimedia Sensor Networks (WMSN 2010)

A Review of Redundancy Elimination Protocols for Wireless Sensor Networks	336
<i>Babak Pazand, Amitava Datta, and Rachel Cardell-Oliver</i>	
A Light-Weighted Misused Key Detection in Wireless Sensor Networks	352
<i>Young-Ju Han, Min-Woo Park, Jong-Myoung Kim, and Tai-Myoung Chung</i>	

Identifying Mudslide Area and Obtaining Forewarned Time Using AMI Associated Sensor Network	368
<i>Cheng-Jen Tang and Miao Ru Dai</i>	

General Track on Information Systems and Information Technologies

Utilization of Ontology in Health for Archetypes Constraint Enforcement	380
<i>Anny Kartika Sari, Wenny Rahayu, and Dennis Wollersheim</i>	
Time-Decaying Bloom Filters for Efficient Middle-Tier Data Management	395
<i>Kai Cheng</i>	
Soft Decision Making for Patients Suspected Influenza	405
<i>Tutut Herawan and Mustafa Mat Deris</i>	
Personal Identification by EEG Using ICA and Neural Network	419
<i>Preecha Tangkraingki, Chidchanok Lursinsap, Siripun Sanguansintukul, and Tayard Desudchit</i>	
A Formal Concept Analysis-Based Domain-Specific Thesaurus and Its Application in Document Representation	431
<i>Jih-Charng Jehng, Shihchieh Chou, and Chin-Yi Cheng</i>	
On the Configuration of the Similarity Search Data Structure D-Index for High Dimensional Objects	443
<i>Arnoldo José Müller-Molina and Takeshi Shinohara</i>	
Automatic Chinese Text Classification Using N-Gram Model	458
<i>Show-Jane Yen, Yue-Shi Lee, Yu-Chieh Wu, Jia-Ching Ying, and Vincent S. Tseng</i>	
Genetic Algorithms Evolving Quasigroups with Good Pseudorandom Properties	472
<i>Václav Snášel, Jiří Dvorský, Eliška Ochodková, Pavel Krömer, Jan Platoš, and Ajith Abraham</i>	
Software Openness: Evaluating Parameters of Parametric Modeling Tools to Support Creativity and Multidisciplinary Design Integration . . .	483
<i>Flora Dilys Salim and Jane Burry</i>	
Dynamic and Cyclic Response Simulation of Shape Memory Alloy Devices	498
<i>Yutaka Toi and Jie He</i>	
Adaptive Fuzzy Filter for Speech Enhancement	511
<i>Chih-Chia Yao and Ming-Hsun Tsai</i>	

Risk Prediction for Postoperative Morbidity of Endovascular Aneurysm Repair Using Ensemble Model	526
<i>Nan-Chen Hsieh, Chien-Hui Chan, and Hsin-Che Tsai</i>	
Further Results on Swarms Solving Graph Coloring	541
<i>Manuel Graña, Blanca Cases, Carmen Hernandez, and Alicia D'Anjou</i>	
Data Collection System for the Navigation of Wheelchair Users: A Preliminary Report	552
<i>Yasuaki Sumida, Kazuaki Goshi, and Katsuya Matsunaga</i>	
Author Index	565

A Control Loop Reduction Scheme for Wireless Process Control on Traffic Light Networks^{*}

Junghoon Lee¹, Gyung-Leen Park¹, In-Hye Shin¹,
Choel Min Kim^{2,**}, and Sang-Wook Kim³

¹ Dept. of Computer Science and Statistics, Jeju National University

² Dept. of Computer Education, Jeju National University

³ College of Information and Communications, Hanyang University
{jhlee,glpark,ihshin76,chkim}@jejunu.ac.kr, wook@hanyang.ac.kr

Abstract. This paper designs a loop reduction scheme and measures its performance for the wireless process control application running on the grid-style traffic light network, aiming at improving the response time of the control system. Based on the WirelessHART standard, which assigns a slot to each (sender, receiver) pair according to the routing schedule, the allocation scheme puts all the pairs having no interference to a single slot, reducing the loop length. For further reduction, the classic Dijkstra's shortest path algorithm is modified such that the number of end-to-end paths starting from the horizontal links and the vertical link, respectively, is almost same. The transmission of the controller, which initiates all message delivery and is the main bottleneck point, will not be blocked. The simulation result demonstrates that the proposed scheme can significantly reduce the control loop length, and the modified path finding algorithm further achieves about 9.8 % improvement, just sacrificing the 3.2 % of transmission success ratio.

1 Introduction

Vehicular telematics networks can keep the vehicle connected to a network even while it is on the move. This network exploits matured wireless communication technologies and essentially includes static elements which provide the access point to moving vehicles. Correspondingly, each component communicates in two-level hierarchy in vehicular networks [1]. Level 1 corresponds to connection between static nodes and level 2 between a static node and moving vehicles trying to access the network. Level 1 communication creates a kind of wireless mesh networks which have been intensively researched and widely deployed into diverse field areas, exploiting commonly available protocols such as IEEE 802.11, Zigbee, and the like [2]. Based on this network, many vehicular applications can be developed and serviced.

^{*} This research was supported by the MKE, Korea, under the ITRC support program supervised by the NIPA. (NIPA-2009-C1090-0902-0040).

^{**} Corresponding author.

Traffic lights are found in every street, and they can desirably install a wireless communication interface as they have sufficient power provision and their locations are highly secure. On such a traffic light network, it is possible to implement a monitor-and-control application when the network includes sensors and actuators [3]. Practically, a lot of traffic-related or environmental devices such as speed detectors, pollution meters, and traffic signal controllers, can be covered by the vehicular network. Moreover, vehicles can also carry a sensor device and report the collected sensor data to a static node when connected. The message of process control applications is necessarily time critical and how to schedule messages is the main concern. It depends on the node distribution, namely, topology, and the communication protocol. First, traffic lights are placed in each intersection of Manhattan-style road networks, so the traffic light network has grid topology. Additionally, the network protocol must provide predictable network access.

The WirelessHART standard provides a robust wireless protocol for various process control applications [4]. First of all, the protocol standard exploits the slot-based access scheme to guarantee a predictable message delivery, and each slot is assigned to the appropriate (sender, receiver) pair. In addition, for the sake of overcoming the transient instability of wireless channels, a special emphasis is put on reliability by mesh networking, channel hopping, and time-synchronized messaging. This protocol has been implemented and is about to be released to the market [5]. It can accommodate a split-merge operation to mask channel errors [6] as well as an efficient routing scheme to find the path that is most likely to successfully deliver the message [7].

For the timely control action, the network must deliver the sensor and control message not only reliably but also timely, or as fast as possible. Thus, every message transmission schedule, or slot assignment, is decided in priori. Moreover, to speed up the response time of a control process, the length of the slot schedule must be minimized. In the mesh network, more than one transmission can be accomplished simultaneously as long as transmitters are far way enough not to interfere each other. If the node is equipped with an directional antenna, more transmissions can be overlapped. However, the slot assignment permitting multiple transmissions is a NP hard problem, which has no polynomial time solution. In this regard, this paper is to design and evaluate the performance of a slot allocation scheme to the the sensor and control messages for the grid-style traffic light control network.

The paper is organized as follows: After defining the problem in Section 1, Section 2 introduces the background of this paper focusing on the WirelessHART protocol and slot allocation schemes. Section 3 designs control loop reduction schemes on the target network. The simulation result is discussed in Section 4, and finally Section 5 summarizes and concludes this paper.

2 Background and Related Work

The WirelessHART standard is defined over the IEEE 802.15.4 GHz radioband physical link, allowing up to 16 frequency channels spaced by 5 MHz guard

band [8]. The link layer provides deterministic slot-based access on top of the time synchronization primitives carried out continuously during the whole network operation time. According to the specification, the size of a single time slot is 10 *ms*, and a central controller node coordinates routing and communication schedules to meet the robustness requirement of industrial applications. For more reliable communication, CCA (Clear Channel Assessment) [9] before each transmission and channel blacklisting is exploited to avoid specific area of interference and also to minimize interference to others. In some WirelessHART implementation, the network manager traverses a communication graph or a grid by the breadth-first search and allocates slots according to this order [8].

In each slot, a sender can try another channel if the CCA result of the channel on the primary schedule is not clear, even though such an operation is not yet defined in the current standard. For each destination node, a controller may reserve two alternative paths having sufficient number of common nodes. Two paths split at some nodes and meet again at other nodes. When two paths split, a node can select the path according to the CCA result in a single slot by switching to the channel associated with the secondary route. When two paths merge, the node can receive from two possible senders by a timed switch operation. Besides, the WirelessHART protocol can be reinforced by many useful performance enhancement schemes such as the virtual-link routing scheme that combines split-merge links and estimates the corresponding error rate to apply the shortest path algorithm.

Wireless mesh networks are cost-effective solutions for ubiquitous high-speed services, and its performance depends on the routing strategy. Routing has been extensively studied in wireless mesh networks, and most routing schemes are based on the shortest path algorithm [10]. The path cost is different according to the main goal of each routing scheme, including actual distance, transmission rate, and error characteristics. In most wireless process control systems, the central routing mechanism is generally taken, as not only the traffic characteristics are accurately known, but also each network access is predictable. Hence, it is possible to jointly allocate channels, route a message set, and schedule each transmission [11]. In addition, the routing procedure is executed repeatedly according to the link condition change and the schedule is distributed to each node [12].

3 Routing and Scheduling Scheme

3.1 System Model

The traffic lights form a grid network in modern cities, as a traffic light node is placed at each crossing of the Manhattan-style road network, as shown in Figure 1(a) [13]. Each node can exchange messages directly with its vertical and horizontal neighbors. Two nodes in the diagonal of a rectangle do not have a direct connection, as there may be obstacles like a tall building that blocks the wireless signal propagation. In this network, the central controller is assumed to be located at the fringe of a rectangular area, for this architecture makes the

determination of the communication schedule simple and systematic. In Figure 1(a), $T_{0,0}$ is the controller node. Any grid network can be transformed into this network by partition and rotation. In the example of Figure 1(b), four 4×4 grids are generated and each of them can be mapped to a grid shown in Figure 1(a), regardless of the grid dimension.

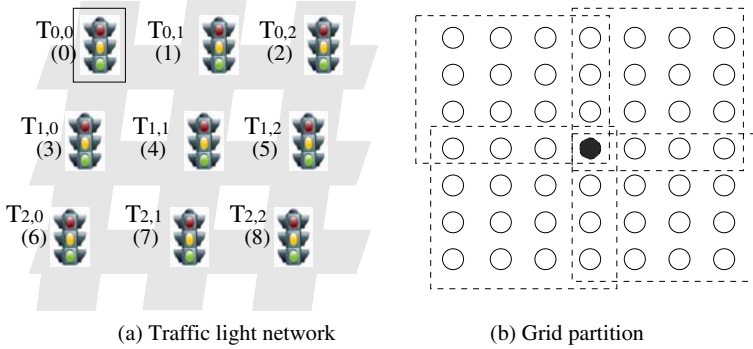


Fig. 1. Traffic light network

A control loop consists of three phases, namely, reading state variables from sensors, deciding the control action, and sending the value of control variables to the actuators. The first and third steps involve the network and have their own communication schedules.

3.2 Loop Reduction

The allocation step basically begins with the well-known Dijkstra's shortest path algorithm which can find the optimal path from a node to the other nodes in the grid given as a cost graph. In this graph, each link can be specified with a cost value such as error rate, available network bandwidth, or pass time. As this paper focuses on the reliability of each end-to-end path, we assume that the slot error rate is given as a link cost. Each link has its own error characteristics due to different power level, obstacle distribution, and so on. The change of link error characteristics can be estimated in many ways [14], but we assume that the probing result is always correct, as the correctness of channel probing is not our concern. Figure 1(a) can be modeled as a cost graph for the 3×3 grid consist of 9 nodes from $T_{0,0}$ to $T_{2,2}$ and 12 links connecting each adjacent nodes. Node 0 and $T_{0,0}$ point the same node, and this paper will selectively use an appropriate notation.

Now, Dijkstra's (one-to-many) shortest path algorithm can find the best route having the lowest error rate. First, the error rate of each link can be replaced by the success probability by subtracting it from 1.0. Actually, they can be used interchangeably. Then, each iteration compares the product of success probabilities of a path extended by a expansion node, instead of the sum of link costs as

in the original algorithm. In each node scan, the node having the highest success probability will be selected. The path from the controller, $T_{0,0}$, to each node is shown in Figure 2(a). The number in the parenthesis means the hop length of the path. Each column is a time slot and each slot has a (sender, receiver) pair. As each node has just one transmitter and one receiver module, neither the same sender nor the same receiver can appear twice in the same column. So, the initial allocation should be modified to avoid such a conflict.

To this end, the scheduler checks if two entries interfere. They interfere if 1) two senders are the same, or 2) two receivers are the same, or 3) the sender of one slot and the receiver of the other slot are the same. This criteria is valid just for the directional antenna which can increase spatial reuse and reduce interference by directing radio beam towards a desired direction [15]. When the node transmits using an omnidirectional antenna, which radiates or receives electromagnetic energy in all directions, the probability of interference increases, as the two transmissions will collide when the two senders are less than two hops away from each other. If two transmissions are not compatible in the same slot, the one with a shorter path will be moved one slot back, as the longer path is more critical to the whole control loop length. Such a move-back procedure is repeated until there is no change. Figure 2 (b) show the final result. In this figure, 0-3 means $T_{0,0} \rightarrow T_{1,0}$.

Slot	1	2	3	4	Slot	1	2	3	4	5	6	7	8	9	10	11
8 (4)	0-3	3-6	6-7	7-8	8 (4)	0-3	3-6	6-7	7-8							
5 (3)	0-3	3-4	4-5		5 (3)			0-3	3-4	4-5						
7 (3)	0-3	3-6	6-7		7 (3)					0-3	3-6	6-7				
2 (2)	0-1	1-2			2 (2)		0-1	1-2								
4 (2)	0-3	3-4			4 (2)							0-3	3-4			
6 (2)	0-3	3-6			6 (2)									0-3	3-6	
1 (1)	0-1				1 (1)				0-1							
3 (1)	0-3				3 (1)											0-3

(a) Initial setting

(b) Final allocation

Fig. 2. Traffic light network

3.3 Routing Scheme

In the control loop scenario, traffic goes from and to $T_{0,0}$. Namely, each node sends and receives a message to and from the controller node once in the control period specified by the system requirement. Even if it is desirable to take the route which has the minimum number of hops to the destination, another detour can be advantageous in terms of delivery ratio and transmission delay. In process control applications, transmission reliability is most important. However, in addition to the end-to-end performance, it is also necessary to reduce the control message exchange time. In the example of Figure 2(b), it takes 11 slots to complete the delivery of all control messages to every node in a control round.

The control loop can be reduced by overlapping transmissions. As long as the network coordinator sticks to the classic shortest path algorithm, it can hardly expect more loop reduction. As shown in Figure 2 (a), each end-to-end route to a node starts from either $T_{0,0} \rightarrow T_{0,1}$ or $T_{0,0} \rightarrow T_{1,0}$ for the 3×3 grid. If all routes start from the same link, say, $T_{0,0} \rightarrow T_{0,1}$ the next slot cannot be overlapped as two transmission will inevitably collide. The first one is from the controller to $T_{0,1}$ (the receiver of the newly starting transmission) and the second one is from $T_{0,1}$ (the sender of the message arrived in the previous slot) to some other nodes. Table 1 describes this situation. For two destinations, namely, *Dest 1* and *Dest 2*, suppose that *Dest 1* takes the route $T_{0,0} \rightarrow T_{0,1} \rightarrow T_{0,2}$ while *Dest 2* $T_{0,0} \rightarrow T_{0,1} \rightarrow T_{1,1}$, as show in Table 1.

Table 1. Example assignment

(a) Invalid assignment			
	t	$t + 1$	$t + 2$
Dest 1	$T_{0,0} \rightarrow T_{0,1}$	$T_{0,1} \rightarrow T_{0,2}$	
Dest 2		$T_{0,0} \rightarrow T_{0,1}$	$T_{0,1} \rightarrow T_{1,1}$

(b) Valid assignment			
	t	$t + 1$	$t + 2$
Dest 1	$T_{0,0} \rightarrow T_{0,1}$	$T_{0,1} \rightarrow T_{0,2}$	
Dest 2		$T_{0,0} \rightarrow T_{1,0}$	$T_{1,0} \rightarrow T_{1,1}$

As $T_{0,1}$ appears twice in slot $t + 1$, one as a sender in the first row, and the other as a receiver in the second row. On the contrary, if the first link for *Dest 2* is $T_{0,0} \rightarrow T_{1,0}$ as in Table 1(b), and two transmissions can be done simultaneously on slot $t + 1$, reducing the control loop length. In case of Figure 2(a), 6 destinations takes $T_{0,0} \rightarrow T_{1,0}$ while just 2 destinations takes $T_{0,0} \rightarrow T_{0,1}$ for the first step. So, the length reduction is not so significant. To make the number of two first links as close as possible, the network can sacrifice the optimality of the shortest path algorithm.

We classify grid nodes into 3 groups. The first group, G_D , has the nodes on the diagonal from top left to bottom right. The second group, G_L , has the nodes on the lower triangle of the grid, while the third group, G_U , has the node on the upper triangle. Accordingly, each group can be specified as follows:

$$\begin{aligned}
 G_D &: \{ T_{i,j} \mid i = j \} \\
 G_L &: \{ T_{i,j} \mid i > j \} \\
 G_U &: \{ T_{i,j} \mid i < j \}
 \end{aligned}$$

For nodes in G_L , we make every route start from the link $T_{0,0} \rightarrow T_{1,0}$, and for those in G_U , from the link $T_{0,0} \rightarrow T_{0,1}$. G_D nodes can take the route according to the normal Dijkstra's algorithm. For G_L nodes, after setting the slot error rate of $T_{0,0} \rightarrow T_{0,1}$ to 1.0, the scheduler runs the shortest path algorithm, and vice versa for the G_U nodes. Likewise, the number of appearances for $T_{0,0} \rightarrow T_{0,1}$ and

$T_{0,0} \rightarrow T_{1,0}$ becomes almost even. While this allocation cannot find the optimal path to each node, it can maximize the simultaneous transmissions in each slot and reduce the control loop length.

4 Performance Measurement

This section measures the performance of the proposed loop reduction scheme via simulation using SMPL which provides abundant functions and libraries for discrete event scheduling [16]. Only the downlink graph was considered for simplicity, as uplink and downlink communications are symmetric and calculated completely in the same way. For each parameter setting, 500 sets of link error rates are generated and the success ratio and the loop length ratio are averaged.

The first experiment measures the performance of our loop reduction scheme and compares with the upper and lower bounds. This experiment is conducted for the classic Dijkstra's algorithm. First, the lower bound accounts for the case that each end-to-end transmission can start one slot after another without being blocked from the largest hop-length destination. In this case, the control loop length is ideally smallest and thus the lower bound can be calculated as in Eq. (1).

$$n \times n - 1 + \min(\text{hop count}), \quad (1)$$

where $n \times n - 1$ is the number of noncontroller nodes in the grid, and for the last message initiated by the controller, it takes as many slots as the hop counts to the destination.

Next, the upper bound of the control loop length is estimated by assuming that there is no overlapped transmission. Hence, it corresponds to the sum of total hops to each node. For the $n \times n$ grid, the number of total transmissions for a control round is calculated as in Eq. (2).

$$\sum_{i=1}^{2(n-1)} (n - |n - i - 1|) \cdot i, \quad (2)$$

where each iteration indexed by i corresponds to the hop count from the controller. Nodes having the same hop counts need the same number of slots. The number of nodes having i hops increases by one until i reaches $(n-1)$, and then decreases also by one until i reaches $2(n-1)$.

Figure 3 shows the comparison results according to the grid dimension and the slot error rate. First, in Figure 3(a), the average slot rate is 0.1, distributing exponentially, while the grid dimension ranges from 3 to 15. The curves plots the ratio of each loop length to the upper bound. Hence, the curve marked as *UpperBound* always remain 1.0. The ratio of the packet length to the upper bound is just 5.7 % larger at maximum, compared with the lower bound to upper bound ratio. When the dimension is 3, both ratios are around 0.58, as nodes are not sufficiently apart from each other and not so many transmissions can be overlapped. The gap between 2 curves reaches 5.7 % at the grid dimension of 10,

but decreases afterwards. It indicates that our scheme can achieve almost ideal loop reduction without complex and time-consuming space search iterations. In addition, Figure 3(b) shows the effect of the slot error rate, which can lead to the variation in the length of each end-to-end path. Here, the grid dimension is set to 5. Even though the gap between the lower bound and our packed scheme gets smaller according to the increase of the slot error rate, its effect is insignificant.

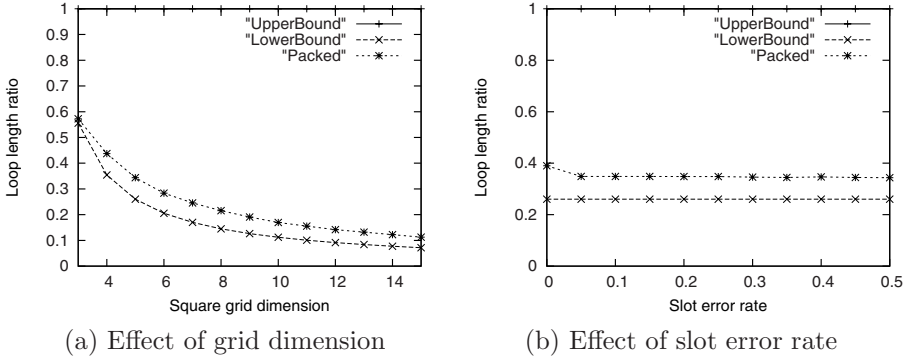


Fig. 3. Loop length reduction

Figure 4 demonstrates the performance of the proposed routing scheme, which sacrifices a little bit the path optimality for loop reduction. Here, the average slot error rate is again set to 1.0 and the dimension is changed from 3 to 15. Figure 4(a) measures how much the success ratio is lost. For all square grid dimension range, the difference is less than 3.2 %. Figure 4(b) measures how much loop reduction we gain. The improvement reaches 9.8 % at the dimension of 3 and decreases up to 3.2 % when the dimension is 15.

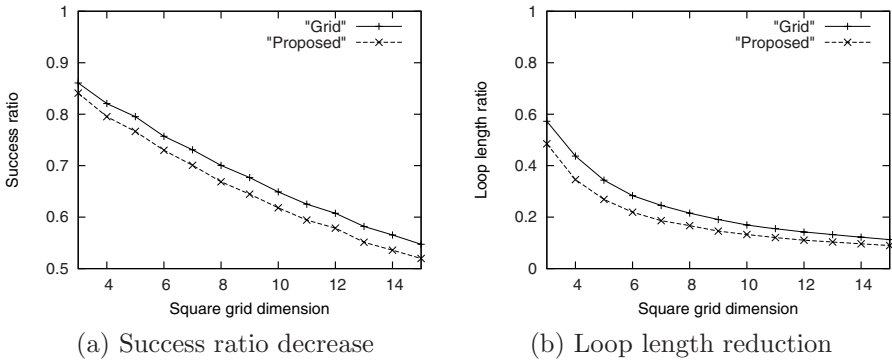


Fig. 4. Loop reduction vs. success ratio according to grid dimension

Finally, Figure 5 shows the effect of the slot error rate in the proposed routing scheme. In this experiment, the grid dimension is also set to 5, while the slot error rate changes from 0 to 0.5. The success ratio gap increases along with the error rate, reaching 4.4 % at maximum. However, the proposed scheme can achieve consistently about 8 % loop reduction in all error ranges.

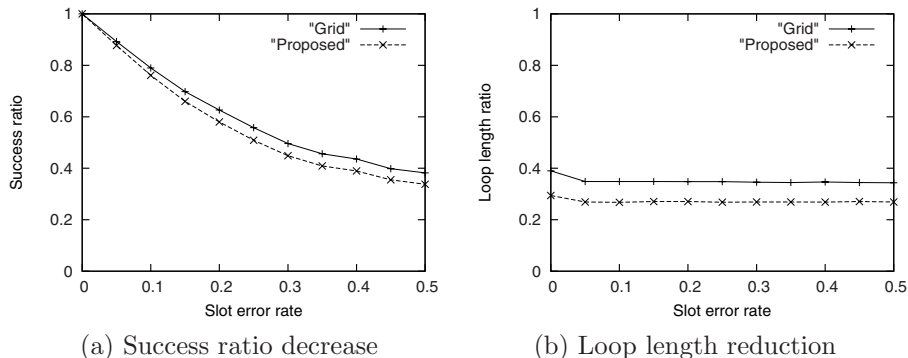


Fig. 5. Loop reduction vs. success ratio according to slot error rate

5 Conclusion

This paper has designed a loop reduction scheme and measured its performance for the wireless process control application running on the grid-style traffic light network, aiming at improving the response time of the control system. Based on the WirelessHART standard, which assigns a slot to each (sender, receiver) pair according to the routing schedule, the allocation scheme puts all the pairs having no interference to a single slot, reducing the loop length. For further reduction, the classic Dijkstra's shortest path algorithm is modified such that the number of end-to-end paths starting from the horizontal links and the vertical link, respectively, is almost same. The transmission of the controller, which initiates all message delivery and is the main bottleneck point, will not be blocked. The simulation result demonstrates that the proposed scheme can significantly reduce the control loop length, and the modified path finding algorithm further achieves about 9.8 % improvement, just sacrificing the 3.2 % of transmission success ratio.

As for future work, a fault-tolerant flooding scheme is expected to be very useful in wireless process control [17]. So, we will design a slot assignment scheme combining the split-merge operation to overcome node or link failure for the control message broadcast.

References

1. Bucciol, P., Li, F.Y., Fragoulis, N., Vandoni, L.: ADHOCYSYS: Robust and service-oriented wireless mesh networks to bridge the digital divide. In: IEEE Globecom Workshops, pp. 1–5 (2007)
2. Gislason, D.: ZIGBEE Wireless Networking. Newnes (2008)

3. IEC/PAS 62591: Industrial communication networks - Fieldbus specifications - WirelessHART communication network and communication profile (2008)
4. Hart Communication Foundation, Why WirelessHARTTM? The Right Standard at the Right Time (2007), <http://www.hartcomm2.org>
5. Han, S., Song, J., Zhu, X., Mok, A.K., Chen, D., Nixon, M., Pratt, W., Gondhalekar, V.: Wi-HTest: Compliance test suite for diagnosing devices in real-time WirelessHART network. In: The 15th IEEE Real-Time and Embedded Technology and Applications Symposium, pp. 327–336 (2009)
6. Lee, J., Shin, I., Kim, C.: Design of a reliable traffic control system on city area based on a wireless network. In: Gervasi, O., Taniar, D., Murgante, B., Laganà, A., Mun, Y., Gavrilova, M.L. (eds.) Computational Science and Its Applications – ICCSA 2009. LNCS, vol. 5592, pp. 821–830. Springer, Heidelberg (2009)
7. Lee, J., Song, H., Mok, A.K.: Design of a reliable communication system for grid-style traffic control networks. Submitted to IEEE RTAS 2010 (2010)
8. Song, S., Han, S., Mok, A.K., Chen, D., Nixon, M., Lucas, M., Pratt, W.: WirelessHART: Applying wireless technology in real-time industrial process control. In: The 14th IEEE Real-Time and Embedded Technology and Applications Symposium, pp. 377–386 (2008)
9. Ramchandran, I., Roy, S.: Clear channel assessment in energy-constrained wide-band wireless networks. IEEE Wireless Magazine, 70–78 (2007)
10. Kodialam, M., Nandagopal, T.: Characterizing the capacity region in multi-radio multi-channel wireless mesh networks. In: ACM MobiCom, pp. 73–87 (2005)
11. Wang, W., Liu, X., Krishnaswamy, D.: Robust routing and scheduling in wireless mesh networks. In: IEEE Conference on Sensor, Mesh and Ad Hoc Communications and Networks, pp. 471–480 (2007)
12. Zaidi, Z., Landfeldt, B.: Monitoring assisted robust routing in wireless mesh networks. In: Mobile Network Applications, pp. 54–66 (2008)
13. Jaap, S., Bechler, M., Wolf, L.: Evaluation of routing protocols for vehicular ad hoc networks in city traffic scenarios. In: Proceedings of the 5th International Conference on Intelligent Transportation Systems Telecommunications (2005)
14. Chang, N., Liu, M.: Optimal channel probing and transmission scheduling for opportunistic spectrum access. In: Proc. ACM International Conference on Mobile Computing and Networking, pp. 27–38 (2007)
15. Dai, H., Ng, K., Wu, M.: An overview of MAC protocols with directional antennas in wireless ad hoc networks. In: Proc. International Conference on Computing in the Global Information Technology, pp. 84–91 (2006)
16. MacDougall, M.: Simulating Computer Systems: Techniques and Tools. MIT Press, Cambridge (1987)
17. Wan, P., Huang, S., Wang, L., Wan, Z., Jia, X.: Minimum latency aggregation scheduling in multihop wireless networks. In: MobiHoc, pp. 185–193 (2009)

Performance Measurement of a Dual-Channel Intersection Switch on the Vehicular Network^{*}

Junghoon Lee¹, Gyung-Leen Park^{1,**}, In-Hye Shin¹,
Ji-Ae Kang¹, Min-Jae Kang², and Ho-Young Kwak³

¹ Dept. of Computer Science and Statistics,
Jeju National University, 690-756, Jeju Do, Republic of Korea

² Dept. of Electronic Engineering,
Jeju National University, 690-756, Jeju Do, Republic of Korea

³ Dept. of Computer Engineering,
Jeju National University, 690-756, Jeju Do, Republic of Korea
{jhlee,glpark,ihshin76,minjk,kwak}@jejunu.ac.kr, loveria2@nate.com

Abstract. This paper designs an intensively measures the performance of an efficient message switch scheme for the intersection area where routing decision may be complex due to traffic concentration, aiming at enhancing end-to-end delays and reliability for the information retrieval service in the vehicular network. Installed at the corner of an intersection, each switch node opens an external interface to exchange messages with vehicles proceeding to the intersection as well as switches the received messages via the dual-channel internal interfaces. Based on the slot-based MAC, at each slot time of internal interfaces, the sender node probes the channel status of the two different destinations and then dynamically selects the appropriate channel according to the probing result. The simulation result shows that the proposed scheme improves the delivery ratio by up to 18.5 % for the experimental channel error rate range as well as up to 8.1 % for the given network load distribution.

1 Introduction

According to the ongoing deployment of in-vehicle telematics devices and vehicular networks, wireless vehicular communications have become an important priority for car manufactures [1]. The vehicular network can be built in many different ways according to the variety of available wireless communication technologies, for example, infrastructure-based cellular networks, VANET (Vehicular Ad hoc NETWORK), and roadside networks. Moreover, not just the vehicle but also diverse communication entities are participating in this network, including traffic lights, roadside units, sensors, gas stations, and many other facilities. Correspondingly, abundant application scenarios are now possible. In addition to the standard safety-related applications such as traffic accident propagation

^{*} This research was supported by the MKE, Korea, under the ITRC support program supervised by the NIPA. (NIPA-2009-C1090-0902-0040)).

^{**} Corresponding author.

and emergency warning [2], information retrieval is very useful for drivers and passengers [3]. Specifically, a driver may want to know the current traffic condition of specific road segments lying on the way to his/her destination, query several shops to decide where to go, and check parking lot availability [4].

In the information retrieval service, queries must be delivered to the destinations and the result is sent back to the query issuer. A message proceeds to its destination according to a routing protocol each vehicle cooperatively runs. As each vehicle can move only along the road segment and static nodes such as gas stations and traffic lights are generally placed on the roadside [5], the message delivery path must trail the actual road layout [6]. Just as in the vehicle's path, intersection areas having many vehicles are important for the delay and reliability of message transmissions, because the large number of vehicles makes the routing decision very complex and raises intervehicle interference. Moreover, the carry and forward strategy is preferred to cope with the disconnection problem in the sparse vehicular network [7]. As will be described in the next section, this behavior further increases the complexity of message routing near the intersection area.

Definitely, there exists an intersection which has much more traffic than others and this hot area may be a bottleneck for both vehicle and data traffic. As an end-to-end path is highly likely to involve one or more such intersection areas, how to manage the communication in this area is very critical to the communication performance such as the transmission delay and the delivery ratio. In this regard, this paper is to design and measure the performance of a reliable and high-speed message switching scheme around the hot intersection area for the vehicular network, taking advantage of multiple channels. It is possible to create multiple channels in a cell in many available wireless protocols such as IEEE 802.11 series [2], Zigbee, and WirelessHart [8]. Existing researches have also pointed out that multiple network interfaces does not cost too much [9]. Moreover, the slot-based collision-free access scheme can be employed for predictable channel access, short delay, and reliability, while the probing-based access can efficiently reduce the effect of the channel error inherent in the wireless frequency channel.

This paper is organized as follows: After issuing the problem in Section 1, Section 2 describes the background of this paper and related work. Section 3 designs a wireless switch for intersection areas and performance measurement results are demonstrated and analyzed in Section 4. Finally, Section 5 summarizes and concludes this paper with a brief introduction of future work.

2 Background and Related Work

Jeju taxi telematics system keeps track of each taxi for the purpose of providing an efficient taxi dispatch service to customer [1]. To this end, each taxi periodically reports its GPS reading via the CDMA (Code Division Multiple Access) network, an instance of cellular networks serviced in the Republic of Korea. This system is rather an expensive vehicular network in the sense that every member taxi should pay the monthly communication fee. In the mean time, in November 2005, the U.S. Department of Transportation began to develop and test an integrated, vehicle-based crash warning system under the program of IVBBS (Integrated

Vehicle-Based Safety Systems) [10]. Here, data transmitted from the roadside to the vehicle could warn a driver that it is not safe to enter an intersection. Vehicles could serve as data collectors and anonymously transmit traffic and road condition information from every major road within the transportation network. After all, more vehicular networks are being deployed into our everyday life.

As for the wireless channel, the IEEE 802.11b standard specifies 11 channels operating in the 2.4 GHz band with 80 MHz of reusable spectrum [11]. Even though the number of simultaneous channels in a cell is limited to 3 due to the channel overlap problem, it is possible to create multiple channels from the wireless spectrum in a cell. In addition, Zigbee and Wireless HART (Highway Addressable Remote Transducer) support channels spaced by 5 MHz guard band, making it possible for each node to hop over channels to reduce the effect of channel errors [8]. The link layer is based on a TDMA (Time Division Multiple Access) style access scheme which runs on top of the time synchronization mechanism carried out continuously during the whole operation time by means of MAC PDUs. The time axis is divided into 10 *ms* time slots and a group of consecutive slots are defined to be a superframe. This access scheme can provide bounded access time for each node.

High mobility and nonuniform distribution of vehicles prevent the existing routing schemes from applying to the vehicular network. Quite a lot of researches have been recently conducted for vehicular networks to deal with such problems. Basically, the carry and forward strategy is preferred to cope with disconnection in the sparse network part [7]. Here, when a node cannot find a receiver, it stores the message in its buffer until it enters the range of a new receiver. Even though this scheme increases the transmission delay, it is better not to discard a message. This leads to a large traffic load around the intersection area. As a variant of carry and forward scheme, VADD (Vehicle-Assisted Data Delivery) exploits predictable vehicle mobility model in which vehicle movement is limited by the traffic pattern and road layout. Based on the traffic pattern estimation, a vehicle decides the best next node to forward a packet [12].

For the information retrieval method, VITP (Vehicular Information Transfer Protocol) has been proposed as an application-layer communication protocol, which is designed to support the establishment of a distributed, ad hoc service infrastructure over VANET [3]. The VITP infrastructure can be used to provide location-based, traffic-oriented services to drivers, using information retrieved from vehicular sensors and taking advantage of on-board GPS navigation systems. In this system, vehicular-service queries must be location-sensitive, specifying explicitly the target location of their inquiry. VITP communication entities follow a best-effort approach in their operation, while a VITP request is transported between protocol peers until some return condition is satisfied. The protocol peer that detects the upholding of the return condition, creates the VITP reply and posts it toward source-region and broadcast so that the issuer can receive even if it has moved. In the message exchange, many protocols can be exploited for information retrieval such as Zhao et al.'s V2VR, which select forward and backward proxies based on the mobility pattern of the vehicle [4].

3 Wireless Switch Design

Fig. 1 shows the basic idea of this paper. At the intersection, four static switch nodes from *A* to *D* are installed at each corner as shown in Fig. 1(a). The number of switch nodes is equal to the number of branches, whether they cross on the same plane or in different layers. Each switch node is bound to the vehicles heading for the intersection from the preassigned branch. For the sake of speeding up the transit time in the intersection area, each switch node opens an external interface to receive and send messages with other vehicles, while internally forwarding messages to an appropriate direction via the 2 additional channels not exposed to the vehicle. As contrast to the external interface that must follow the standard vehicular communication protocol, we can define a new protocol for the internal access. For the internal message exchange, this paper suggests that each node be connected to dual frequency channels which run slot-based MAC as shown in Fig 1. (b).

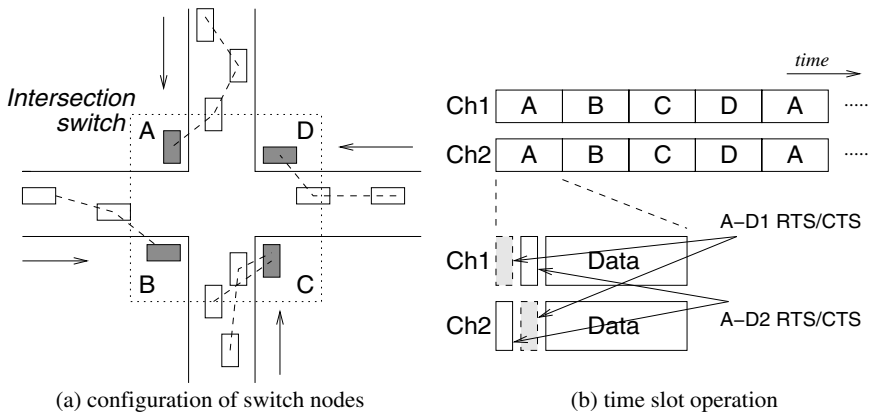


Fig. 1. Basic concept

When a vehicle having messages to send or relay approaches to an intersection, it checks if it can reach a switch node. If so, the message is sent directly to the switch node without contacting any other vehicles. Otherwise, the message is just forwarded to a vehicle in its moving direction just like in the normal ad hoc routing protocol. A switch node, receiving a message via the external interface, decides the next switch node to forward the message in its internal interface according to the final destination. In the intersection area, the complex routing among vehicles is skipped, reducing the number of relays and erroneous transmissions between moving vehicles. For this internal operation, each switch node transmits according to the TDMA manner to provide the predictable access. Even though this access demands that every node have the common clock, current GPS technology can easily achieve the global clock synchronization of reasonable accuracy [13]. So, many ad hoc networks also exploit the slot-based access scheme.

Many vehicles passing by the intersection area may cause channel interference, so an efficient method is necessary to cope with this problem. This problem can be solved by dual channels. The slots of two channels are synchronized, and two slots of the same time instance must be assigned to the same node. At the beginning of each slot, the transmitter node gets two messages having different destinations from its message queue. Our previous work has proposed a channel probing mechanism from two different nodes [14]. On the contrary, this paper makes a single node be assigned time slots on each channel and select the appropriate one. A channel probing result its corresponding action is similar in both schemes.

Let these two destinations be D_1 and D_2 , respectively. The sender probes, for example, with RTS/CTS methods, D_1 and D_2 on the two channels as shown in Fig 1 (b). The transmission on the bad channel is meaningless not solely for the power consumption aspect. It may extend the instable energy and its duration. Table 1 shows the probing result and corresponding actions. As shown in row 1, D_1 on channel 1 and also D_2 on channel 2 are both in good state, D_1 and D_2 send as scheduled. In row 2, every channel status is good except the one for D_2 on channel 2. If we switch $\langle D_1, D_2 \rangle$ to $\langle D_2, D_1 \rangle$, both nodes can successfully send their messages. Otherwise, only A can succeed. In this case, the channel switch can save one transmission loss. Row 8 describes the situation that D_1 is good only on channel 2 while D_2 also only on channel 1. Switching two channels saves 2 transmissions that might fail on the regular schedule.

Table 1.Channel status and transmission

No.	Ch1-D1	Ch2-D2	Ch1-D2	Ch2-D1	Ch1	Ch2	save
1	Good	Good	X	X	D1	D2	0
2	Good	Bad	Good	Good	D2	D1	1
3	Good	Bad	Good	Bad	D1	-	0
4	Good	Bad	Bad	X	D1	-	0
5	Bad	Good	Good	Good	D2	D1	1
6	Bad	Good	Good	Bad	-	D2	0
7	Bad	Good	Bad	X	-	D2	0
8	Bad	Bad	Good	Good	D2	D1	2
9	Bad	Bad	Good	Bad	D2	-	1
10	Bad	Bad	Bad	Good	-	D1	1
11	Bad	Bad	Bad	Bad	-	-	0

X : don't care

Each switch node maintains a message queue. The buffer space can be assumed to have no limitation, but it is desirable to discard a message which stayed in the buffer too much so as to obviate undue delay for the subsequent messages. After all, we assume that a message is automatically deleted from its queue when its waiting time exceeds the given discard interval. In addition, the routing decision outside the intersection area can be exploited by other work and

there are many available schemes especially targeting at the urban area [15]. For example, VADD can cope with the situation that a vehicle moves after issuing a retrieval request. Finally, how to assign slots is not our concern, and we just consider the round-robin style allocation and assume that each slot is assigned to a switch node. There also exist some channel allocation schemes for the give traffic requirement [16].

4 Performance Measurement

We evaluate the performance of our scheme, focusing on the effect of channel switches, via simulation using SMPL, which provides a simple and robust library for discrete event trace similar to ns-2 event scheduler [17]. The main performance metrics are the delivery ratio and the access delay. The delivery ratio measures the ratio of the number of successfully delivered messages through the intersection area to the number of total messages arrived at the 4 switch nodes. The access delay is the time interval from the instant a message arrives at a switch node to the instant it gets out of the intersection switch. The average of access delays just takes the messages that have been successfully transmitted, so it is likely to increase when the number of successful transmissions increases. After all, this section shows the effect of the channel error probability, offered load, and the discard interval, one by one. In the following experiments, the message is assumed to arrive at each switch node according to the exponential distribution. In addition, for simplicity, each message fits a single time slot, making simple to estimate the traffic load. Moreover, every time is aligned to a slot length.

The vehicle distribution is taken from the Jeju taxi telematics system which keeps track of each member taxi [1]. Fig. 2 shows a snapshot of vehicle distribution on the real road network. This history is very useful for the vehicular

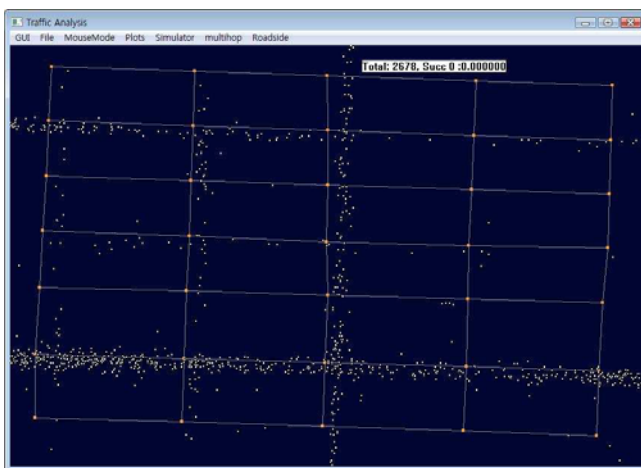


Fig. 2. Vehicle distribution at intersections

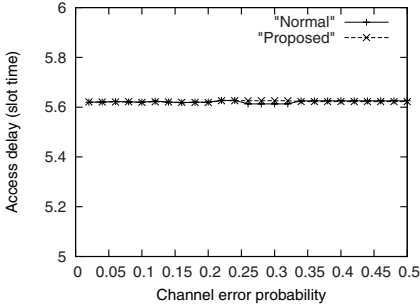


Fig. 3. Access delay at low load

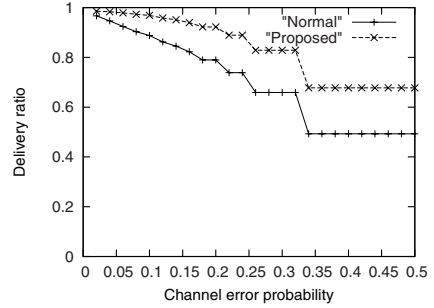


Fig. 4. Delivery ratio at low load

network design and the traffic pattern analysis. In this figure, the intersection is represented by a rectangle and the road segment is by a line, while the location of taxis during the specific time interval is marked with a small dot. The intersection area has more dots, which appear just along the road segment, even though the number of dots is different for each segment which meets at an intersection.

Fig. 3 and Fig. 4 show the effect of the channel error probability when the offered load is relatively low, namely, 0.5. The curve marked as *proposed* is the performance of the case of channel switching. The channel error probability is a probability that a channel is clean when a switch node is to send its message. This experiment assumes that a message is discarded if it is not served until 16 time slots from its arrival. Fig. 3 reveals that a message can exit the intersection area around in 5.6 time slots, that is, less than 1.5 round. In Fig. 4 the delivery ratio is almost 1.0 when the channel error rate is close to 0. At the higher error rate, we can achieve better improvement, as there are more cases channel switching is possible. When the channel error probability is greater than 0.4, the proposed scheme can improve the delivery ratio by up to 18.5 %.

Fig. 5 and Fig. 6 plot the effect of the channel error probability when the offered load is 1.0. On high load, a switch node can almost always find two

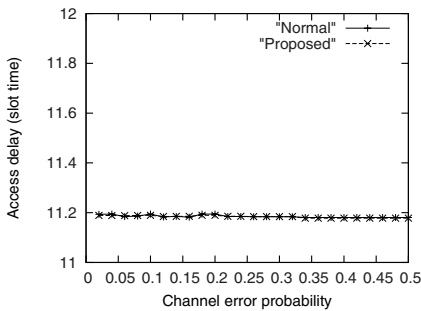


Fig. 5. Access delay at high load

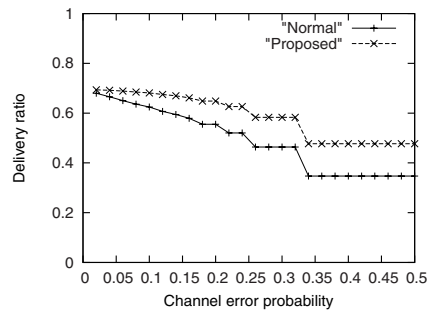


Fig. 6. Delivery ratio at high load

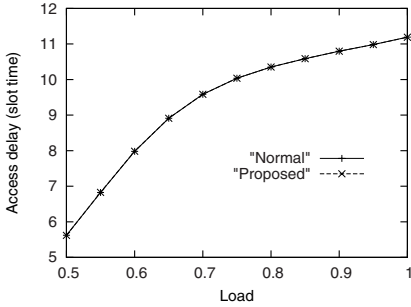


Fig. 7. Access delay vs. load

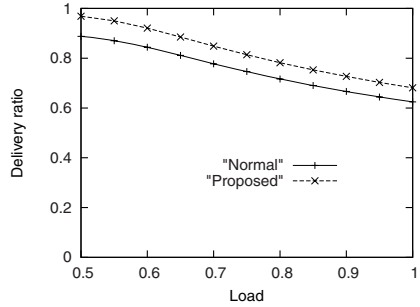


Fig. 8. Delivery ratio vs. load

messages having different destinations. As can be found in Fig. 5, the access delay for both cases stays at 11.2 time slots, namely, less than 3 rounds. This indicates that the discard interval can control the access delay quite well regardless of channel errors. When the channel error probability gets larger than 0.4, the proposed scheme can improve the delivery ratio by up to 13.0 %. Fig. 6 shows that both the delivery ratio and its performance gap are less than those in Fig. 4, indicating that the message discarded in the queue outnumbers that saved by the channel switch.

Fig. 7 and Fig. 8 demonstrate the effect of message load to the access delay and the delivery ratio, respectively. In this experiment, the channel error probability is set to 0.1 and the discard interval is fixed to 16 time slots. Fig. 7 shows that the access delays for two cases are almost same and the degree of inclination gets smaller after the load is 0.8, as more messages are discarded for the range of 0.8 through 1.0. In Fig. 8, the two curves show the constant gap for all load range, while the gap is 8.1 % at maximum, indicating that the early discard more affects the delivery ratio on higher load.

Fig. 9 and Fig. 10 show the effect of discard intervals. It can be simply expected that the shorter the discard interval, the shorter the access delay. Fig. 9 exhibits

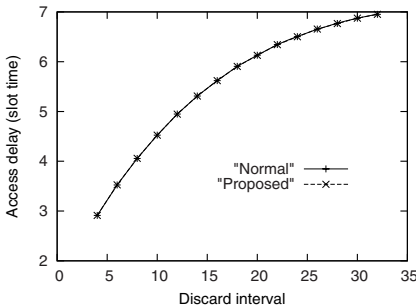


Fig. 9. Effect of discard interval

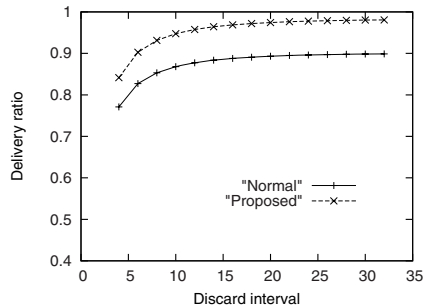


Fig. 10. Effect of discard interval

the largely linear increase of access delays according to the discard interval. Two curves also have no conceivable difference as in the other cases. However, Fig. 10 reveals that the effect of discard intervals is not so significant to the delivery ratio, less than 18 % over the interval range of 5 through 35. Two curves show a similar pattern, but the gap gets a little bit larger according to the increase of discard intervals. For all experiments, it can be pointed out by our experiment that the improved delivery ratio doesn't affect the transit time.

5 Conclusion

This paper has proposed and measured the performance of a message switch scheme for a hot intersection area in vehicular telematics network. Installed at each corner of an intersection, switch nodes cooperatively exchange messages according to the current channel status, speeding up the intersection transit time and improving the delivery ratio. The simulation result shows that the proposed scheme can improve the delivery ratio by up to 18.5 % for the given range of channel error probability, and 8.1 % for the given range of network load. This result comes from the fact that our scheme can avoid the complex routing and message interference in the intersection of intensive traffic. Next, dual channels can efficiently cope with channel errors. Even though such employment of additional equipments brings extra cost and loses the advantage of autonomous operations, access predictability, enhanced delivery speed, and better delivery ratio can compensate for the cost. In short, the wireless switch can be installed at the intersections, which have a lot of traffic especially in the urban area, to enhance the communication quality and reliability.

As future work, we are planning to assess the performance of our scheme with the real-life location history data obtained by the Jeju taxi telematics system [11] and analyze the end-to-end performance characteristics to revise our algorithm and test the diverse message scheduling policy.

References

1. Lee, J., Park, G., Kim, H., Yang, Y., Kim, P., Kim, S.: A telematics service system based on the Linux cluster. In: Shi, Y., van Albada, G.D., Dongarra, J., Sloot, P.M.A. (eds.) ICCS 2007. LNCS, vol. 4490, pp. 660–667. Springer, Heidelberg (2007)
2. Society of Automotive Engineers: Dedicated short range communication message set dictionary. Technical Report Standard J2735 (2006)
3. Dikaiakos, M., Iqbal, S., Nadeem, T., Iftode, L.: VITP: An information transfer protocol for vehicular computing. In: ACM VANET, pp. 30–39 (2005)
4. Zhao, J., Arnold, T., Zhang, Y., Cao, G.: Extending drive-thru data access by vehicle-to-vehicle relay. In: ACM VANET, pp. 30–39 (2008)
5. Lochert, C., Scheuermann, B., Wewetzer, C., Luebke, A., Mauve, M.: Data aggregation and roadside unit placement for a VANET traffic information system. In: ACM VANET, pp. 58–65 (2008)

6. Korkmaz, G., Ekici, E., Ozguner, F., Ozguner, U.: Urban multihop broadcast protocol for inter-vehicle communication systems. In: ACM VANET, pp. 76–85 (2004)
7. Davis, J., Fagg, A., Levine, B.: Wearable computers as packet transport mechanism in highly-partitioned ad-hoc networks. In: Int'l Symposium on Wearable Computing (2001)
8. Song, S., Han, S., Mok, A., Chen, D., Nixon, M., Lucas, M., Pratt, W.: Wireless HART: Applying Wireless Technology in Real-Time Industrial Process Control. In: The 14th IEEE Real-Time and Embedded Technology and Applications Symposium, pp. 377–386 (2008)
9. Naumov, V., Gross, T.: Connectivity-aware routing (CAR) in vehicular ad hoc networks. In: IEEE Infocom, pp. 1919–1927 (2007)
10. US Department of Transportation: Integrated Vehicle-Based Safety Systems (IVBSS) Phase I Interim Report. DOT HS 810 952 (2008), http://www.itsdocs.fhwa.dot.gov/JPODOCS//REPTS_TE/14434.htm
11. Kodialam, M., Nandagopal, T.: Characterizing the capacity region in multi-radio multi-channel wireless mesh networks. In: ACM MobiCom, pp. 73–87 (2005)
12. Zhao, J., Cao, G.: VADD: Vehicle-assisted data delivery in vehicular ad hoc networks. In: IEEE INFOCOM (2006)
13. Lee, J., Kang, M., Park, G., Shin, S.: Design of a reliable real-time scheduling policy for dual-channel networks. *Journal of Information Science and Engineering*, 1407–1419 (2007)
14. Lee, J., Park, G., Shin, S., Kang, M.: Design of intersection switches for the vehicular network. In: Hong, C.S., Tonouchi, T., Ma, Y., Chao, C.-S. (eds.) APNOMS 2009. LNCS, vol. 5787, pp. 523–526. Springer, Heidelberg (2009)
15. Jaap, S., Bechler, M., Wolf, L.: Evaluation of Routing Protocols for Vehicular Ad Hoc Networks in City Traffic Scenarios. In: Proceedings of the 5th International Conference on Intelligent Transportation Systems Telecommunications (2005)
16. Li, H., Shenoy, P., Ramamritham, K.: Scheduling communication in real-time sensor applications. In: Proc.10th IEEE Real-time Embedded Technology and Applications Symposium, pp. 2065–2080 (2004)
17. MacDougall, M.: *Simulating Computer Systems: Techniques and Tools*. MIT Press, Cambridge (1987)

A Rapid Code Acquisition Scheme for Optical CDMA Systems

Youngyoon Lee, Dahae Chong, Chonghan Song, Youngpo Lee,
Seung Goo Kang, and Seokho Yoon*

School of Information and Communication Engineering, Sungkyunkwan University,
300 Chunchun-dong, Jangan-gu, Suwon, Gyeonggi-do, 440-746, Korea
{news8876,lvjs1019,starsong83,leeyp204,
general185,syoon}@skku.edu

Abstract. In this paper, we propose a novel code acquisition scheme called modified multiple-shift (MMS) for optical code division multiple access (CDMA) systems. By using the modified local code whose sign is positive or negative, the proposed MMS scheme provides a shorter mean acquisition time (MAT) than the conventional multiple-shift (MS) scheme. The simulation results demonstrate that the MAT of the MMS scheme is shorter than that of the MS scheme in both single-user and multi-user environments.

Keywords: Optical code division multiple access, optical orthogonal code, multiple-shift, code acquisition, mean acquisition time.

1 Introduction

In code division multiple access (CDMA)-based systems, the data demodulation is possible only after a code synchronization is completed. Therefore, the code synchronization is one of the most important tasks in CDMA-based systems [1]. Generally, the code synchronization consists of two stages: code acquisition and tracking. Achieving the code synchronization is called code acquisition and maintaining the code synchronization is called tracking [2], of which the former is dealt with in this paper. In code acquisition process, the most significant performance measure is the mean acquisition time (MAT), which is a mean time that elapses prior to acquisition.

An optical CDMA system uses a spreading code called optical orthogonal code (OOC) proposed by Salehi [3]. Due to its good auto-correlation and cross-correlation properties, the OOC has been widely used for various CDMA-based systems including optical CDMA systems [4], [5]. Keshavarzian and Salehi introduced the serial-search (SS) [4] scheme using the OOC, which is simple; however, its MAT increases as the code length becomes longer. Thus, the SS scheme is not suitable for rapid acquisition of a long code used in multi-user environments. In order to overcome this drawback, in [5], the same authors proposed the multiple-shift (MS) scheme using the OOC, which consists of two stages and offers a shorter MAT compared with that of the SS scheme.

* Corresponding author.

In this paper, we propose a novel code acquisition scheme called modified multiple-shift (MMS). The MMS scheme comprises two stages similarly to the MS scheme, however, by using the modified local code differently from the MS scheme, it provides a shorter MAT than that of the MS scheme.

The remainder of this paper is organized as follows. Section 2 describes the system model. In Section 3 we present the conventional MS and proposed MMS schemes. Section 4 analyzes the MAT performance of the MMS scheme. In Section 5 the simulation results show the MATs of the MS and MMS schemes in single-user and multi-user environments. Section 6 concludes this paper.

2 System Model

In an optical CDMA channel, there exist various kinds of impairments such as noise, multipath signals, and multiple access interference (MAI). The influences of noise and multipath signals can be almost completely mitigated by using fiber-optic medium; however, that of MAI should be alleviated in the receiver [6], [7]. In this paper, thus, we consider a multi-user environment without the influences of noise and multipath signals. Then, the received signal $r(t)$ can be written as

$$r(t) = \sum_{n=1}^N s^{(n)}(t - \tau^{(n)}), \quad (1)$$

where $s^{(n)}(t)$ is the transmitted signal of the n -th user; $\tau^{(n)} \in [0, T_b)$ denotes the time delay of the n -th user with the bit duration T_b ; and N is the number of users. We consider the on-off-keying (OOK) modulation and assume that the bit rate is the same for all users. Thus, the transmitted signal $s^{(n)}(t)$ can be expressed as

$$s^{(n)}(t) = \sum_{i=-\infty}^{\infty} b_i^{(n)} c^{(n)}(t - iT_b), \quad (2)$$

where $b_i^{(n)}$ is the i -th binary data bit of the n -th user and $c^{(n)}(t) = \sum_{j=0}^{F-1} a_j^{(n)} p(t - jT_c)$ is the OOC of the n -th user with the chip duration T_c and sequence $a_j^{(n)} \in \{0, 1\}$ of length F and weight K (the total number of ‘1’s in $a_j^{(n)}$), where the rectangular pulse $p(t)$ over $[0, T_c]$ is defined as

$$p(t) = \begin{cases} 1, & 0 \leq t < T_c, \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

The OOC has four parameters F , K , λ_a , and λ_c , where λ_a and λ_c are auto-correlation and cross-correlation constraints, respectively [3]. In order to maintain the strict orthogonality of the OOC, both λ_a and λ_c have to be zero; however, since an OOC consists of 0 and 1, the ideal strict orthogonality cannot be satisfied. Thus, in this paper, both λ_a and λ_c are set to 1.

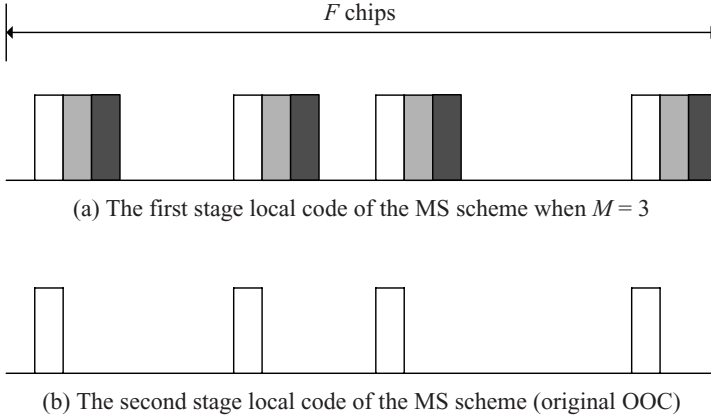


Fig. 1. The local codes of the MS scheme when OOC of (32,4,1,1) is used

3 Code Acquisition Schemes

3.1 MS Scheme

In the MS scheme, total F cells in the search space are divided into Q groups, each of which contains M cells. The relation of Q and M is given by

$$Q = \left\lceil \frac{F}{M} \right\rceil, \quad (4)$$

where $\lceil \cdot \rceil$ denotes the ceiling operation.

The MS scheme consists of two stages. In the first stage, the received signal $r(t)$ is correlated with the first stage local code shown in Fig. 1. The correlation is repeated on a group-by-group basis. If the correlation value corresponding to a certain group exceeds a given threshold $TH_{MS,first}$, the group is declared to be the correct group having the time delay $\tau^{(n)}$ and the process is transferred to the second stage. In the second stage, the correlation-based search is performed again with the second stage local code (original OOC) on a cell-by-cell basis over M cells in the correct group. As in the first stage, when the correlation value corresponding to a certain cell exceeds a given threshold $TH_{MS,second}$, which is different from $TH_{MS,first}$, the cell is declared to be an estimate of the time delay $\tau^{(n)}$.

Using the definition in [5], the MAT T_{MS} of the MS scheme is given by

$$T_{MS} = \frac{Q + 1}{2} + \frac{M + 1}{2}. \quad (5)$$

From (4) and (5), we can find that the minimum value of T_{MS} equals to \sqrt{F} when $M = \sqrt{F}$.

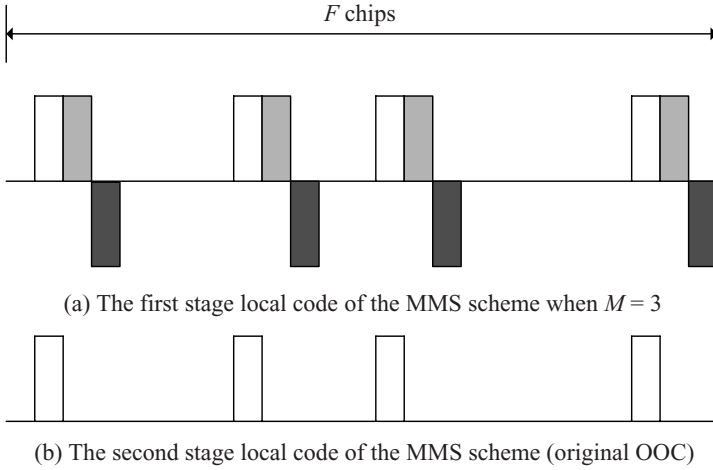


Fig. 2. The local codes of the MMS scheme when OOC of (32,4,1,1) is used

3.2 MMS Scheme

The MMS scheme proposed in this paper also consists of two stages similarly to the MS scheme. However, using the modified local code, the MMS scheme can provide a shorter MAT than that of the MS scheme.

In the first stage, the first stage local code shown in Fig. 2 is used instead of the conventional local code used for MMS scheme. The first stage local code consists of several chips whose sign is positive or negative. When M is an even number, the first half of the group template is positive and its second half is negative; otherwise, the number of positive chips is larger by one than that of negative chips.

In the first stage, the local and the received codes are correlated and the local code is updated by M chips at every operation where the correlation value is calculated. Unlike the MS scheme, the first stage of the MMS scheme finishes the search when the absolute value of the correlation value is equal to or larger than the $TH_{MMS,first}$. If the correlation value of the correct group is positive, we only need to search the first half of the M cells of the correct group in the second stage; otherwise, the search is performed over the second half of the M cells. Thus, we only need to search the maximum $\lceil M/2 \rceil$ chips regardless of the sign of the correlation value. In other words, the number of the operations in the second stage is reduced by nearly half than that of the MS scheme by using the modified local code. In the second stage of the MMS scheme, the correlation is compared to the threshold $TH_{MMS,second}$ and the maximum number of the updates of the local code is $\lceil M/2 \rceil$. Therefore, the MAT of the MMS scheme T_{MMS} is

$$T_{MMS} = \frac{Q + 1}{2} + \frac{M + 2}{4}. \tag{6}$$

We can notice from (5) and (6) that the MAT in the second stage is reduced by nearly half than that of the MS scheme and the minimum value of T_{MMS} equals $\sqrt{F/2} + 1$ when $M = \sqrt{2F}$.

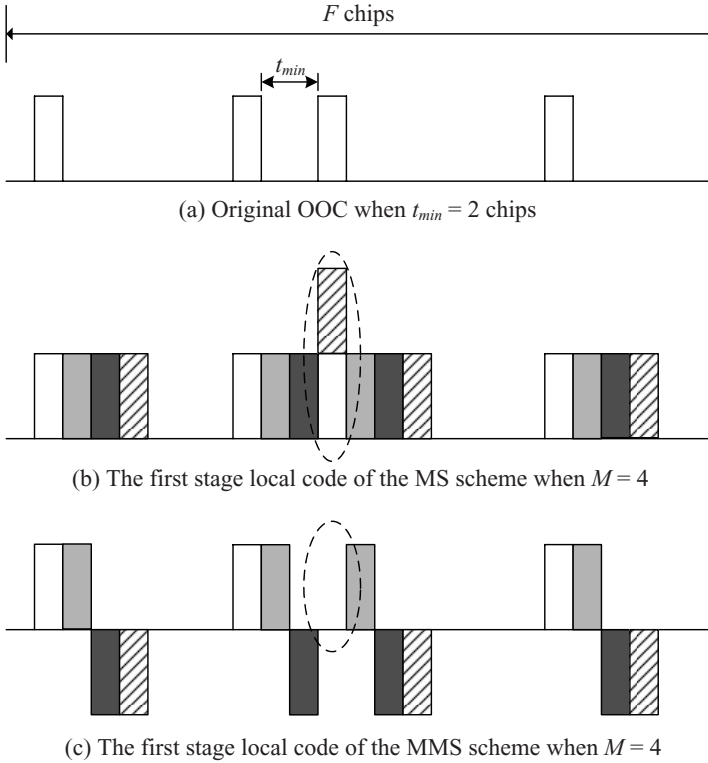


Fig. 3. The local codes of the MS and MMS schemes when $M > t_{min}$

For a correct operation of the proposed MMS and conventional MS schemes, M has to be equal to or smaller than t_{min} , where t_{min} is the minimum chip interval of the OOC. Otherwise, chips of local code are overlapped as shown in Fig. 3, where $t_{min} = 2$ and $M = 4$. In this case, the MMS scheme cannot work appropriately.

Fig. 3 shows the original OOC and the first stage local codes of the MS scheme and the MMS schemes. In Fig. 3, the dashed ellipse means the overlapped chips of the local code. Though the code is not synchronized, the absolute value of the correlation value is increased (decreased) when the amplitude of the overlapped chips is increased (decreased). Thus, the false alarm can be occurred, increasing the MAT.

4 Performance Analysis

In this section, we derive the MAT expression of the MMS scheme by modeling the acquisition process as a discrete-time Markov process [1] with the circular flow graph diagram shown in Figs. 4 and 5. In Fig. 4, ‘ACQ’ and ‘FA’ represent the acquisition and false alarm states, respectively, and $P_D(p_d)$ and $P_{FA}(p_{fa})$ denote the detection and false alarm probabilities in the first and second stages, respectively. Fig. 5 represents the transition process diagrams of $H_{det}(z)$ in Fig. 4.

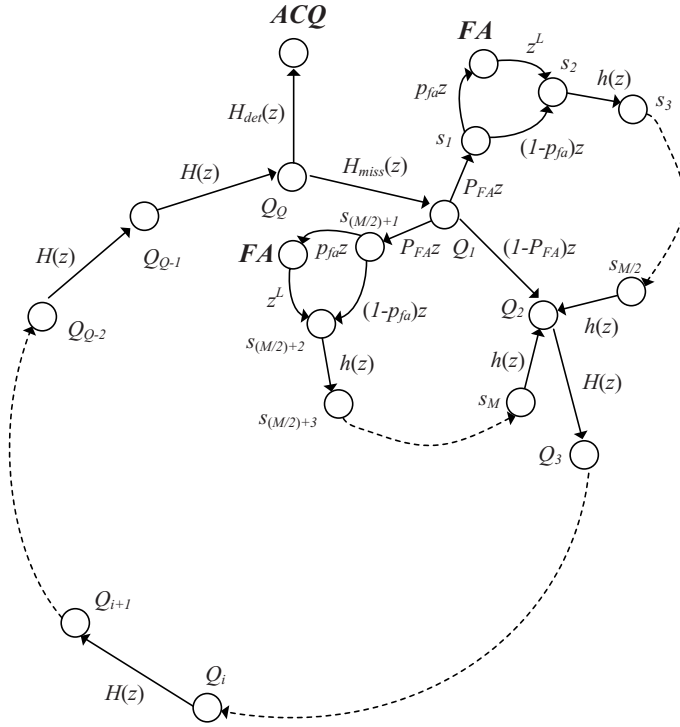


Fig. 4. Markov state model for the MMS scheme

In Fig. 4 states Q_1, Q_2, \dots, Q_{Q-1} correspond to the incorrect groups and a state Q_Q corresponds to the correct group, where the acquisition can be achieved. States $s_1, s_2, \dots, s_{M/2}$ and states $s_{(M/2)+1}, s_{(M/2)+2}, \dots, s_M$ correspond to cells in the second stage, when the correlation value of the correct group is positive and negative, respectively. The gains $H_{det}(z)$ and $H_{miss}(z)$ represent the transition gains from Q_Q to ACQ (detection) and from Q_Q to Q_1 (miss), respectively. The gains $H(z)$ and $h(z)$ represent the transition gains between two consecutive successful incorrect groups and between two consecutive successful incorrect cells, respectively. L is the penalty time factor due to the false alarm. The gains $h(z), H(z), H_{det}^{(v)}(z)$, and $H_{miss}(z)$ can be expressed as follows:

$$h(z) = (1 - p_{fa})z + p_{fa}z^{L+1}, \tag{7}$$

$$H(z) = (1 - P_{FA})z + P_{FA}zh^{M/2}(z), \tag{8}$$

$$H_{det}^{(v)}(z) = p_d P_D z^2 h^{v-1}(z), \tag{9}$$

and

$$H_{miss}(z) = (1 - P_D)z + P_D(1 - p_d)z^2 h^{(M/2)-1}(z), \tag{10}$$

where $v \in \{1, 2, \dots, M/2\}$ is distributed uniformly over $[1, M/2]$ and represents the correct state.

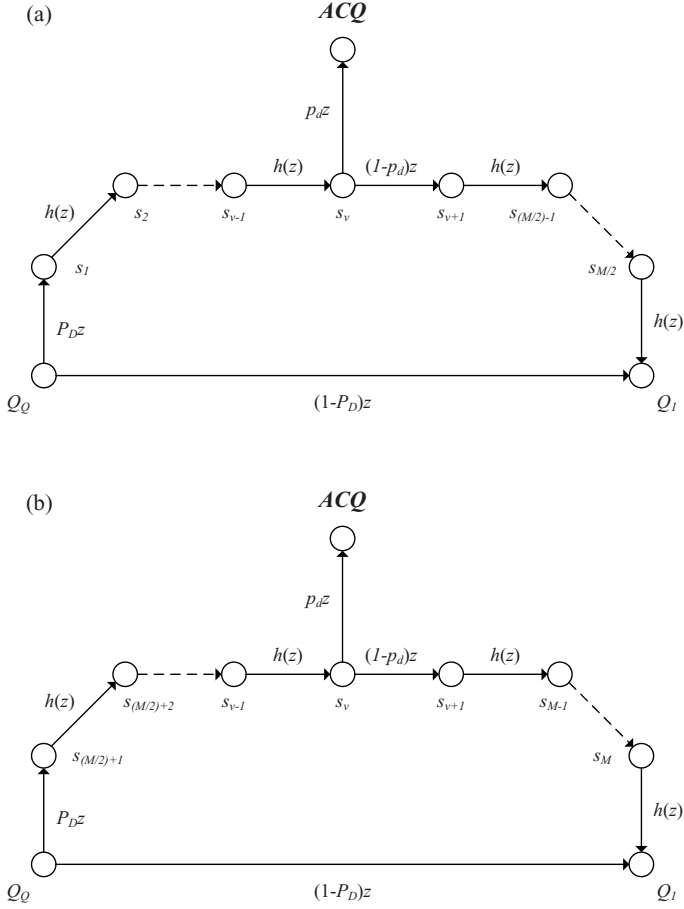


Fig. 5. (a) Positive group from Q_Q to Q_1 (b) Negative group from Q_Q to Q_1

Let us assume that the search is started at Q_i , then the transfer function between Q_i and ACQ nodes can be written as

$$U_i^{(v)}(z) = \frac{H^{Q-i}(z)H_{det}^{(v)}(z)}{1 - H_{miss}(z)H^{Q-1}(z)}. \quad (11)$$

In (11), $i \in \{1, 2, \dots, Q\}$ is assumed to be distributed uniformly over $[1, Q]$. Thus, after averaging $U_i^{(v)}(z)$ over the probability density function (pdf) of i and v , we can re-write (11) as

$$U(z) = \frac{H_{det}(z)}{1 - H_{miss}(z)H^{Q-1}(z)} \frac{1}{Q} \sum_{i=1}^Q H^{Q-i}(z), \quad (12)$$

where $H_{det}(z) = \frac{2}{M} \sum_{v=1}^{M/2} H_{det}^{(v)}(z)$. Using the moment generating function, we can obtain the following relationship between $U(Z)$ and T_{MMS} .

$$U(z) = E(z^{T_{MMS}}), \quad (13)$$

where $E(\cdot)$ denotes the statistical expectation operation. In (13), $E(T_{MMS})$ can be computed by its moment generating function as

$$E(T_{MMS}) = \left. \frac{dU(z)}{dz} \right|_{z=1} = U'(1), \quad (14)$$

and thus, can be obtained as

$$E(T_{MMS}) = \frac{H'_{det}(1) + H'_{miss}(1)}{p_d P_D} + (Q - 1)H'(1) \frac{2 - p_d P_D}{2p_d P_D}, \quad (15)$$

where $H'(1)$, $H'_{det}(1)$, and $H'_{miss}(1)$ can be expressed as

$$H'(1) = 1 + \frac{M}{2} P_{FA}(1 + Lp_{fa}), \quad (16)$$

$$H'_{det}(1) = 2p_d P_D + \frac{(M/2) - 1}{2} p_d P_D (1 + Lp_{fa}), \quad (17)$$

and

$$H'_{miss}(1) = (1 - P_D) + P_D(1 - p_d)[(M/2) + 1 + \{(M/2) - 1\}Lp_{fa}], \quad (18)$$

respectively.

When $p_d = P_D = 1$ and $p_{fa} = P_{FA} = 0$ (i.e., single-user case), (15) is re-written as

$$E(T_{MMS}) = \frac{Q + 1}{2} + \frac{M + 2}{4}. \quad (19)$$

Differentiating (19), we finally obtain the optimum M and minimum T_{MMS} as $\sqrt{2F}$ and $\sqrt{\frac{F}{2}} + 1$, respectively, when $F \gg 1$.

5 Simulation Results

In this section, we compare the MAT of the conventional MS scheme with that of the proposed MMS scheme in single-user and multi-user environments. Simulation parameters are as follows: $F = 200$; $K = 5$; $\lambda_a = \lambda_c = 1$; $N = 1, 2, 3$, and 4 ; $TH_{MS,first} = TH_{MS,second} = TH_{MMS,first} = TH_{MMS,second} = K$; and $L = 10$. We assume that each user transmits the data 1 or 0 with equal probability and chip synchronization is perfect.

Fig. 6 shows the MAT performance of the MS and MMS schemes as a function of M in a single-user environment. In Fig. 6, the dotted and solid lines represent the simulation results of the MS and MMS schemes, respectively, and, the markers ∇ and

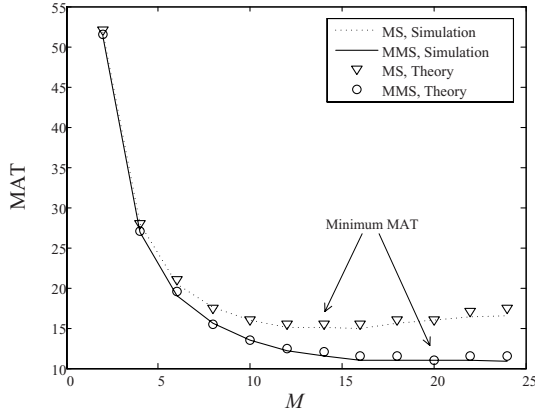


Fig. 6. The MATs of the MS and MMS schemes as a function of M in a single-user environment

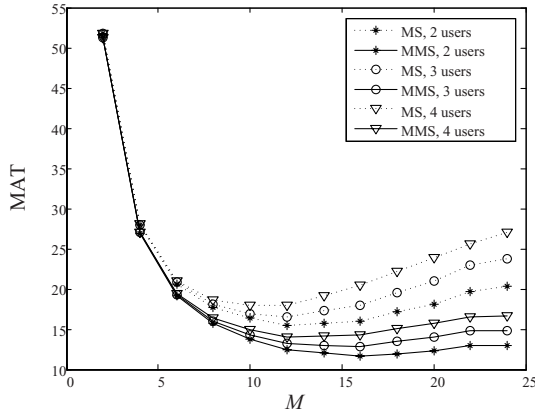


Fig. 7. The MATs of the MS and MMS schemes as a function of M in a multi-user environment

O represent the theoretical results of the MS and MMS schemes, respectively. From the figure, we can observe that the MAT of the MMS scheme is shorter than that of the MS scheme, and confirm that the MS and MMS schemes provide the minimum MAT when $M = 14$ ($\sqrt{F} \simeq 14$) and $M = 20$ ($\sqrt{2F} = 20$), respectively. The difference between the MATs of the MS and MMS schemes increases as M increases. In Fig. 7 we can see the MAT performance of the MS and MMS schemes as a function of M in a multi-user environment. From Fig. 7 we can confirm that the MAT of the MMS scheme is shorter than that of the MS scheme, and the difference between the MATs of the MS and MMS schemes increases as M increases as in the single-user environment.

6 Conclusion

In this paper, we have proposed a novel code acquisition scheme called MMS for optical CDMA systems. Exploiting the modified local code, the proposed MMS scheme can provide a shorter MAT compared with that of the MS scheme. The performance of the MMS scheme has been analyzed using the circular flow graph diagram. The simulation results have confirmed that the MMS scheme offers a shorter MAT compared with that of the MS scheme in both single-user and multi-user environments.

Acknowledgments. This research was supported by the MKE (The Ministry of Knowledge Economy), Korea, under the ITRC (Information Technology Research Center) support program supervised by the NIPA (National IT Industry Promotion Agency) (NIPA-2009-(C1090-0902-0005)).

References

1. Polydoros, A., Weber, C.L.: A unified approach to serial-search spread spectrum code acquisition-Part I: General theory. *IEEE Trans. Commun.* 32(5), 542–549 (1984)
2. Chong, D., Lee, B., Kim, S., Joung, Y.-B., Song, I., Yoon, S.: Phase-shift-network-based differential sequential estimation for code acquisition in CDMA systems. *Journal of Korea Inform. and Commun. Society* 32(3), 281–289 (2007)
3. Salehi, J.A.: Code division multiple-access techniques in optical fiber networks - Part I: Fundamental principles. *IEEE Trans. Commun.* 37(8), 824–833 (1989)
4. Keshavarzian, A., Salehi, J.A.: Optical orthogonal code acquisition in fiber-optic CDMA systems via the simple serial-search method. *IEEE Trans. Commun.* 50(3), 473–483 (2002)
5. Keshavarzian, A., Salehi, J.A.: Multiple-shift code acquisition of optical orthogonal codes in optical CDMA systems. *IEEE Trans. Commun.* 53(3), 687–697 (2005)
6. Salehi, J.A., Brackett, C.A.: Code division multiple-access techniques in optical fiber networks - Part II: Systems performance analysis. *IEEE Trans. Commun.* 37(8), 834–842 (1989)
7. Stok, A., Sargent, E.H.: Lighting the local area: Optical code-division multiple access and quality of service provisioning. *IEEE Network* 14(6), 42–46 (2000)

Optimal and Suboptimal Synchronization Schemes for Ultra-Wideband Systems

Dahae Chong, Chonghan Song, Youngpo Lee, Myungsoo Lee,
Junhwan Kim, and Seokho Yoon*

School of Information and Communication Engineering, Sungkyunkwan University,
300 Chunchun-dong, Jangan-gu, Suwon, Gyeonggi-do, 440-746, Korea
{lvjs1019, starsong83, leeyp204, maxls813,
ourlife3, syoon}@skku.edu

Abstract. The conventional ultra-wideband (UWB) synchronization schemes could cause poor performance in receiver operations such as demodulation after the synchronization, since unreliable timing information corresponding to low-power multipath components might be transferred to the next stage. To solve this problem, in this paper, we propose novel synchronization schemes for UWB systems. We first derive an optimal scheme based on the maximum likelihood (ML) criterion and then develop a simpler suboptimal scheme. Simulation results show that both proposed schemes provide a better synchronization performance over the conventional scheme.

Keywords: Maximum likelihood (ML), Optimal, Suboptimal, Synchronization, Ultra-wideband (UWB).

1 Introduction

Ultra-wideband (UWB) is a technology that has distinct features characterized by its wide bandwidth, and has been adopted as a standard for wireless personal area networks (WPANs) such as IEEE 802.15.3a and IEEE 802.15.4a, due to its high data rate and low-power transmission [1]-[4].

In UWB systems, the received signal must be first synchronized with a locally generated template prior to data demodulation, and thus, a rapid synchronization is one of the most important technical issues in UWB systems [5], [6]. For the fast synchronization in UWB systems, several schemes have been investigated, where timing information corresponding to any of the received multipath components is regarded to be correct [7]-[10].

In [7], Homier and Scholts presented synchronization schemes with various search methods, and analytically derived their mean synchronization time in a noise-free environment, and in [8], Vijayakumaran and Wong introduced a fast synchronization scheme with the permutation search method. Ramachandran and Roy [9] analyzed the schemes of [7] with the tapped-delay-line channel and IEEE 802.15.3a channel models, and Arias-de-Reyna and Acha-Catalina [10] proposed a new serial search synchronization scheme using specific search spacing based on Fibonacci sequence. However, these

* Corresponding author.

schemes could cause poor performance in receiver operations such as demodulation after the synchronization, if the locally generated template would be synchronized with a low-power multipath component. Generally, the first one among multipath components has the largest average power. Thus, the synchronization with the first multipath component is expected to offer a good performance in receiver operations after the synchronization.

In this paper, we propose new synchronization schemes for UWB systems, allowing us to obtain the timing information of the first multipath component efficiently. First, an optimal scheme is proposed based on a maximum-likelihood (ML) criterion. However, the optimal scheme requires the channel information, and thus, is difficult to implement practically. Thus, a simpler suboptimal scheme is developed. The simulation results show that the proposed schemes provide a substantial performance improvement over the conventional scheme.

This paper is organized as follows. Section 2 describes the system model. In Section 3 the optimal and suboptimal schemes are proposed. Section 4 presents the simulation results of the conventional and proposed schemes over several multipath channel models. Section 5 concludes this paper with a brief summary.

2 System Model

In this paper, we consider a direct sequence (DS)-UWB system [9]. During the acquisition process, we assume that an unmodulated signal (acquisition preamble) is transmitted. Then, the transmitted DS-UWB signal $s(t)$ can be expressed as

$$s(t) = \sqrt{E_c} \sum_{i=0}^{N-1} c_i p(t - iT_c), \quad (1)$$

where E_c is the signal energy; $c_i \in \{1, -1\}$ is the i -th chip of a pseudo noise (PN) sequence with a period of N chips; T_c is the chip duration; and $p(t)$ is the second derivative Gaussian pulse [7], [9] with a duration of T_c . At the receiver, the received signal $r(t)$ can be expressed as

$$r(t) = \sqrt{E_c} \sum_{i=0}^{N-1} \sum_{j=0}^{L_p-1} \alpha_j c_i p(t - iT_c - jT_c - \tau T_c) + w(t), \quad (2)$$

where L_p denotes the number of multipaths; α_j is the channel coefficient of the j -th multipath with the average power $\frac{1-e^{-\mu}}{1-e^{-\mu L_p}} e^{-(j-1)\mu}$, where μ is a decay factor; τ is the time delay; and $w(t)$ is an additive white Gaussian noise (AWGN) process with mean zero and two-sided power spectral density $\frac{N_0}{2}$.

In this paper, we consider a parallel receiver for UWB synchronization shown in Fig. 1. The receiver first yields the m -th correlator output y_m by correlating $r(t)$ with the locally generated template signal $g_m(t) = \sum_{n=0}^{N-1} c_n p(t - (n+m)T_c)$, where $m = 0, 1, \dots, N-1$, over a correlation time NT_c , and then, takes the absolute value of y_m to remove the influence of the signal conversion due to the channel reflections.

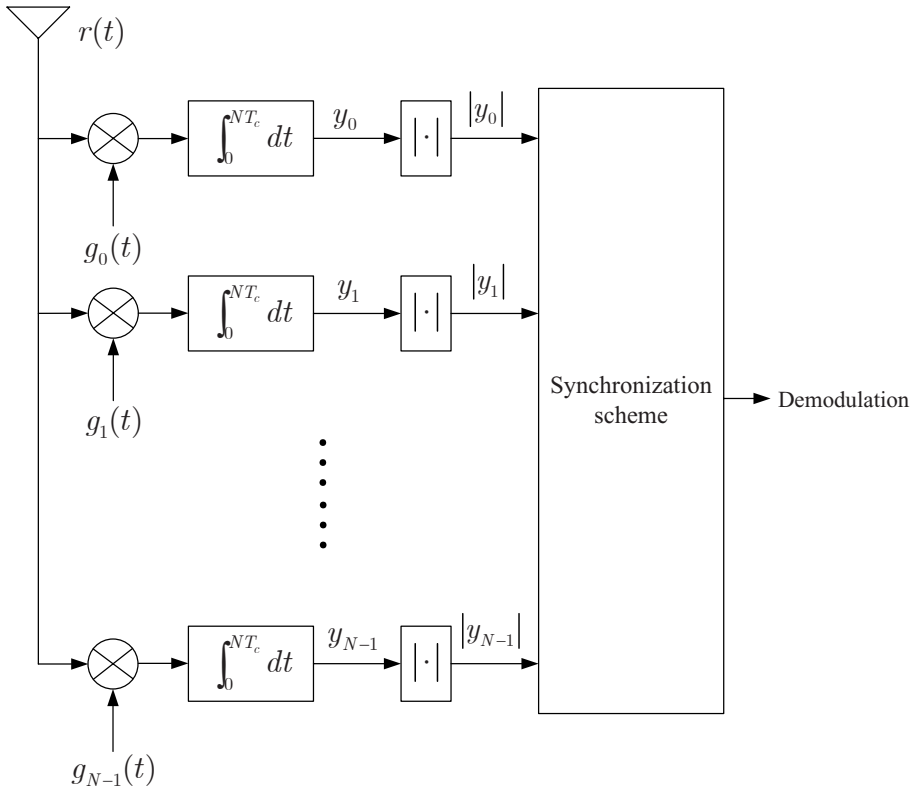


Fig. 1. A parallel receiver for UWB synchronization

Subsequently, the synchronization process produces an estimate of the time delay τ employing the absolute correlator outputs, and finally, the demodulation process is started. In the conventional scheme, an estimate $\hat{\tau}_c$ of the time delay is obtained as

$$\hat{\tau}_c = \arg \max_{0 \leq m \leq N-1} |y_m|. \quad (3)$$

However, the conventional scheme is inefficient to obtain the timing information of the first multipath component, since it does not sufficiently exploit the signal energy spread over multipath components. Thus, we propose novel schemes for selecting the first multipath component efficiently.

3 Proposed Schemes

In UWB channels, the first multipath component generally has the largest power [11], [12]. Based on which, we propose novel schemes based on the ML approach for selecting the first multipath component efficiently.

We first calculate the mean and variance of y_m to obtain the distribution of the correlator output y_m expressed as

$$\begin{aligned} y_m &= \int_0^{NT_c} r(t)g_m(t)dt \\ &= \sqrt{E_c} \int_0^{NT_c} \sum_{i=0}^{N-1} \sum_{n=0}^{N-1} \sum_{j=0}^{L_p-1} \alpha_j c_i c_n \\ &\quad \cdot p(t - (i + j + \tau)T_c)p(t - (n + m)T_c)dt + w_m, \end{aligned} \quad (4)$$

where w_m is the noise component of the m -th correlator output with mean zero and variance $\frac{NN_0}{2}$ and can be written as

$$w_m = \int_0^{NT_c} \sum_{n=0}^{N-1} w(t)c_n p(t - (n + m)T_c)dt. \quad (5)$$

The value of y_m changes according to the phase difference (i.e. $m - \tau$) between the received and template signals as follows:

$$y_m = \begin{cases} N\sqrt{E_c}\alpha_0 + w_0, & m - \tau = 0 \\ N\sqrt{E_c}\alpha_1 + w_1, & m - \tau = 1 \\ \vdots & \vdots \\ N\sqrt{E_c}\alpha_{L_p-1} + w_{L_p-1}, & m - \tau = L_p - 1 \\ w_{L_p}, & m - \tau = L_p \\ \vdots & \vdots \\ w_{N-1}, & m - \tau = N - 1. \end{cases} \quad (6)$$

Then, the probability density function (pdf) f_j of the correlator output y_m corresponding to the j -th multipath component can be written as

$$f_j(y_m) = \frac{1}{\sqrt{\pi NN_0}} \exp\left(-\frac{(y_m - N\sqrt{E_c}\alpha_j)^2}{NN_0}\right), \quad (7)$$

where $N\sqrt{E_c}\alpha_j$ is the mean of the j -th multipath component. The pdf f_w of the correlator output which does not correspond to any of the multipaths can be written as

$$f_w(y_m) = \frac{1}{\sqrt{\pi NN_0}} \exp\left(-\frac{y_m^2}{NN_0}\right). \quad (8)$$

The set of the correlator outputs is denoted by

$$y = [y_0, y_1, \dots, y_{N-1}]^T \text{ with } y_i = y_{(i \bmod N)}. \quad (9)$$

Then, the pdf f of y , given that $m - \tau = 0$, can be expressed as

$$f(y|m) = \prod_{b=0}^{N-1} f_w(y_b) \prod_{j=0}^{L_p-1} \frac{f_j(y_{m+j})}{f_w(y_{m+j})}. \quad (10)$$

Since τ is distributed uniformly over $[0, N - 1]$, the optimal scheme based on the ML approach can be expressed as

$$\hat{\tau}_o = \arg \max_{0 \leq m \leq N-1} f(y|m). \quad (11)$$

In (11), we can remove the term which is independent from m , which results in

$$\hat{\tau}_o = \arg \max_{0 \leq m \leq N-1} \prod_{j=0}^{L_p-1} \frac{f_j(y_{m+j})}{f_w(y_{m+j})}. \quad (12)$$

Using (7) and (8), (12) becomes

$$\hat{\tau}_o = \arg \max_{0 \leq m \leq N-1} \exp \left(\sum_{j=0}^{L_p-1} -\frac{(y_{m+j} - N\sqrt{E_c}\alpha_j)^2}{NN_0} + \frac{y_{m+j}^2}{NN_0} \right). \quad (13)$$

In (13), we can remove the terms which are independent from m and constant. Then, the optimal scheme can be reduced as

$$\hat{\tau}_o = \arg \max_{0 \leq m \leq N-1} \sum_{j=0}^{L_p-1} y_{m+j}\alpha_j, \quad (14)$$

which is optimal to select the first multipath component in UWB channel. However, the optimal scheme requires information on channel coefficients $\{\alpha_j\}_{j=0}^{L_p-1}$, and thus, is difficult to implement. Hence, we develop a simpler suboptimal scheme.

When $m = \tau = 0$, $\sum_{j=0}^{L_p-1} y_{m+j}\alpha_j$ of (14) is computed as

$$\begin{aligned} \sum_{j=0}^{L_p-1} y_j\alpha_j &= \sum_{j=0}^{L_p-1} N\sqrt{E_c}\alpha_j^2 \\ &= \sum_{j=0}^{L_p-1} |N\sqrt{E_c}\alpha_j||\alpha_j| \\ &= \sum_{j=0}^{L_p-1} |y_j||\alpha_j| \end{aligned} \quad (15)$$

in a noise-free environment. Using (15), the optimal scheme can be re-written as

$$\hat{\tau}_o = \arg \max_{0 \leq m \leq N-1} \sum_{j=0}^{L_p-1} |y_{m+j}||\alpha_j|, \quad (16)$$

which requires the absolute values of channel coefficients. By removing $\{|\alpha_j|\}_{j=0}^{L_p-1}$ and using a parameter Q instead of L_p in (16), we obtain the suboptimal scheme as

$$\hat{\tau}_s = \arg \max_{0 \leq m \leq N-1} \sum_{j=0}^{Q-1} |y_{m+j}| \quad (17)$$

for UWB synchronization, where Q is the combining size of the correlator outputs. Combining the correlator outputs corresponding to the paths with low-power may lead to an increase in the noise variance rather than in the signal strength. Thus, Q is set to the average number of paths with power attenuated by at most 10 dB from the power of the first path, and can be predetermined if a channel model is given.

4 Simulation Results

In this section, we compare the performance of the optimal, suboptimal, and conventional schemes in terms of the probability of false synchronization defined as the probability that the estimate of the scheme does not correspond to a timing information of the first multipath component. We define the signal to noise ratio (SNR) as E_c/N_0 and use a PN sequence with a period of 255 chips and a chip duration of 0.5 ns. We simulated the synchronization process both in the tapped-delay line channel model with a tap spacing of 0.5 ns and in the IEEE 802.15.3a channel model [11], [12], where four different environments (CM1, CM2, CM3, and CM4) are defined. As the channel environment changes from CM1 to CM4, the multipaths spread more widely.

Figs. 2[5] show the probabilities of false synchronization of the schemes in the tapped-delay line channel model. As shown in figures, the proposed schemes have a better performance compared with that of the conventional scheme. Since the optimal scheme uses the channel information, it has the best performance. On the other hand,

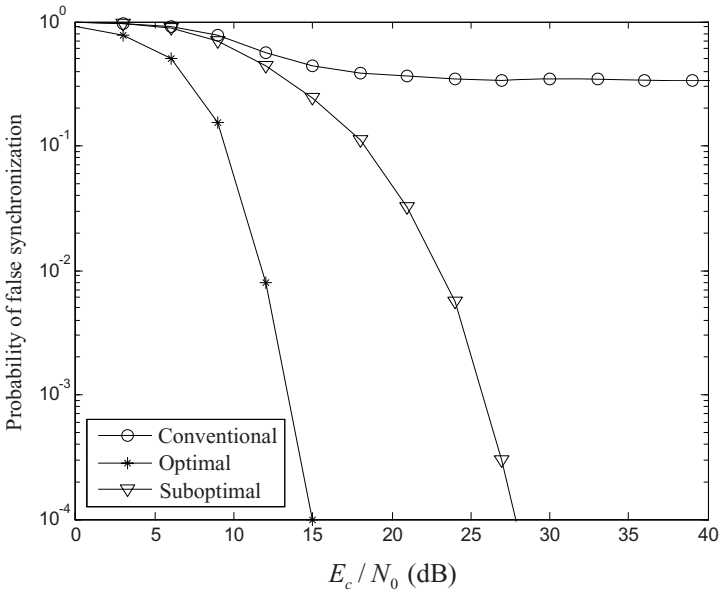


Fig. 2. Probabilities of false synchronization of the conventional, optimal, and suboptimal schemes in the tapped-delay line channel model with CM1 coefficients

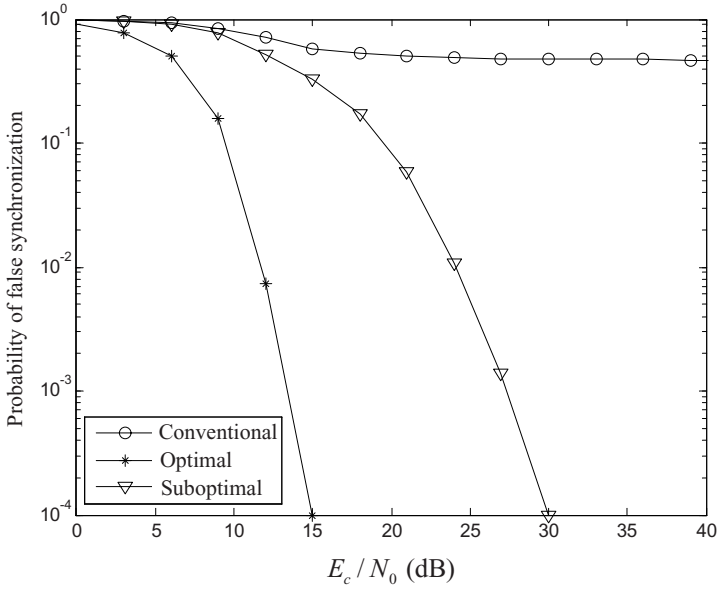


Fig. 3. Probabilities of false synchronization of the conventional, optimal, and suboptimal schemes in the tapped-delay line channel model with CM2 coefficients

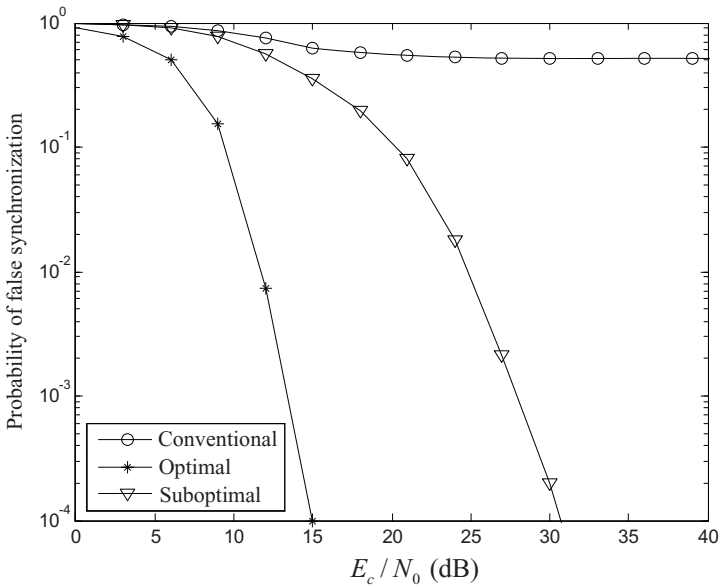


Fig. 4. Probabilities of false synchronization of the conventional, optimal, and suboptimal schemes in the tapped-delay line channel model with CM3 coefficients

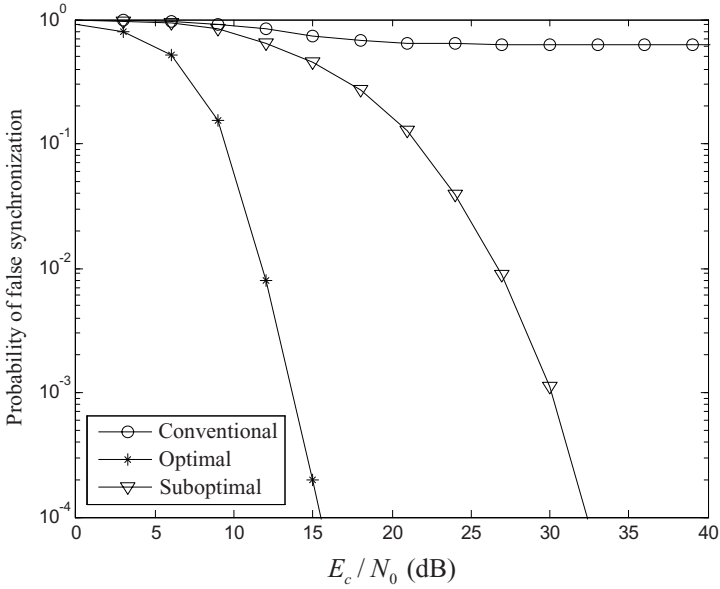


Fig. 5. Probabilities of false synchronization of the conventional, optimal, and suboptimal schemes in the tapped-delay line channel model with CM4 coefficients

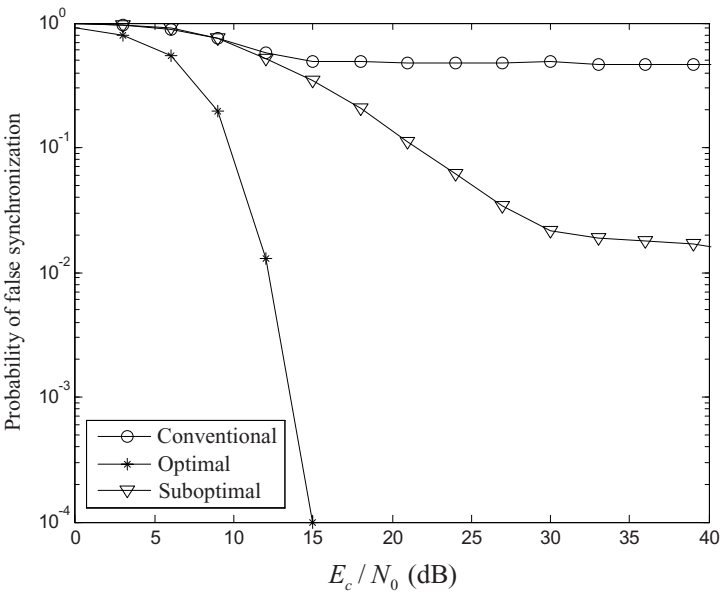


Fig. 6. Probabilities of false synchronization of the conventional, optimal, and suboptimal schemes in the IEEE 802.15.3a CM1 channel model

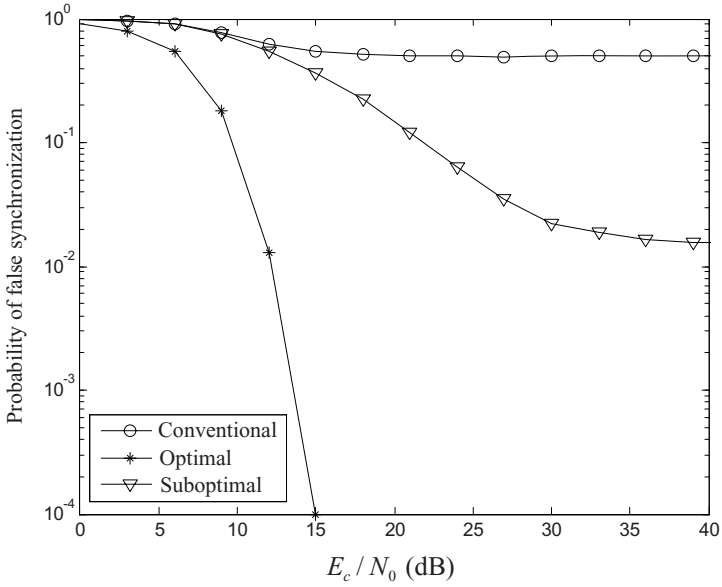


Fig. 7. Probabilities of false synchronization of the conventional, optimal, and suboptimal schemes in the IEEE 802.15.3a CM2 channel model

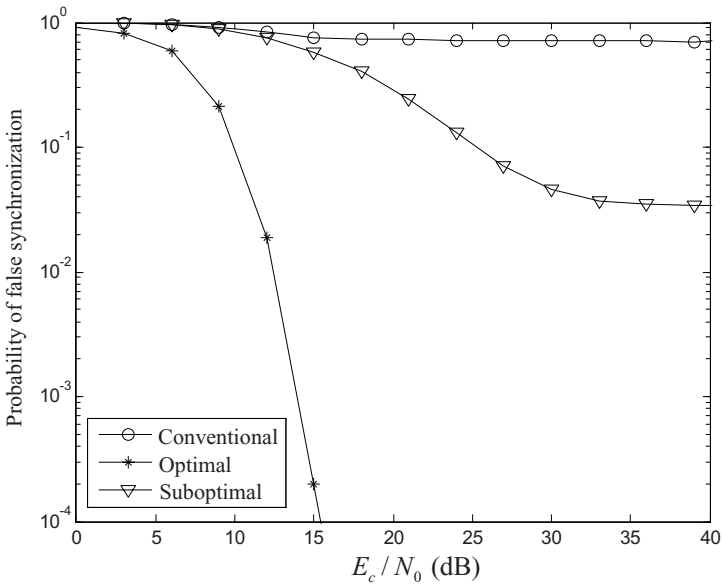


Fig. 8. Probabilities of false synchronization of the conventional, optimal, and suboptimal schemes in the IEEE 802.15.3a CM3 channel model

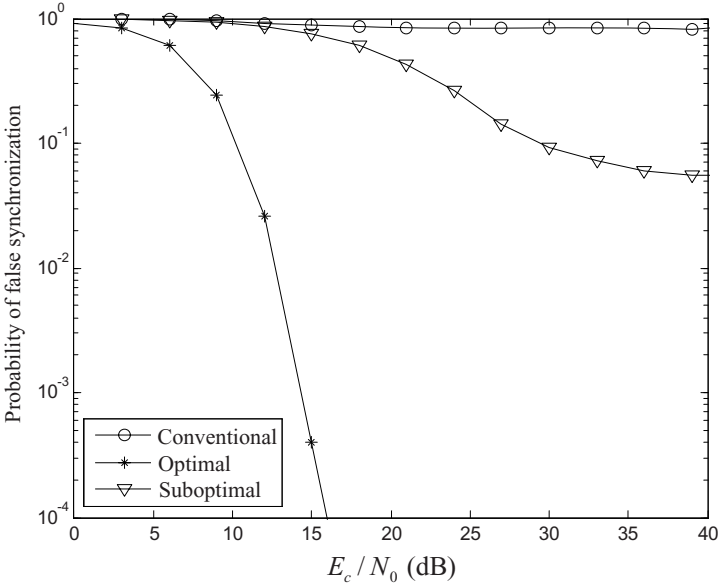


Fig. 9. Probabilities of false synchronization of the conventional, optimal, and suboptimal schemes in the IEEE 802.15.3a CM4 channel model

since the suboptimal scheme does not use the channel information, it exhibits some performance degradation compared with the optimal scheme.

Figs. 6-9 show the probabilities of false synchronization of the schemes in the IEEE 802.15.3a channel model. As shown in figures, the proposed schemes have a better performance compared with that of the conventional scheme. Another important observation is that the performance of the suboptimal scheme is saturated. This is because each correlator output in the IEEE 802.15.3a channel model includes the interference occurred by the other correlator outputs differently from that in the tapped-delay line channel model. However, the optimal scheme utilizes all channel information, and thus, provides the best synchronization performance regardless of the influence of the interference.

5 Conclusion and Future Works

In this paper, we have proposed novel synchronization schemes for UWB systems. We first derive an optimal scheme based on the ML criterion and then develop a simpler suboptimal scheme. Simulation results have shown that both proposed schemes have a better synchronization performance than the conventional scheme. The optimal scheme has the best performance; however, it requires the channel information. On the other hand, since the suboptimal scheme does not require the channel information, it has some performance degradation compared with the optimal scheme.

In the practical channel environment (i.e., IEEE 802.15.3a channel model), the sub-optimal scheme cannot overcome the influence of the interference, and thus, its performance is saturated. In the future, we will deal with the synchronization scheme robust to the influence of the interference.

Acknowledgments. This research was supported by the MKE (The Ministry of Knowledge Economy), Korea, under the ITRC (Information Technology Research Center) support program supervised by the NIPA (National IT Industry Promotion Agency) (NIPA-2009-(C1090-0902-0005)).

References

1. Win, M.Z., Scholtz, R.A.: Impulse radio: how it works. *IEEE Commun. Lett.* 2(2), 36–38 (1998)
2. Porcino, D., Hirt, W.: Ultra-wideband radio technology: potential and challenges ahead. *IEEE Commun. Mag.* 41(7), 64–74 (2003)
3. Yang, L., Giannakis, G.B.: Ultra-wideband communications: an idea whose time has come. *IEEE Sig. Process. Mag.* 21(6), 26–54 (2004)
4. Gezici, S., Kobayashi, H., Poor, H.V., Molisch, A.F.: Performance evaluation of impulse radio UWB systems with pulse-based polarity randomization. *IEEE Trans. Sig. Process.* 53(7), 2537–2549 (2005)
5. Aedudodla, S.R., Vijayakumaran, S., Wong, T.F.: Timing acquisition in ultra-wideband communication systems. *IEEE Trans. Veh. Technol.* 54(5), 1570–1583 (2005)
6. Wu, L., Wu, X., Tian, Z.: Asymptotically optimal receivers with noisy template: design and comparison with RAKE. *IEEE J. Sel. Areas Commun.* 24(4), 808–814 (2006)
7. Homier, E.A., Scholtz, R.A.: Rapid acquisition of ultra-wideband signals in the dense multipath channel. In: *Proc. IEEE UWBST*, Baltimore, MD, pp. 105–109 (2002)
8. Vijayakumaran, S., Wong, T.F.: A search strategy for ultra-wideband signal acquisition. *IEEE Trans. Commun.* 53(12), 2015–2019 (2005)
9. Ramachandran, I., Roy, S.: On acquisition of wideband direct-sequence spread spectrum signals. *IEEE Trans. Commun.* 5(6), 1537–1546 (2006)
10. Arias-de-Reyna, E., Acha-Catalina, J.J.: Blind and efficient serial search strategy for ultrawideband signal initial acquisition. *IEEE Trans. Veh. Technol.* 58(6), 3053–3057 (2009)
11. Molisch, A.F., Foerster, J.R.: Channel models for ultrawideband personal area networks. *IEEE Wireless Commun.* 10(6), 14–21 (2003)
12. Molisch, A.F.: Ultrawideband propagation channels - theory, measurement, and modeling. *IEEE Trans. Veh. Technol.* 54(5), 1528–1545 (2005)

Partial Information Relaying with Multi-Layered Superposition Coding*

Jingyu Kim and Dong In Kim

School of Information and Communication Engineering, Sungkyunkwan University
Suwon 440-746, Korea
{greengyu, dongin}@skku.edu

Abstract. In cooperative communications, relay usually forwards a full information of source data as received. Unlike this approach, this paper proposes a new relaying scheme based on multi-layered superposition coding with multiple antennas which forwards only partial information of superposed layer(s) on top of basic layer(s) of the source data. Fast forwarding of the partial information results in an increase in overall capacity. Simulation results are presented to show the superiority of the proposed scheme over full information relaying and direct transmission.

Keywords: Cooperative relaying, multi-layered superposition coding, multiple antennas, partial information relaying, spectral efficiency.

1 Introduction

Recently there have been many proposals on cooperative relaying as a useful means to achieve better diversity and reliable link gains. Especially, when source-to-destination (S-D) link condition is not sufficient enough to reliably transmit data, it is possible to detour source-to-relay (S-R) and relay-to-destination (R-D) link. Typically it is more likely to have (S-R) and (R-D) links at line-of-sight (LOS), and it is assumed that they provide more reliable link gains than (S-D) link. Under this assumption, relay can play an important role in achieving a better overall data rate [2]. In half-duplexing mode, conventional two-hop relaying transmission forwards a full information of source data as received, which requires appreciable time duration for the second-hop transmission. Such longer second-hop time duration may cause some capacity loss and degrade spectral efficiency. Although direct transmission (DT) does not need any second-hop transmission, we can not offer a required overall data rate when (S-D) link is so weak.

To address the above issue, we propose partial information relaying that utilizes multi-layered superposition coding (MLSC) when multiple antennas are

* This research was supported by the MKE (The Ministry of Knowledge Economy), Korea, under the ITRC (Information Technology Research Center) support program supervised by the NIPA (National IT Industry Promotion Agency) (NIPA-2009-(C1090-0902-0005)).

available at source, relay and destination nodes. It is reported that partial decode-&-forward (DF) relaying with superposition coding can increase the capacity [1], [3], [6], but these contributions considered the partial DF under single antenna scenario. If multiple antennas are used at source, relay and destination nodes, it increases the degree-of-freedom to form the partial information by the MLSC more efficiently so as to increase the overall data rate. Here, the MLSC translates source data into two types of linearly combined data, such as basic layer(s) and superposed layer(s), where relay forwards only the superposed layer(s). The superposed layer is decoded at destination node after relay decodes and forwards it. Since the basic layer is not relayed, it is decoded at destination node after first decoding the superposed layer and successfully cancelling it out. If all data streams (layers) are superposed ones, it is same as full information relaying (FIR) via two-hop transmission. In the other case, if all of those streams are basic ones, it is same as DT. Hence the proposed MLSC can be viewed as the hybrid of FIR and DT.

If the amount of partial information to be relayed is changed in a controlled manner, the second-hop time can be adjusted adaptively [3]. As the basic layer is dominant than the superposed layer, the second-hop time can be made shorter than usual and the spectral efficiency can be improved. However, if (S-D) link condition is not sufficient enough to increase the amount of basic layer, the amount of superposed layer should be increased accordingly. Therefore, how to split the total power between the two layers depends on individual link gains, and a low-rate channel state information (CSI), in terms of the signal-to-interference-plus-noise ratio (SINR) is required at source node to perform the proposed MLSC for partial information relaying.

The rest of the paper is organized as follows. The MIMO cooperative system model with MLSC is described in Section 2. Section 3 explains how to set the key design parameters of MLSC in order to maximize the overall capacity. Simulation results and discussions are given in Section 4, and then concluding remark in Section 5.

2 System Model

As shown in Fig.1, there are single source, relay and destination nodes. Each node has N_S, N_R, N_D antennas, respectively. $\mathbf{H}_0, \mathbf{H}_1$ and \mathbf{H}_2 are the channel matrices of size $N_D \times N_S, N_R \times N_S, N_D \times N_R$ associated with (S-D), (S-R) and (R-D) links, respectively. Because of half-duplexing mode assumed here, relay can not transmit and receive the data streams at the same time.

In the first hop, source node broadcasts the signal \mathbf{x} to relay and destination nodes. \mathbf{x} is composed of \mathbf{x}_s and \mathbf{x}_b where \mathbf{x}_s are the superposed streams being broadcast from $L \leq N_S$ source antennas, and \mathbf{x}_b are the basic streams from the remaining $(N_S - L)$ source antennas. Destination node stores the received signal in the buffer. Relay node decodes only the superposed streams \mathbf{x}_s .

In the second hop, relay node forwards \mathbf{x}_r which is the same as \mathbf{x}_s if decoding is correct. Then, destination node first decodes the received signal from relay

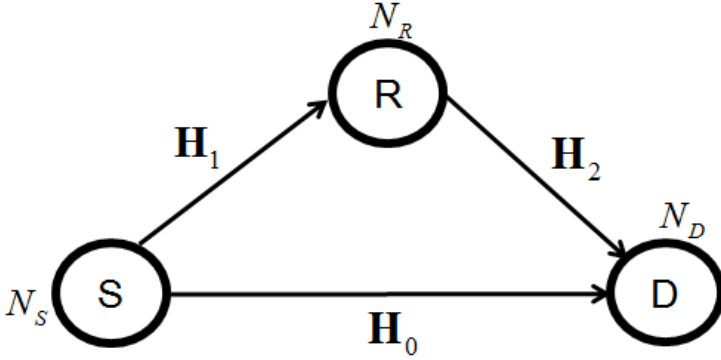


Fig. 1. The wireless channel model with relay

node for obtaining \mathbf{x}_s , and cancels out from the received signal in the first hop. After the successive interference cancellation (SIC), destination node is able to decode the received signal from source node for obtaining \mathbf{x}_b .

Mathematically, the above transmission and reception can be modeled as follows: \mathbf{P}_0 and \mathbf{P}_1 represent the diagonal power matrices at source and destination nodes, respectively. $\mathbf{n}_0, \mathbf{n}_1, \mathbf{n}_2$ are additive white Gaussian noise (AWGN) vectors with zero mean and $\sigma_0^2, \sigma_1^2, \sigma_2^2$ variances.

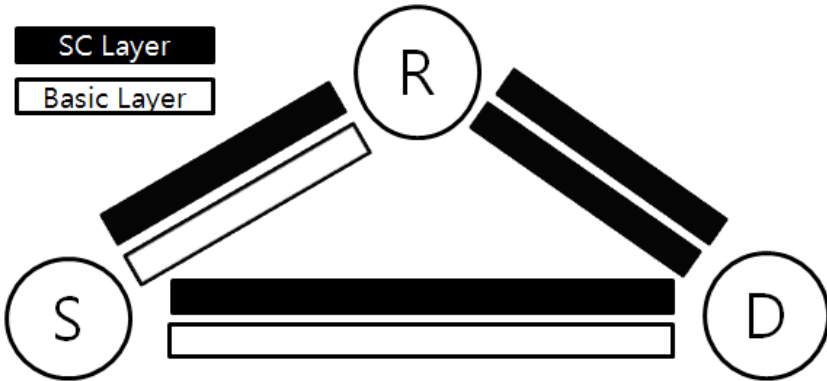


Fig. 2. The illustration of partial information relaying with MLSC

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_b \end{bmatrix}$$

$$\text{1st hop: } \begin{cases} \mathbf{y}_r = \mathbf{H}_1 \mathbf{P}_0 \mathbf{x} + \mathbf{n}_1 \\ \mathbf{y}_{d,1} = \mathbf{H}_0 \mathbf{P}_0 \mathbf{x} + \mathbf{n}_0 \end{cases}$$

$$\text{2nd hop: } \mathbf{y}_{d,2} = \mathbf{H}_2 \mathbf{P}_1 \mathbf{x}_r + \mathbf{n}_2$$

Fig. 2 shows how to forward partial information of the superposed data stream ($L = 1$) on top of the basic one when $N_S = N_R = N_D = 2$.

2.1 MMSE-SIC

The minimum mean-square-error (MMSE) receiver with SIC is employed to decorrelate inter-stream interferences between basic and superposed layers after which the basic layer is decoded using the SIC [4], [5]. To evaluate the overall capacity, we define R_b, R_s and R_2 , namely R_b is the achievable data rate of the basic layer on (S-D) link, R_s the achievable rate of the superposed layer on (S-R) link, and R_2 the channel capacity of (R-D) link as follows:

$$R_b = \sum_{m=1}^{N_S-L} R_{b,m} = \sum_{m=L+1}^{N_S} \log_2 \left(1 + p_{0,m} \mathbf{h}_{0,m}^H \right. \\ \left. \times \left[\mathbf{H}_0(m+1) \mathbf{P}_0(m+1) \mathbf{H}_0^H(m+1) + \sigma_0^2 \mathbf{I}_N \right]^{-1} \mathbf{h}_{0,m} \right) \quad (1)$$

$$R_s = \sum_{m=1}^L R_{s,m} = \sum_{m=1}^L \log_2 \left(1 + p_{1,m} \mathbf{h}_{1,m}^H \right. \\ \left. \times \left[\mathbf{H}_1(m+1) \mathbf{P}_1(m+1) \mathbf{H}_1^H(m+1) + \sigma_1^2 \mathbf{I}_N \right]^{-1} \mathbf{h}_{1,m} \right) \quad (2)$$

$$R_2 = \sum_{n=1}^{N_R} R_{2,n} \\ = \log_2 \det \left(\mathbf{I}_{N_D} + \mathbf{H}_2 \mathbf{P}_1 \mathbf{H}_2^H \right) \quad (3)$$

In the above $\mathbf{h}_{i,m}$ denotes the m th column of matrix \mathbf{H}_i ($i = 0, 1$) with N rows and M columns, $\mathbf{P}_i(m) = \text{diag}\{p_{i,m}, p_{i,m+1}, \dots, p_{i,M-1}, p_{i,M}\}$, and $\mathbf{H}_i(m)$ is the $N \times (M - m + 1)$ matrix $[\mathbf{h}_{i,m}, \mathbf{h}_{i,m+1}, \dots, \mathbf{h}_{i,M-1}, \mathbf{h}_{i,M}]$. The superscript $(\cdot)^H$ denotes the hermitian transpose operation, and \mathbf{I}_N is the identity matrix of size $N \times N$. Note that R_b is evaluated under the assumption of perfect SIC while R_s in the presence of the inter-stream interferences from basic layers.

2.2 Power Allocation

The source data is divided into a number of superposed layers and the rest of basic layers with a total power constraint of P_T . If we define the power division factor to allocate the total power to a set of superposed layers and the rest of basic layers, they are allocated the powers $\alpha \times P_T$ and $(1 - \alpha) \times P_T$, respectively. The MLSC should share a total of source antennas N_S among the superposed and basic layers, and if L source antennas are used for the superposed layers where $0 \leq L \leq N_S$, the remaining $(N_S - L)$ antennas are assigned to the basic layers. Assuming that the SINR estimates are available at source and relay nodes,

an optimal power allocation based on water-filling algorithm can be performed. With the total power constraint applied, the above power allocation can be formulated as

$$P_{0,m} = \begin{cases} \left(\frac{1}{\lambda_0^{sc}} - \frac{\sigma_1^2}{|h_{1,m}|^2} \right)^+ & \text{for the superposition layers } (0 \leq m \leq L) \\ \left(\frac{1}{\lambda_0^b} - \frac{\sigma_1^2}{|h_{1,m}|^2} \right)^+ & \text{for the basic layers } (L + 1 \leq m \leq N_s) \end{cases} \quad (4)$$

$$P_{1,n} = \left(\frac{1}{\lambda_1} - \frac{\sigma_2^2}{|h_{2,n}|^2} \right)^+ \quad (5)$$

$$\text{subject to } \sum_{m=1}^L P_{0,m} \leq \alpha \times P_T, \quad \sum_{m=L+1}^{N_s} P_{0,m} \leq (1 - \alpha) \times P_T, \quad \sum_{n=1}^{N_R} P_{1,n} \leq P_T \quad (6)$$

where $(x)^+ = \max(0, x)$ and the Lagrange multipliers λ_0^{sc} , λ_0^b and λ_1 are chosen to meet the total power constraints.

2.3 Overall Capacity

Once R_b , R_s and R_2 are evaluated, we can derive the overall capacity R_{sc} for the partial information relaying with MLSC. A total number of bits that source transmits to destination is $N \times (R_b + R_s)$ where N denotes the duration of the first hop. Since relay forwards only $N \times R_s$, the duration of the second hop can be $N_2 = \frac{NR_s}{R_2}$ so that the overall capacity can be derived as [3]

$$R_{sc} = \frac{N(R_b + R_s)}{N + N_2} = \frac{(R_b + R_s)R_2}{R_2 + R_s} \quad (7)$$

To compare R_{sc} with conventional schemes, such as full information relaying (FIR) via two-hop transmission and direct transmission (DT), the corresponding overall rates are evaluated as

$$R_0 = \log_2 \det \left(\mathbf{I}_{N_D} + \mathbf{H}_0 \mathbf{P}_0 \mathbf{H}_0^H \right) \quad (8)$$

$$R_1 = \log_2 \det \left(\mathbf{I}_{N_R} + \mathbf{H}_1 \mathbf{P}_0 \mathbf{H}_1^H \right) \quad (9)$$

$$R_{02} = \log_2 \det \left(\mathbf{I}_{N_D} + \mathbf{H}_0 \mathbf{P}_0 \mathbf{H}_0^H + \mathbf{H}_2 \mathbf{P}_1 \mathbf{H}_2^H \right) \quad (10)$$

$$R_{FIR} = \frac{R_1 \times R_{02}}{R_1 + R_{02}} \quad (11)$$

$$R_{DT} = R_0 \quad (12)$$

where R_0 is the capacity of (S-D) link, R_1 the capacity of (S-R) link, and R_{02} the capacity for maximal ratio combining (MRC) of (S-D) and (R-D) links in the second hop.

3 Precoding and Power Division

3.1 Precoding

In (2) we notice that relay node decodes \mathbf{x}_s in the presence of the inter-stream interferences from \mathbf{x}_b . Given that relay node does not need to decode \mathbf{x}_b , a precoding matrix to decorrelate the interferences may be inserted to increase the capacity R_{sc} . For this the singular value decomposition is applied to \mathbf{H}_1 as

$$\mathbf{H}_1 = \mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^H \quad (13)$$

where \mathbf{U}_1 and \mathbf{V}_1 are unitary matrices and $\mathbf{\Sigma}_1$ is a diagonal matrix as defined in [7]. To decorrelate the interferences from basic layer(s), the precoding matrix at source node could be set to \mathbf{V}_1 . The resulting channel matrices $\hat{\mathbf{H}}_i$ ($i = 0, 1$) are formulated as

$$\hat{\mathbf{H}}_i = \mathbf{H}_i \mathbf{V}_1 \quad (14)$$

The overall capacity with precoding is evaluated after \mathbf{H}_i is replaced by $\hat{\mathbf{H}}_i$ in (1) and (2). Note that precoding at relay node is not performed since the power allocation via water-filling boosts up the capacity of (R-D) link.

3.2 Choice of the Power Division Factor

The overall capacity is adjusted by two variables α and L . α is a power division factor between superposition and basic layers and can be selected any real value between 0 and 1. L is the number of transmit antennas for superposed layers and modulates transmit diversity of MLSC system and can be selected any discrete value between 0 and N_s . Since our objective is to maximize the overall capacity, α and L can be found by solving following optimization problem.

$$\max_{\alpha, L} R_{sc}(\alpha, L) = \frac{(\sum_{m=1}^L R_{s,m}(\alpha) + \sum_{m=L+1}^{N_s} R_{b,m}(1 - \alpha)) R_2}{\sum_{m=1}^L R_{s,m}(\alpha) + R_2} \quad (15)$$

To lower the complexity of above equation, reduce the optimization problem from two variables to one variable. Notice that L is the number of transmit antennas for the superposition layers and it can be decided by the allocated power properly. Assume that large alpha (e.g. $\alpha = 0.8$) is chosen. Because it indicates that more power are allocated to the superposition layers than the basic layers, antennas for the superposition layers should be increased. Surely opposition can be established. For easy explanation, when there are 4 antennas at the source, L can be determined as follows.

When there are 4 antennas at the source node

$$\alpha = 0 \rightarrow L = 0$$

$$0 < \alpha < 0.375 \rightarrow L = 1$$

$$0.375 \leq \alpha < 0.625 \rightarrow L = 2$$

$$0.625 \leq \alpha < 1 \rightarrow L = 3$$

$$\alpha = 1 \rightarrow L = 4$$

If N_s is another value (e.g. 2,3,5,...), we can set the number L as similar way.

Since L is obtained by α , optimizing problem (15) is a function of only α . To determine the optimal α , the above problem is simplified as follows.

$$\begin{aligned} & \arg \max_{\alpha} R_{sc}(\alpha) \\ & \text{subject to } 0 \leq \alpha \leq 1 \end{aligned} \quad (16)$$

Since this optimization problem belongs to a convex problem, this is solved by some numerical search algorithm (e.g. bisection method) or exhaustively [8].

We have already assumed that source knows CSI of (S-D) and (S-R). But source usually does not know exact (R-D) channel gain. If source only knows statistical information of (R-D) channel, R_2 should be replaced with the expected value \bar{R}_2 . Although \bar{R}_2 is not the same as real R_2 value, using it as suboptimal does not cause significant error to find optimal α .

4 Simulation Results

For the simulation we consider the multiple antennas cooperation system ($N_s = N_r = N_d = 4$ antennas). SNR is defined by the ratio of total power P_T and noise variance σ^2 . Fig. 3 shows the simulation result when SNRs of (S-R) and (R-D) are all 20dB. The result of MLSC is compared with results of FIR and DT in this graph. It is seen that MLSC shows much better performance than other schemes.

Especially when SNR of (S-D) is around 10dB, performance gap of MLSC and the others is significant. To explain the reason, we illustrate Fig. 6 showing a trend of optimized α value. It is seen that the value of α gets smaller as the SNR of (S-D) approaches high value. From this result, MLSC operates similarly as DT at the low SNR of (S-D) and FIR at high. Because small alpha indicates that MLSC operates like DT, however large alpha does like FIR. Around the middle of SNR of (S-D) axis, the value of alpha is close to 0.5. This result means that when channel gain of (S-D) is some normal value (in this case 6-12dB region), signal power of the superposed and basic layers are similar. Hence the duration of second hop time is adaptively modulated and that causes significant performance difference between MLSC and the other schemes.

For the asymmetric link, we illustrate Fig. 4 and Fig. 5. In Fig. 4, SNRs of (S-R) and (R-D) are 15 and 20dB respectively and the capacity gain of MLSC at low SNR region shows more prominent compared with Fig. 3 and Fig. 5. From this result, MLSC can be more useful when SNRs of (S-D) and (S-R) are not so different.

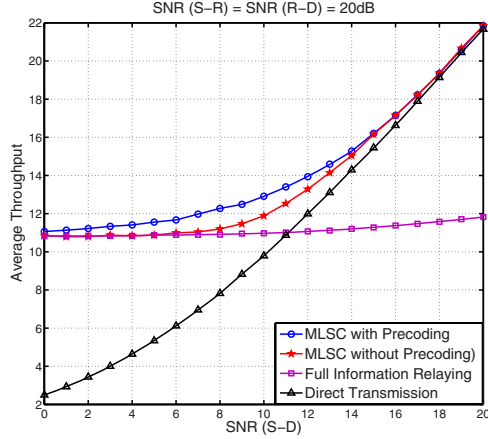


Fig. 3. Average capacity when both SNRs of (S-R) and (R-D) are 20dB

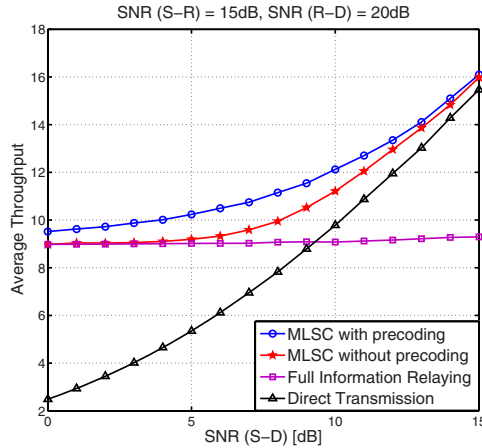


Fig. 4. Average capacity when SNRs (S-R) of and (R-D) are 15 and 20dB respectively

Furthermore, it is apparent that MLSC with precoding outperforms without it. It is seen that decorrelating interferences from basic layer increases overall capacity R_{sc} . This gap is more prominent when SNR of (S-D) is around 8dB because more interferences exist without precoding.

As you see the result, the capacity of FIR is mostly affected by relay link channels (S-R) and (R-D). On the other hand, the capacity of DT is determined by (S-D) channel only. So the key idea of this proposed scheme is to use all the channels adaptively depending on instant channel gains. So as we see the result, the proposed scheme named MLSC is very useful to achieve high transmission rate.

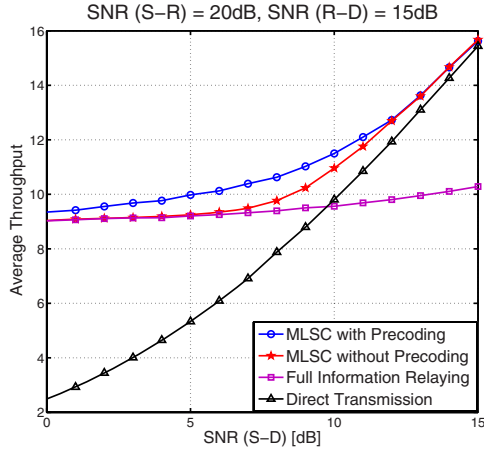


Fig. 5. Average capacity when SNRs (S-R) of and (R-D) are 20 and 15dB respectively

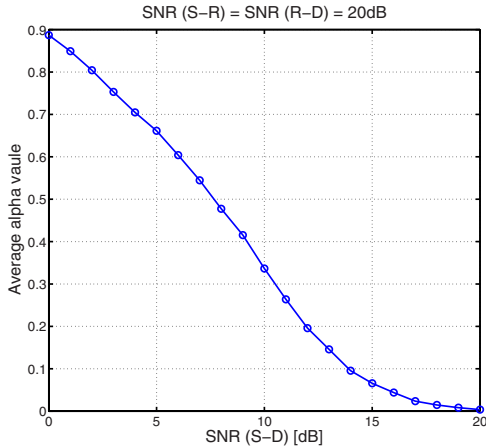


Fig. 6. Average α value versus SNR of (S-D) when both SNRs of (S-R) and (R-D) are 20dB

5 Conclusion

We have proposed the partial information relaying based on MLSC where relay forwards only a part of source data. By controlling the amount of information to be relayed, considering the asymmetric link conditions, the second-hop time at relay can be adjusted adaptively. Although finding an optimized power division factor causes some overhead to the proposed relaying scheme, we can achieve much improved overall capacity. Moreover, it was shown that the precoding method can be effective to decorrelate the inter-stream interference and yield better capacity. Because we have only considered a single relay node, multiple-relay

cooperative system based on the MLSC can be a future work. Here, which and how many links are allocated to the basic and superposed layers would be an interesting issue to be further addressed.

References

1. Yuksel, M., Erkip, E.: Broadcast Strategies for the Fading Relay Channel. In: IEEE Proc. MILCOM 2004, October 2004, vol. 2, pp. 1060–1065 (2004)
2. Laneman, J.N., Tse, D.N.C.: Cooperative diversity in wireless networks: efficient protocols and outage behavior. *IEEE Trans. Inform. Theory* 50, 3062–3080 (2004)
3. Popovski, P., de Carvalho, E.: Improving the rates in wireless relay systems through superposition coding. *IEEE Trans. Wireless Commun.* 7, 4831–4836 (2008)
4. Wolniansky, P.W., Foschini, G.J., Golden, G.D., Valenzuelar, R.A.: V-BLAST: an architecture for realizing very high data rates over the rich-scattering wireless channel. In: Proc. of URSI International Symposium on Signals, Systems and Electronics (ISSSE 1998), Pisa, Italy, September 1998, pp. 295–300 (1998)
5. Chung, S.T., Lozano, A., Huang, H.C., Sutivong, A., Cioffi, J.M.: Approaching the MIMO Capacity with a Low-Rate Feedback Channel in V-BLAST. *EURASIP Journal on Applied Signal Processing* 5, 762–771 (2004)
6. Goparaju, A.K., Wei, S., Liu, Y.: On superposition coding based cooperative diversity schemes. In: Proc. 65th IEEE Vehicular Technology Conference (VTC Spring 2007), Dublin, Ireland, April 2007, pp. 1046–1050 (2007)
7. Andrea Goldsmith: *Wireless Communication*, Cambridge (2005)
8. Boyd, S., Vandenberghe, L.: *Convex Optimization*, Cambridge (2003)

Performance Measurement of the Hybrid Prefetch Scheme on Vehicular Telematics Networks^{*}

Junghoon Lee¹, Gyung-Leen Park^{1, **}, Youngshin Hong¹,
In-Hye Shin¹, and Sang Joon Lee²

¹ Dept. of Computer Science and Statistics, ² Dept. of Computer Engineering
Jeju National University, 690-756, Jeju Do, Republic of Korea
{jhlee, glpark, yshong, ihshin76, sjlee}@jejunu.ac.kr

Abstract. This paper measures and analyzes the performance of the hybrid prefetch scheme for gateways on the vehicular telematics network where every information is indexed by the link (street) ID, via simulation using a discrete event simulator. Combining the classic LRU and FAR techniques, the hybrid scheme groups links according to whether they are referenced during the given time interval or not, orders by the Euclidean distance in each group, and fetches the referenced set first. The extensive experiment results demonstrate that 1) the hybrid scheme can improve the request-level instant reply ratio by up to 20.3 %, overcoming 20 % loss of record-level hit ratio, compared with LRU, 2) this scheme outperforms other spatial-locality-based schemes at all satisfiability ranges and LRU when the satisfiability is above 0.8, also showing stable reply ratio for the practical time tolerance range, 3) it quickly stabilizes after the prefetch memory initialization and is less affected by the hourly data access pattern change, maintaining steady instant reply ratio for the number of records in each request.

1 Introduction

The vehicular telematics network makes it possible to provide an information service to drivers or passengers by keeping vehicles connected to the global network even while they are driving. The most common application is data retrieval. Traffic information, weather update, parking lot availability, local restaurant information, and the like is provided to the user [1]. A vehicle connects to a gateway such as the RSU (Road Side Unit), which relays information from the server to mobile hosts, maintaining a copy of original data created at the server [2]. Each server usually updates its information at regular intervals, making the data copy in the local memory, or cache, of gateways get stale. Gateways must reload the data from each information server. If the information a vehicle requested is in the

^{*} This research was supported by the MKE, Korea, under the ITRC support program supervised by the NIPA. (NIPA-2009-C1090-0902-0040)).

^{**} Corresponding author.

gateway and valid, the gateway can send it immediately. Otherwise, the gateway must refresh its local copy, either proactively or reactively.

Prefetching is a technique in which a data consumer proactively fetches from the server the data item that is predicted to be needed in the near future [3]. The main concern of the prefetch scheme is necessarily how to decide the prefetch order. Not to mention, the data item many vehicles commonly need must be prefetched first. However, the size of such information is too large to refresh the whole set. For example, the amount of traffic information for a city or country is too huge for a gateway to reload all data items for each update period. In the location-dependent information service, drivers around a specific gateway are generally interested in such information as train schedule at nearby terminals and up-to-date traffic conditions of the route to their destinations [4]. If multiple vehicles access the same item, the prefetch policy for the gateway must take into account the spatial and temporal locality between the requests from different vehicles.

Prefetch schemes are basically limited by both prediction accuracy and the penalty for misprediction [5]. The accuracy depends on the efficiency of a prefetch policy and the size of cache. However, considering the cost and size of installable memory these days, the cache size is not a restriction any more. As contrast, as the amount of data records for a city is very large, it is impossible for each gateway to download from the server immediately, while all records are updated at the server almost simultaneously. Hence, the refresh order most impacts the data validity in the gateway. The prefetch order is based on the prediction which item is likely to be requested most and soon. In addition, the penalty of misprediction is definitely the waste of network bandwidth, or sometimes communication cost, and the increase in response time.

There have been a lot of cache/prefetch strategies for single mobile host, and their concern mainly lie in mobility prediction and personalization. As contrast, prefetch policies for the gateway must be based on the common access pattern from the vehicles in the vicinity [6]. So, prefetch in gateways must consider the temporal and spatial locality of data access, possibly exploiting classical prefetch policies such as LRU (Least Recently Used) and FAR (Furthest Away Replacement) [7]. Data prefetching and caching have many features in common, and each policy in cache management has its counterpart in the prefetch policy. First, in cache management, LRU evicts the item that is least recently used. This name can be sustained in the prefetch domain if we change first to last. That is, LRU in the prefetch policy means that the data item least recently used will be updated last. Similarly, FAR in cache management replaces the item which is farthest away from the reference point. In the prefetch, the item that is closest to a reference point will be fetched first. After all, LRU and FAR can be interchangeably used in cache and prefetch.

Combining LRU and FAR, our previous paper has proposed a hybrid prefetch scheme which groups the data records into two sets according to whether a record was referenced during the given previous interval or not [8]. In each set, records are ordered by the Euclidean distance from the gateway. The hybrid scheme

can fetch the closer item first in each group. Even though the record-level hit ratio is less than the LRU scheme, the hybrid scheme can enhance the request-level instant reply ratio. However, user requirement has a variety, not just the basic hit or reply ratio. In this regard, this paper is to measure and analyze the performance of the hybrid prefetch scheme in terms of diverse user satisfiability criteria. For example, a user satisfies if at least 90 % of information is instantly available or a user can tolerate up to 10 second delay.

This paper is organized as follows: After defining the problem in Section 1, Section 2 describes the background of this paper and related work. Section 3 explains the experimental setup such as the target network model and the data access pattern. After the performance measurement results are analyzed in Section 4, Section 5 summarizes and concludes this paper with a brief introduction of future work.

2 Background and Related Work

2.1 Classic and Hybrid Prefetch Schemes

This section introduces some previous work on data prefetching in the mobile network, focusing on the hybrid prefetch scheme. A prefetch scheme is built on top of prediction on the future data access from a single or multiple clients. Prediction depends on a heuristic which has been intensively developed in the cache management domain, focusing on how to take advantage of the locality of data access. One of the most basic cache schemes is LRU and this scheme is known to best take advantage of temporal locality, that is, a recently referenced item will be referenced again soon. In the mobile network, spatial locality has been also considered to account for the characteristic that an item near the recently referenced one will be referenced soon.

Let's assume a service vehicles retrieve information on all links comprising each trip. The information of the link many vehicles commonly take is necessarily referenced frequently. Intuitively, more vehicles will likely take the road segment closer to a gateway, as all trips pass the road near to gateway. However, not all of close roads are taken by drivers. Some roads, even if they are close to the gateway, may lead to a dead end or always have heavy traffic, so the navigation software does not select them. FAR schemes cannot filter such road segments, possibly creating not a little misprediction. However, the prefetch module can calculate the distance, namely, either Euclidean or network distance, in priori, so this heuristic has very small run-time overhead. Oppositely, LRU is known to have highest hit ratio for most cases of cache management, however, it cannot give precedence to the road segment closer to the gateway. In the road network, a road segment, especially far away from the gateway, which has been just referenced, cannot be referenced for a long time.

Based on the observation for LRU and FAR, the hybrid prefetch scheme groups the links into two sets [8]. The first is the group of links referenced during the given previous interval, say, 1 hour, and the other, those not referenced. This categorization takes advantage of LRU. Actually, in the mid-size

city level scenario, for a gateway on which most traffic is concentrated, the links taken by the vehicle starting from a specific gateway location for a whole day are below 30 %. In each set, links are ordered by the Euclidean distance to take into account spatial locality. The first set is able to gather the links that many drivers commonly take and will be prefetched prior to the second set. The link far away from the gateway will be fetched later within this group. In basic LRU, when two links are not referenced yet, it is not possible to have an item closer to the gateway fetched first. The hybrid scheme can first fetch the closer item to account for the higher probability to be referenced within the second group.

For a car to reach a place far away, it must start from a link near to it. In the vehicular telematics network, the movement of a vehicle is restricted to the road network. So, the distance between the two points can be represented either by the Euclidean distance or by the network distance. The Euclidean distance is the length of the straight line connecting the two end points, while the network distance is the sum of each road segment distance along the path. In the road network, even if a location is very close in the Euclidean distance, its network length can be very large when there is no direct path from the start position. FAR-E denotes the policy which prefetches first the item closest to the given point in the Euclidean distance, while FAR-N uses the network distance instead of the Euclidean distance. Spatial locality is well reflected by FAR schemes.

2.2 Other Prefetch/Cache Schemes

Besides, there are some important work on prefetch for mobile host case. To begin with, a power-controlled prefetching scheme has been designed, focusing on how to handle the intermittent connectivity of radio links from data servers to mobile clients and how to manage power and bandwidth resources [9]. The system decides the item to prefetch based on the future request estimation, and this estimation employs a statistical profile which models the user and channel behaviors for constant time intervals based on a Markovian model. For each time slot, the global system state is defined by user, channel interference, and terminal buffer states. This scheme focused on which item to evict from the local buffer of a mobile host.

Liang and et al. considered mobile users in a two-tier network, comprised of WLANs surrounded by a ubiquitous cellular network [4]. The authors claimed that successfully prefetched data, when a user is about to leave the coverage area, potentially save expensive future cellular network access. Users may roam anywhere and are not constrained to any one network. The mobility pattern of each vehicle is known in priori, making it possible to estimate which networks a vehicle may roam to in the near future. This work defined the total cost for accessing a document as the sum of access cost and the penalty for access delay. A mathematical model was developed for the cost function, by which the WLAN gateway decides whether to prefetch an item and the prefetch order.

P. Deshpande et al. have proposed a cooperative prefetch scheme based on the location prediction of vehicles [10]. Their scheme begins with an argument that people drive on familiar routes frequently, so the mobility and connectivity during their drives can be predicted quite accurately by their driving history [11]. By means of this prediction, multiple roadside APs on the route of a vehicle cooperatively download a portion of a large object. When the vehicle approaches and connects to an AP, the prefetched data can be directly provided. All these schemes are developed based on the prediction of the future location of a vehicle.

Additionally, some strategies integrate an index mechanism to cache/prefetch policies. First, a semantic cache scheme builds a Voronoi diagram and stores the data item along with its semantic description in the mobile client [12]. By this semantic description, it is possible to provide partial answers to the given query, even if it does not match exactly. This scheme works quite well for the continuous queries which keeps receiving updated information while a vehicle is moving.

3 Experiment Setup

3.1 Network Architecture and Information Service

Prefetching is used in widely different contexts in computer systems to speed up data retrieval. So, the design must be preceded by the specification of the target network architecture and service. This paper shares the underlying framework with our previous work [8]. Figure 1 illustrates the vehicular telematics network architecture. Static gateways are installed at the place with high vehicle density. A vehicle can connect to a gateway by means of a pre-defined network protocol such as IEEE 802.11, DSRC (Dedicated Short Range Communication), Zigbee, and so on [13]. Each gateway has a stable connection to the global network, be it wired or wireless. Servers reside in the Internet domain and can send data to the gateway via the Internet protocol, while the service provider creates and provisions useful contents.

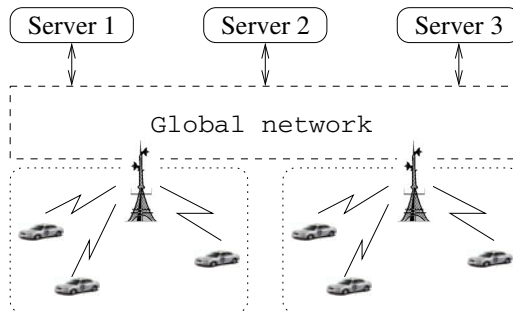


Fig. 1. Vehicular telematics service network

Apparently, one of the most frequently used services for the telematics network is path finding. In the advanced telematics service scenario, besides the simple route selection, a driver may want to know detailed live information along the path, including traffic condition update, seat availability at the roadside restaurant, and an instant discount coupon of shopping stores that a driver may visit during his travel. With the data, he can make a detailed plan or even try to another route. Those data are stored in the central information server, periodically updated in general, and retrieved by the vehicle by way of gateways. In the meantime, all information stored in the information server is associated and indexed by the link ID, so the link ID is the search key to the local cache of gateways as well as server database. Link ID, which is analogous to the street name, is a reasonable key index, as most POI (Point Of Interest) information is associated with a street.

In this service scenario, a vehicle first decides a path plan by the start and the destination using a path finding algorithm such as A* embedded in the telematics device [14]. Then, an entire route consist of link IDs is submitted to the gateway, which searches its local cache for each link ID and waits for missing items to be valid. The gateway cannot reply until every requested link record is ready. Hence, the more the data miss, the poorer the response time, as more interaction with the remote server is needed. The vehicle can request the query only when it can reach a gateway and receive the reply when it remains connected until the gateway returns the reply. In depends on the driver's decision whether he will keep waiting when the response time gets too long. However, the driver at least knows the unavailability of information before his start.

3.2 Access Pattern Data

The realistic data access pattern is very important for the performance assessment of data prefetch schemes. This paper employs the location history data collected by the the Jeju taxi telematics system which constantly tracks the location of vehicles [15]. Each vehicle reports its GPS reading every minute to the central server via the cellular network. During the 2.3 months' test period which has the largest number of enrollments, up to 1.3 M records have been collected. The road network of Jeju area, being represented as a graph data structure, consists of about 17,000 nodes, or intersections and about 25,000 links, or the road that connects two nodes. Our analysis system processed the commercially used digital map to make a graph structure and detailed road shape easily accessible. We can estimate the route taken by a vehicle and assume that this vehicle retrieves the information on the links belonging to the route before it has started its trip.

The experiment is based on the trajectory of each travel obtained from the location history data. The path, the series of links, taken by a vehicle is submitted to the data retrieval service. The snapshot locations of a vehicle do not exactly represent the path actually taken, however, it is possible to find start and end points of a trip by investigating the change of taxi status connected to the tachometer. By applying the well-known A* algorithm [14] with the start and

destination of a trip, the whole set of links involved in the path taken by the car can be generated. Next, in deciding the location of a gateway, the start positions of all trips are first examined to find the area which has the maximum number of starts. Actually, the city hall area is such a spot. The experiment focuses on this single gateway, as gateways run independently in general. The simulation takes 2,190 trips (94,546 records) that start on one day from the location within 300 km radius from the gateway location, considering the transmission range of wireless communication interface.

4 Performance Measurement

This section measures and assesses the performance of the hybrid prefetch scheme via simulation using SMPL [16] which provides a lot of functions and libraries for discrete event scheduling, easily combined with the commonly used compilers. The experiment first measures the hit ratio and the instant response ratio for the proposed scheme, comparing with LRU, FAR-N, FAR-E, and random prefetch described in Section 2. Then, the reply ratio according to the satisfiability and the time tolerance will be measured. Finally, the instant reply ratio is further analyzed according to the time flow and the hop count, or the length of each request.

To focus on the prefetch strategy, the experiment also assumes that there is just one server and its contents are updated every ten minutes. It can be necessarily much shorter. After all, request arrivals, periodic updates at the server, and record refreshes at the gateway are integrated into a discrete event scheduler. In addition, the refresh time of each record distributes exponentially, considering the server-side load as well as the network delay between the gateway and an information server. To parameterize the refresh time, we define *update load*. If this value is 1.0, it takes one update period to refresh all link records, that is, it takes 10 minutes to refresh 25,000 records in the Jeju city scenario. If it is 5.0, only 20 % of entire link information can be downloaded to a gateway.

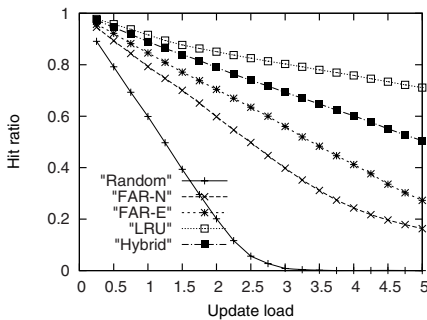


Fig. 2. Hit ratio

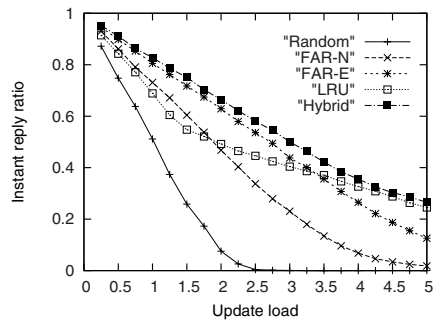


Fig. 3. Instant reply ratio

Figure 2 plots the record level hit ratio for each policy. For all schemes, the hit ratio decreases according to the increase of update load. LRU shows the highest hit ratio and decreases slowly according to the load increase, as it can best reflect temporal locality and refreshes first the link referenced at least once. Even in the road network, there exists a link the drivers are most likely to take. The hit ratio of the random scheme gets down to zero when the update load reaches 3.0, namely, only 1/3 of link data can be refreshed on average. The hybrid scheme is the next to LRU, as it can partially exploit the temporal locality. The hit ratio of the hybrid scheme is better than that of FAR-E by up to 23.1 % and worse than that of LRU by up to 20 % when the update load is 5.0.

Figure 3 shows the instant reply ratio for the respective strategies. The responsiveness of a gateway depends on how many link data are already refreshed when a request arrives. The gateway can reply instantly when all the requested link data are valid. This rate can be considered as the request-level hit ratio. Even though LRU has the highest record-level hit ratio, the hybrid scheme shows a better instant reply ratio by up to 20.3 % when the update load is 1.5. The gap gets smaller by up to 1.9 % along with the increase of the update load. Other spatial-locality-based schemes such as FAR-E and FAR-N also show better performance than LRU until the update load becomes 1.5 and 3.5, respectively. We can recognize that geographical affinity can better cover the whole link set in a trip than temporal locality. The gap between FAR-E and the hybrid scheme is 13.1 % at maximum.

Figure 4 plots the reply ratio according to the satisfiability, which denotes the bound a user satisfies. For the satisfiability of 0.9, a user is satisfied with the result if 90 % of link information is returned instantly. This scenario is more practical to the real-life service model. The experiment changes the satisfiability from 0.6 to 1.0, and the update load is set to 0.5. The hybrid scheme shows the best performance when the satisfiability lies between 0.8 and 1.0. However, LRU, having the highest link-level hit ratio, gets better when the satisfiability goes below 0.8. The hybrid scheme outperforms other spatial-locality-based schemes for all ranges. All schemes except LRU show a linear increase in the reply ratio.

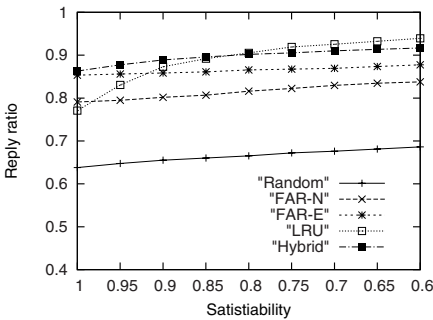


Fig. 4. Effect of satisfiability

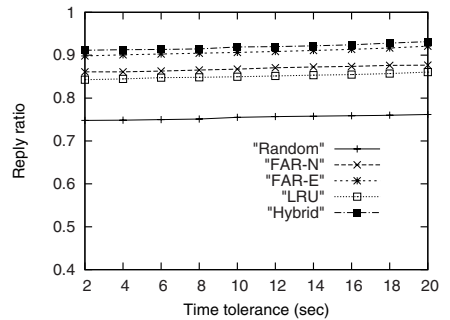


Fig. 5. Effect of delay tolerance

We think that the satisfiability is generally set to be above 0.9 for a user to get sufficient information for his/her route.

Figure 5 plots the reply ratio according to the time tolerance. If the tolerance is 2 seconds, a user will be satisfied with the result if he/she receives whole data records within 2 seconds. Here, the update load is set to 0.5 again. Interestingly, Figure 5 shows that the time tolerance has no significant influence to the instant reply ratio. This result indicates that refreshing the core set is most important to the performance. Anyway, the hybrid scheme shows the best instant reply ratio for all ranges, particularly outperforming FAR-E by up to 1.2% in this parameter setup.

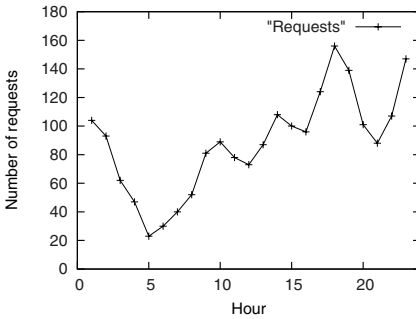


Fig. 6. Hourly number of requests

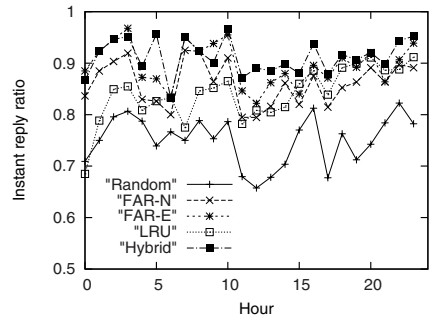


Fig. 7. Hourly instant reply ratio

Figure 6 and Figure 7 show the hourly instant reply ratio. The data access pattern has an important effect on the hit ratio and also the instant reply ratio. In our information service model, a vehicle requests a data record according to the destination of a trip. When trips have a common destination or more common route pieces, the instant reply ratio will be better. For example, during the commute, many vehicles head to an office area in the downtown and the residential area. Figure 6 shows the hourly change in the number of requests or trips. We can find some peaks on two commute times, namely just before and after the office hour around 9 and around 18, respectively, in addition to the shopping hours around 14. Around 5 o'clock, there are just 20 trips, and the performance in this period is a little bit meaningless.

Figure 7 traces the change of the hourly instant reply ratio. The initial reply ratio is relatively low as just a few previous access history data is available. So, the hit ratios of FAR-E and hybrid schemes are almost same, and so are those of random and LRU schemes. As the time passes by, each curve shows different behavior, but the hybrid scheme seems to be less affected by the change of access pattern. For 10 o'clock, just after the morning commute, the instant reply ratio decrease a little bit because the destination of each trip is quite different. Anyway, each scheme shows a slight drop in the reply ratio when the traffic pattern changes, for example, the morning and evening commute, shopping time, and personal time after office hour. The hybrid scheme can reach

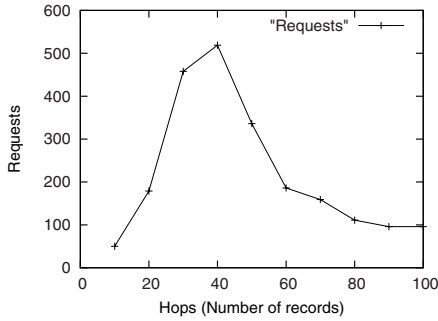


Fig. 8. Hourly number of requests

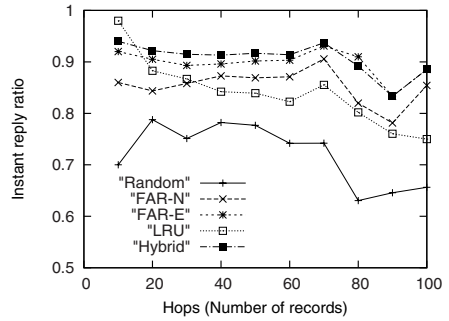


Fig. 9. Hourly instant reply ratio

the stable status after initialization. While LRU shows the 22.6 % of difference in the instant reply ratio, the hybrid scheme shows just 13.2 % difference.

Figure 8 and Figure 9 plot the effect of hop counts, which is analogous to the number of records in a request. The instant reply ratio is expected to decrease when a request has more links, as the gateway can reply instantly only if all of link records are valid. Among 2,190 requests, more than 500 requests contains 30 - 40 records, and the number of records in a request is concentrated in the range of 20 through 50. 50 requests has less than 10 links, while 100 requests consists of more than 100 hops, which means a vehicle must have taken a long trip. Figure 9 shows instant reply ratio according to the hop counts. Here, the update load is set to 0.5. We can recognize that random and LRU schemes are more affected by the number of records. On the contrary, the instant reply ratios of hybrid, FAR-N, FAR-E schemes do not significantly change until the hop count reaches 70. Hence, the number of records in a request is not a critical factor in the spatial locality-based prefetch schemes.

5 Conclusion

This paper has measured and analyzed the performance of the hybrid prefetch scheme for gateways on the vehicular telematics network via simulation using SMPL. The target service indexes every information by the link ID, analogous to the street ID or name, and updates at regular intervals. To take advantage of both LRU and FAR techniques, the hybrid scheme groups links according to whether they are referenced during the given time interval or not, orders by the Euclidean distance in each group, and fetches the referenced set first. The experiment results first have revealed that the hybrid scheme can improve the request-level instant reply ratio by up to 20.3 %, overcoming 20 % loss of record-level hit ratio, compared with LRU. Second, this scheme outperforms other spatial-locality-based schemes at all satisfiability ranges and LRU when the satisfiability is above 0.8, also showing stable behavior for the practical time tolerance range. Third and finally, it quickly stabilizes after the prefetch memory initialization and is less affected by the data access pattern change, demonstrating steady instant reply ratio for the number of records in each request.

As future work, we are to design a cooperative prefetch scheme that can balance the workload of each gateway, as load balancing can improve prefetch performance such as responsiveness and reply quality in most distributed applications [17]. To this end, localizing and identifying the core access set for each gateway is most important.

References

1. Said, E., Omar, E., Robert, L.: Data prefetching algorithm in mobile environments. *European Journal of Scientific Research* 28, 478–491 (2009)
2. Yu, B., Gong, J., Xu, C.: Data aggregation and roadside unit placement for a vanet traffic information system. In: *ACM VANET*, pp. 49–57 (2008)
3. Denko, M.: Cooperative data caching and prefetching in wireless ad hoc networks. *International Journal of Business Data Communications and Networking* 3, 1–15 (2007)
4. Liang, B., Drew, S., Wang, D.: Performance of multiuser network-aware prefetching in heterogeneous wireless systems. *Wireless Networks* 15, 99–110 (2009)
5. Sato, K., Koita, T., Fukuta, A.: Broadcasted location-aware data cache for vehicular applications. *EURASIP Journal on Embedded Systems* (2007)
6. Lee, J., Shin, I., Park, G., Kim, I., Rhee, Y.: Design of a cache management scheme for gateways on the vehicular telematics network. In: Gervasi, O., Taniar, D., Murgante, B., Laganà, A., Mun, Y., Gavrilova, M.L. (eds.) *Computational Science and Its Applications – ICCSA 2009*. LNCS, vol. 5592, pp. 831–840. Springer, Heidelberg (2009)
7. Zheng, B., Xu, J., Lee, D.: Cache invalidation and replacement strategies for location-dependent data in mobile environments. *IEEE Transactions on Computers* 51, 1141–1153 (2002)
8. Lee, J., Park, G., Kim, S., Kim, H., Shin, S.: A hybrid prefetch policy for the retrieval of line-associated information on vehicular networks. In: *ACM SAC (2010)* (accepted)
9. Gitzenis, S., Bambos, N.: Efficient data prefetching for power-controlled wireless packet networks. In: *International Conference on Mobile and Ubiquitous Systems: Networking and Services*, pp. 64–73 (2004)
10. Deshpande, P., Kashyap, A., Sung, C., Das, S.: Predictive methods for improved vehicular WiFi access. In: *ACM SIGMOBILE*, pp. 263–276 (2009)
11. Balasubramanian, A., Mahajan, R., Venkataramani, A., Zahorjan, J.: Interactive WiFi connectivity for moving vehicles. In: *ACM SIGCOMM*, pp. 427–438 (2009)
12. Zheng, B., Lee, W., Lee, D.: On semantic caching and query scheduling for mobile nearest-neighbor search. *Wireless Networks* 10, 653–664 (2004)
13. Society of Automotive Engineers: Dedicated short range communication message set dictionary. Tech. Rep. Standard J2735, SAE (2006)
14. Goldberg, A., Kaplan, H., Werneck, R.: Reach for A*: Efficient point-to-point shortest path algorithms. MSR-TR-2005-132. Microsoft (2005)
15. Lee, J., Park, G., Kim, H., Yang, Y., Kim, P., Kim, S.: A telematics service system based on the Linux cluster. In: Shi, Y., van Albada, G.D., Dongarra, J., Sloot, P.M.A. (eds.) *ICCS 2007*. LNCS, vol. 4490, pp. 660–667. Springer, Heidelberg (2007)
16. MacDougall, M.: *Simulating Computer Systems: Techniques and Tools*. MIT Press, Cambridge (1987)
17. Ting, Y., Chang, Y.: A novel cooperative caching scheme for wireless ad hoc networks: Group caching. In: *IEEE Int'l Conference in Networking, Architecture, and Storage*, pp. 62–68 (2007)

Power Control for Soft Fractional Frequency Reuse in OFDMA System*

Young Min Kwon, Ok Kyung Lee, Ju Yong Lee, and Min Young Chung**

School of Information and Communication Engineering
Sungkyunkwan University
300, Chunchun-dong, Jangan-gu, Suwon, Kyunggi-do, 440-746, Korea

Abstract. In this paper, soft Fractional Frequency Reuse (FFR) is considered to reduce the co-channel interference and also guarantee the cell edge users data rate. We apply power control in soft FFR and investigate the system performance. Soft FFR divides whole frequency band into multiple subbands with different kinds of transmission power levels and allocates subband according to SINR of cell users. By controlling allocated powers, we can increase the frequency efficiency by using the whole frequency band and also improve the data rate for the cell edge users remarkably by reducing the co-channel interference. If the throughput requirement for the outer region is small, we decrease the power allocated to the outer region which can save power energy while high overall system throughput is achieved.

Keywords: OFDMA, Frequency Reuses Scheme, FFR, SFR, Soft FFR.

1 Introduction

The next generation mobile broadband wireless communication systems, such as IEEE 802.16e Mobile WiMAX and 3GPP LTE (Long Term Evolution), are based on Orthogonal Frequency Division Multiple Access (OFDMA) to support the high data rate service [1][2]. Mobile WiMAX and 3GPP LTE also adopt the frequency reuse factor of one, in which each cell serves users with whole system bandwidth. In OFDMA, however, users with the same channel simultaneously suffer Co-Channel Interference (CCI) of neighbor cells, which may cause severe degradation of system performance, especially in cell edge region [3]. The CCI can be reduced with the frequency reuse factor of more than one. However, it may reduce the system capacity [4].

Several frequency reuse schemes have been studied to mitigate the interference of the cell edge region and to increase the system capacity [5]–[7]. The Fractional

* This work was supported by National Research Foundation of Korea Grant funded by the Korean Government (2009-0074466) and by the MKE (The Ministry of Knowledge Economy), Korea, under the ITRC (Information Technology Research Center) support program supervised by the NIPA(National IT Industry Promotion Agency (NIPA-2009-(C1090-0902-0005)).

** Corresponding author.

Frequency Reuse (FFR) [8] is proposed for mobile WiMAX in IEEE 802.16e working group, and Soft Frequency Reuse (SFR) [9] is proposed in the 3GPP working group. FFR scheme and SFR scheme divide the available spectrum into two reserved parts, subbands for the inner region and subbands for the outer region. The subband for inner User Equipment (UE) is common in each cell, and the subband for outer UEs is different among adjacent cells. Thus, UEs in outer region do not suffer the CCI from neighbor cells within first-tier, so that the spectral efficiency of outer region is increased. However, since FFR does not use whole available frequency bandwidth, overall cell throughput in a cell is lower than reuse one. SFR has two types of subband, the subband with high transmission power level and with low transmission power level. It assigns the subband with high power to outer users and the subband with low power to inner users. Since SFR can use whole system band in a cell, overall cell capacity in a cell is higher than FFR. However, overall system capacity of SFR may be lower than that of reuse one environment.

Soft FFR scheme has been proposed to improve overall cell throughput of FFR [10]. While FFR do not use the subbands allocated for outer region in the adjacent cells, soft FFR uses these subbands for the inner UEs with low transmit power. As a result, the soft FFR also use the subband with high transmit power level and with low transmit power level like SFR. Unlike SFR, soft FFR uses the common subband, which can guarantee the throughput of inner users. In this paper, we investigate the performance of soft FFR in the environment, where subbands have several transmit power levels. An outline of the rest of this paper is as follows. In Section 2, we propose a power control in soft FFR scheme. The result and analysis are then presented in Section 3. Finally, conclusions are drawn in Section 4.

2 Power Control in Soft Fractional Frequency Reuse (Soft FFR)

When frequency reuse factor is equal to one, UEs in the edge of cell may suffer high outage probability due to the CCI from adjacent cells. When frequency reuse factor is more than one, adjacent cells can operate on mutually orthogonal frequency, which yields small CCI from adjacent cells. However, since a cell can serve UEs with a part of whole system bandwidth, the spectral efficiency is low. To compromise these approaches, frequency reuse schemes such as FFR or soft FFR have been proposed. Both FFR and soft FFR statically partition the cell into two distinct geographical regions: the inner region and the outer region, as shown in Fig. 1.

FFR divides the system bandwidth into the subband for UEs in inner region and the subband for UEs in outer region. In Fig. 1, the system bandwidth is divided into four subbands, F0, F1, F2, and F3. In FFR, Base Station (BS) allocates the subband with different reuse factor according to the user's location in a cell. The reuse factor of F0 is equal to one and F0 is allocated to users in the cell's inner region. UEs in the inner region are close to the serving BS and

far from interfering BSs, and they are served with the common subband. Also one of three subbands, F1, F2, and F3 is allocated to UEs in the outer region and the other subbands are used for adjacent cells. Thus, there is no co-channel interference among frequency bands in the first-tier outer region. Since each cell cannot use the subband for the other region in adjacent cells, overall capacity of the system may be reduced. There is a trade-off between the spectral efficiency of users in the outer region and overall capacity.

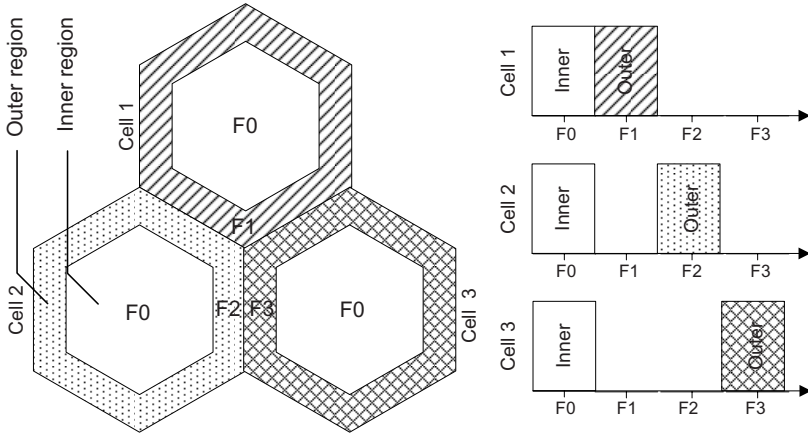


Fig. 1. Frequency allocation and power distribution of FFR scheme

In FFR, since a cell uses partial subbands among overall system frequency bands, the overall capacity may be decreased. To improve the overall capacity, soft FFR scheme has been proposed, which can serve UEs with the subband for the outer region of adjacent cell. As in Fig. 2, it divides the system bandwidth into four subbands, F0, F1, F2, and F3, like FFR. While FFR does not use the subband for outer region in the adjacent cells, soft FFR uses these subbands with transmit power level, P_{in} , for the inner UEs. In soft FFR, UEs in the inner region are served with additional subband with P_{in} as well as subband with P_{comm} . UEs in the outer region of soft FFR are served in the subband with transmission power level, P_{out} . Since soft FFR can use all subbands in each cell, the overall capacity of soft FFR may be larger than that of FFR. Due to the interference by the subband with P_{in} , soft FFR may have lower spectral efficiency of outer UEs as compared with FFR.

In soft FFR, we control the spectral efficiency of outer UEs by the transmission power ratio (ρ), which is the ratio of the inner transmission power level to outer transmission power level, that is, $\rho = P_{in}/P_{out}$. Soft FFR has common subband with transmission power level, P_{comm} . Since the common subband maintains specific transmission power and is only used by inner UEs, it can guarantee a certain amount of capacity for inner region.

Fig. 2 explains the control of the transmission power in soft FFR. In soft FFR, subbands with P_{in} are assigned to the inner UEs. On the other hand, subbands with P_{out} , F1, F2, or F3, are assigned to the UEs in outer region, which do not overlap with the subband of the adjacent cells. The SINR of outer UEs is controlled by P_{out} , which can be dynamically adjusted according to the population of outer UEs and the required data rate in outer region. In the next section, we investigate the effect of power control where P_{comm} has a fixed value and P_{out} is adjusted according to the requirements of outer UEs.

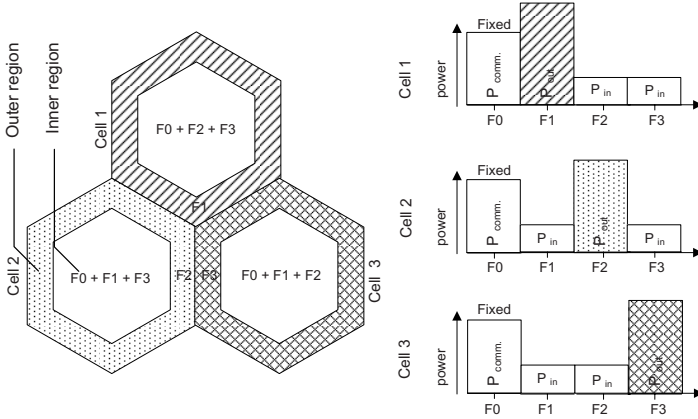


Fig. 2. Power control in soft Fractional Frequency Reuse (soft FFR)

3 Simulation Results

The effect of power control in soft FFR is evaluated through system level simulation. We consider the two-tier hexagonal cellular system with nineteen cells, where the inter-site distance is 500 meters. Each cell is divided into two regions: the inner region and the outer region. We assume that the ratio of the inner region area to the outer region area is two. The 48 UEs are located with a uniform distribution in a cell. Two thirds of whole UEs are in the inner region, and the rest are in the outer region. A BS allocates subband to UE by estimating the signaling from UEs. We assume that the BS knows a serving UE’s location exactly according to signaling. We use simulation parameters in Table 1 [2].

In the established system environment, we evaluate the average cell throughput of soft FFR in inner region and outer region. At a TTI of simulation, the BS in a cell gathers the location information of UEs and divides UEs into two groups, the inner group and the outer group, and it allocates subband resources to each UE using the proportional fair scheduling [11]. In soft FFR, the subband with P_{comm} and P_{in} are assigned to UEs in inner group, whereas the subband with P_{out} is allocated to UEs in the outer group. The throughput of an UE is obtained based on the Signal to Interference and Noise Ratio (SINR) of UE in

Table 1. Simulation Parameters

Parameter	Value
Channel bandwidth	10 MHz
Carrier frequency	2 GHz
FFT size	1024
Number of subcarriers	601
Subcarrier spacing	15kHz
Cellular layout	Hexagonal grid, 19 sites
Inter-site distance	500 m
Bandwidth	10 MHz
Log-normal shadowing	8dB
Penetration loss	20 dB
Propagation loss	$128.1 + 37.6 \log_{10}(R(km))$
BS antenna gain	15 dBi
UE antenna gain	0 dBi
White noise power density	-174dBm/Hz
Scheduling	Proportional Fair
TTI	1 ms

the assigned subband. In system level simulation, SINR is determined by the path loss and lognormal fading measured in the subband. The throughput of an UE m is estimated using the Shannon capacity as

$$T_m = W_{sub} \log_2(1 + SINR_m) \quad (1)$$

where W_{sub} is the bandwidth of a subband assigned to an UE and $SINR_m$ is the SINR of an UE m . The cell throughput in each region is total throughput of UEs in the corresponding region and expressed as

$$T_{cell} = \sum_{m=1}^M T_m \quad (2)$$

where M is the number of UEs in a group.

To compare the cell throughput of soft FFR with that of FFR, we assume that each frequency reuse scheme can use the total transmission power of $40W$. FFR divides the total transmission power into two and allocate the power to inner subband and outer subband. In soft FFR, we allocate a quarter of total transmission power to common subband. The remaining transmission power is distributed to the subband with P_{in} and P_{out} according to the power ratio ρ . Fig. 3 shows the cell throughput of the inner and outer region as the power ratio changes. Since soft FFR can use the overall system band, cell throughputs in soft FFR are larger than those in FFR. As we reduce the power ratio (ρ) to improve the throughput in outer region, the overall throughput and the inner throughput are decreased. With small ρ , the first tier cells use the subband with lower power and the interference of first tier is reduced. However, since the six

cells of the second tier use the same subband with the target cell, the interference of the second tier is increased. As a result, cell throughput is affected more by the interference from the second tier than the first tier. Therefore, as ρ decreases from 0.1 to 0.001, the outer cell throughput of soft FFR is not increased any more although the inner cell throughput is decreased.

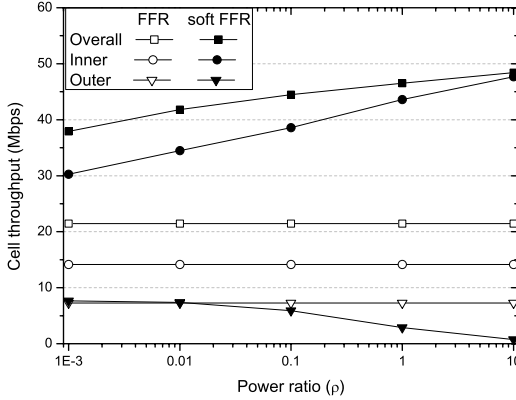
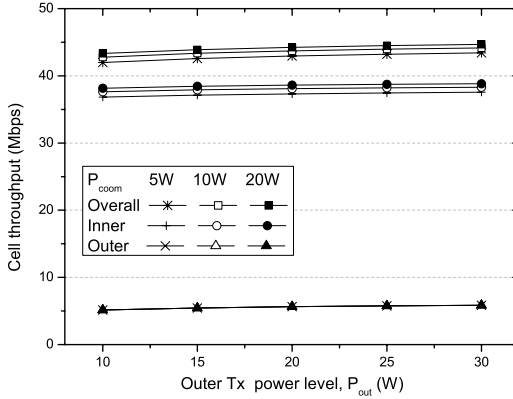


Fig. 3. Cell throughput of FFR and soft FFR according to power ratio, $\rho = P_{in}/P_{out}$ ($P_{comm} = 10W$, total transmission power = $40W$)

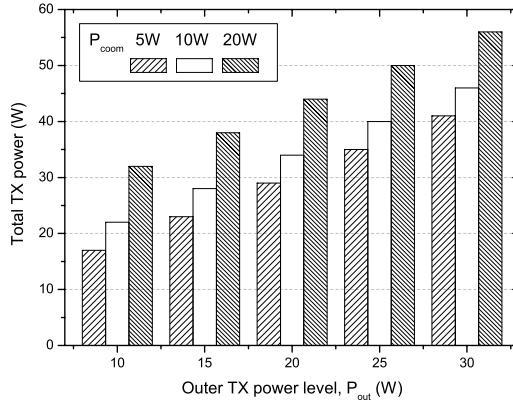
Now we consider the environment that soft FFR has the fixed power ratio, $\rho = 0.1$. When the transmission power of common subband is fixed to $5W$, $10W$, and $20W$, we increase the power of subband for outer UEs, P_{out} , and estimate the cell throughput in each region as P_{out} changes. Since the power ratio ρ is constant, the transmission power level of subband for inner UEs, P_{in} , is also increased as P_{out} increases. Fig. 4 (a) shows the cell throughput of each region in this environment. In the inner region, although P_{comm} is increased, the cell throughput is not increased significantly. Since the common subband has frequency reuse factor with one, the signal power of serving BS and interfering BSs are simultaneously increased as P_{comm} increases. As a result, the SINR of UEs using the common subband is less sensitive to the change of transmission power. In the outer region, since the power ratio (ρ) is fixed, the cell throughput does not vary although the outer transmission power level is increased. The overall cell throughput is less affected by the transmission power level of common subband, P_{comm} , with constant power ratio (ρ).

Fig. 4 (b) shows that total transmission power increases as P_{comm} increases, when ρ is fixed to 0.1. However, the cell throughput remains almost the same. Therefore, while we can achieve similar throughput performance, we can save the transmission power by reducing P_{comm} .

To investigate the power control by adjusting P_{out} , we assign the fixed transmit power to P_{comm} and P_{in} , which are $10W$ and $5W$, respectively. In this



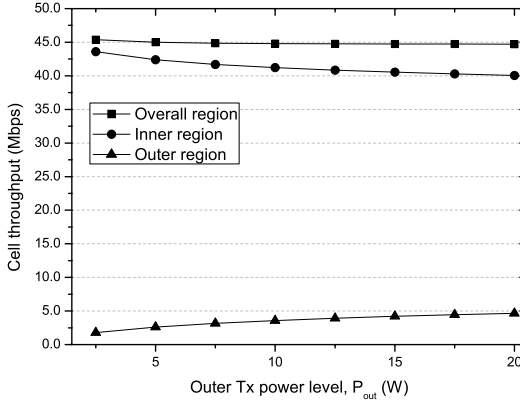
(a) Cell throughput



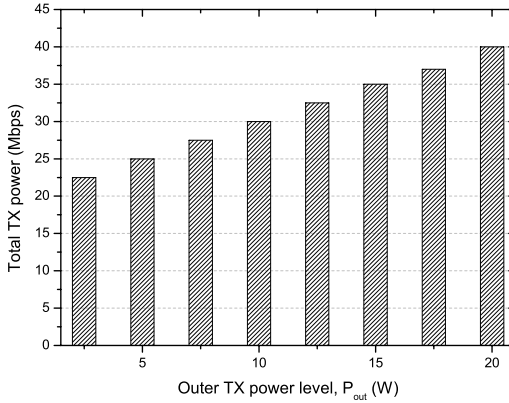
(b) Total transmission power

Fig. 4. Cell throughput and total transmission power of soft FFR according to the outer transmission power level, P_{out} ($\rho = P_{in}/P_{out} = 0.1$, $P_{comm} = 5W, 10W, 20W$)

environment, we estimate the cell throughput in each region and the total transmit power according to outer transmit power level, P_{out} . In Fig. 5 (a), when P_{out} is increased, the cell throughput in outer region is increased due to the increased signal power. On the other hand, the cell throughput in inner region is decreased due to the increased interference. However, the total transmission power increases to improve the cell throughput of outer UEs in Fig. 5 (b). Thus, P_{out} can be controlled to the throughput requirement of outer UEs in soft FFR. If a few UEs exist in outer region, we reduce the transmission power of the sub-band used in outer region, which yields high throughput for inner region due to



(a) Cell throughput



(b) Total transmission power

Fig. 5. Cell throughput and total transmission power according to P_{out} in soft FFR ($P_{comm} = 10W$, $P_{in} = 5W$)

the reduced interfering signal power. As a result, the overall throughput can be increased and the soft FFR scheme can save the total transmission power.

4 Conclusion

In this paper, we investigated the effect when the power of subband for each region changes in soft FFR. In the simulation, we observe that the transmission power level of common subband, P_{comm} , does not influence the overall cell throughput significantly. Thus, we save the total transmission power by reducing

P_{comm} within the required cell throughput. In addition, we also save the total transmission power by using the relation between P_{in} and P_{out} . If the power in the outer region increases, the throughput in the outer region increases and the throughput in the inner region decreases. Otherwise, the throughput in the outer region decreases and the throughput in the inner region increases. As the power in the outer region decreases, the overall capacity increases, since the spectral efficiency in the inner region is higher than that in the outer region. According to the throughput requirements in the outer region, we can reduce the power consumption for the outer region, while maintaining high system throughput.

References

1. IEEE 802.16e-2005: IEEE Standard for Local and Metropolitan Area Networks, Part 16: Air Interface for Fixed and Mobile BroadBand Wireless Systems (October 2005)
2. 3GPP TR 25.814 ver 7.1.0, Physical layer aspects for evolved Universal Terrestrial Radio Access (UTRA) (Release 7) (September 2006)
3. Giuliano, R., Monti, C., Loreti, P.: Wireless technologies advances for emergency and rural communications - WiMAX fractional frequency reuse for rural environments. *IEEE Wireless Communications* 15(3), 60–65 (2008)
4. Wang, Z., Stirling-Gallacher, R.A.: Frequency reuse scheme for cellular OFDM systems. *Electronics Letters* 38(8), 387–388 (2002)
5. IEEE 802.20 MBWA C802.20 – 05 – 69, Air Interface Spec Final Fixed (November 2005)
6. 3GPP R1-050599, Interference Mitigation Considerations and Results on Frequency Reuse, Siemens (June 2005)
7. 3GPP R1-050833, Uplink Interference Mitigation via Power Control, LG Eletrincs (May 2006)
8. Mobile WiMAX Part I: A Technical Overview and Performance Evaluation, WiMAX Forum (February 2006)
9. 3GPP R1-050507, Soft frequency reuse scheme for UTRAN LTE, Huawei, Tech. Rep. (May 2005)
10. IEEE 802.16 Broadband Wireless Access Working Group, IEEE C802.16m-08/782, Fractional Frequency Reuse in Uplink, LG Electronics (August 2008)
11. Andrews, M.: Instability of the proportional fair scheduling algorithm for HDR. *IEEE Transactions on Wireless Communications* 3(5), 1422–1426 (2004)

Authentication – Based Medium Access Control to Prevent Protocol Jamming: A-MAC

Jaemin Jeung, Seungmyeong Jeong, and Jaesung Lim*

Graduate School of Information and Communications, Ajou University, South Korea
{mmsg, aflight, jaslim}@ajou.ac.kr

Abstract. Recently, Wireless Local Area Network (WLAN) is used by enterprises, government, and the military, as well as small office and home offices. Although it is convenient, it has inherent security weaknesses due to wireless characteristics. For this reason, security-sensitive groups are still unwilling to use WLAN. There are several types of attacks to degrade the wireless network throughput using security weakness, especially Protocol Jamming Attacks are critical. These attacks consume little energy and can be easily implemented. In this paper, we introduce Authentication-based Medium Access Control (A-MAC) to prevent Virtual Carrier Sense (VCS) Jamming Attack and Deauthentication / Disassociation Jamming Attack that are typical Protocol Jamming Attacks in 802.11 based wireless systems. The proposed scheme can authenticate frames using Universal Hashing Message Authentication Codes (UMAC-32) and Hidden Sequence Number (SN). A-MAC frequently changes the key and SN using shift row and shift column processes to overcome the weakness of short 32 bits hashing codes. A-MAC can achieve integrity, authentication, and anti-replay attack security features. It can prevent Protocol Jamming Attacks that degrade wireless network throughput. Our simulation shows A-MAC can sustain throughput under Protocol Jamming Attacks.

Keyword: Protocol Jamming, Authentication, Medium Access Control.

1 Introduction

Wireless Local Area Network (WLAN) provides users with mobility within a broad coverage area with connection to networks. In addition, it is cost-efficient, allowing ease of integration with other networks and network components. For these reasons, WLAN is used by enterprises, government, and the military, as well as small office and home offices. However, it has inherent security weaknesses due to wireless characteristics. Furthermore, it uses the Industrial Scientific Medical (ISM) band to be highly susceptible to interference. Security-sensitive groups are still unwilling to use WLAN due to these weaknesses. If we use open source scanning software, such as

* This research was supported by the MKE (The Ministry of Knowledge Economy), Korea, under the ITRC (Information Technology Research Center) support program supervised by the NIPA (National IT Industry Promotion Agency) (NIPA-2009-C1090-0902-0003).

KISMET, we can easily analyze the packets [1]. Many researchers have been studying these issues, since these weaknesses are well known. Thus, WEP, WPA, TKIP, 802.11i standards were formed. However, the standards cannot resolve Protocol Jamming Attacks. Furthermore, Protocol Jamming Attacks consume little energy, and can be easily implemented. An adversary can decrease the wireless networks throughput for a long time, especially, when jammers are hidden.

Virtual Carrier Sense (VCS) jamming and De-authentication / Disassociation jamming are typical Protocol Jamming Attacks. VCS jamming attacks exploit a weakness in the Carrier Sense Multiple Access / Collision Avoidance (CSMA/CA). CSMA/CA systems use VCS to avoid collision. An adversary exploits the VCS process to decrease the wireless networks throughput. De-authentication / Disassociation jamming attacks exploit the protocol concerned with connection in the mobile station and AP. All stations have to go through the process of authentication / association before sending data. If an AP sends a De-authentication / Disassociation frame to a station, the stations have to terminate the connections. Unfortunately, based on the 802.11 standards, the station cannot reject the notification of these frames. Therefore, these attacks are more effective and powerful than VCS attacks [2].

We consider that the cause of Protocol Jamming Attacks is authentication problems. That is, APs and mobile stations cannot authenticate RTS / CTS / De-authentication / Disassociation frames. Even if there is an authentication process, it can be attacked by replay-attacks. Replay attacks also decrease network throughput. Thus, we need authentication, integrity and anti-replay-attacks features to prevent Protocol Jamming Attacks. In this paper, the proposed A-MAC substitutes digested message for Cyclic Redundancy Check (CRC). The digested message is created when the Universal Hashing Message Authentication Codes (UMAC-32) are XORed with a Hidden Sequence Number. UMAC-32 ensures the authenticity and the integrity of transmitted messages. In addition, it can check the transmission error. So we can substitute the digested message for CRC. The Hidden Sequence Number for anti-replay attacks is not seen in the transmitted frames. It is stored in a Sequence Table of each correspondent. In this approach, we can achieve integrity, authentication and anti-replay attacks security features, and prevent Protocol Jamming Attacks that degrade network throughput, as well as checking for transmission errors, without frame size overhead. We also introduce Key Shift & Sequence Number Select Notification (KSSN) that notify key and SN change by shift row and shift column processes, so that we compensate for the short 32bits digested message. Using the KSSN, each correspondent can change the key and SN, if necessary. That is, the KSSN is verified by the A-MAC, each correspondent changes the key and SN. This method is suitable for control and management frames that are short frame and occur frequently, because the KSSN can use a one-way key exchange mechanism, instead of four-way handshaking. Thus, we can prevent VCS and De-authentication / Disassociation jamming attacks using the A-MAC. Our simulation shows its performance and effectiveness.

The remainder of this paper is organized as follows. Section 2 describes the background of the proposed scheme and related work. We introduce adversary models, A-MAC and KSSN process in Section 3. Security analysis and performance analysis are presented in Section 4. Finally we conclude the paper and outline directions for future work in Section 5.

2 Background and Related Work

2.1 Background

802.11 wireless networks use Carrier Sense Multiple Access with Collision Avoidance (CSMA/CA). The Distributed Coordination Function (DCF) is the basis of the standard CSMA/CA access mechanism to avoid collisions. It first checks to see that the radio link is clear before transmitting. Stations use a random back-off after each frame, with the first sender seizing the channel, to avoid collisions. The DCF may use the Request To Send (RTS) / Clear To Send (CTS) clearing technique to further reduce the possibility of collisions.

Carrier sensing is used to determine if the medium is available. Two types of carrier sensing functions in 802.11 manage this process: Physical Carrier Sensing and Virtual Carrier Sensing. If either carrier sensing function indicates that the medium is busy, the MAC reports this to higher layers [3]. Unlike wired communications, Physical Carrier Sensing cannot provide all the necessary information, due to the hidden node problem. In Fig. 1, A is about to send frames to B. However, D cannot recognize the situations, because A's physical signal cannot reach D. Since A and D cannot hear each other, they may sense the media is free at the same time and both try to transmit simultaneously; this causes collision in the network. This is the hidden node problem. To reduce this problem, 802.11 standards adopt RTS, CTS, and VCS mechanisms. Fig. 1 shows these mechanisms. Activity on the medium by stations is represented by the shaded bars, and each bar is labeled with the frame type. Inter-frame spacing is depicted by the lack of any activity. DIFS is Distributed Inter-Frame Space and SIFS is Short Inter-Frame Space. If A wants to send frames to B, A sends the RTS frame to B, and then B replies with the CTS frame. Neighbor stations, C and D can overhear RTS and CTS frames, and then defer access to the medium until the

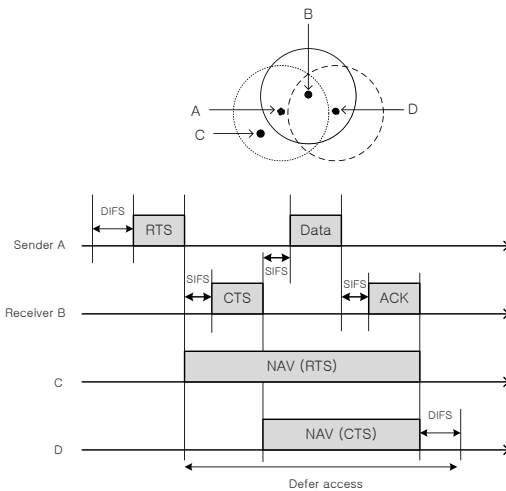


Fig. 1. Virtual Carrier Sense

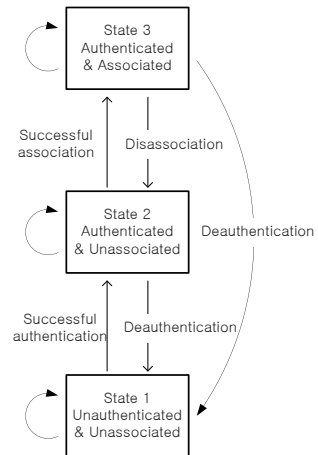


Fig. 2. 802.11 state diagram

Network Allocation Vector (NAV) elapses. VCS is provided by the NAV. The NAV is a timer that indicates the amount of time the medium will be reserved, in microseconds. It is usually presented by 2 bytes. If a station sets the timer for which they expect to use the medium, other stations count down from the NAV to 0. When the NAV reaches 0, the VCS indicates that the medium is idle. These RTS, CTS, and VCS mechanisms will prevent collisions.

Stations have to go through three states for connection to Access Point (AP). Fig. 2 shows the state diagram of 802.11. Stations are either authenticated or unauthenticated and can be associated or unassociated. These situations can be combined into three states. State 1 is not authenticated and not associated, State 2 is authenticated but not associated, and State 3 is authenticated and associated. Each state is a successively higher point in the development of an 802.11 connection. All stations start in State 1, and data can be transmitted only in State 3. That is, any stations cannot send data in State 1 and State 2. When a station transfers from one state to another state, it would be vulnerable.

2.2 Related Work

Any current or upcoming 802.11 standards would not help mitigate the risk of Protocol Jamming [4]. Researchers have tried to resolve these problems recently. Zhou proposed a packet-by-packet authentication method [5]. The author encrypted all frame content by a secret key, and then attached the encrypted content to the end of the original frame. For example; $A \rightarrow B: \{RTS, E(sID, dID, TS, SN)_k\}$, $B \rightarrow A: \{CTS, E(sID, dID, TS, SN+1)_k\}$. User A decrypts the encrypted attachment, verifies sID, dID, TS, and SN+1. If the frame is correctly decrypted, A can commence data transmission. However, this authentication scheme increases transmission overhead. Assume TS and SN are 4 bytes each, the total encrypted attachment is 20 bytes (6 bytes each for sID and dID). Considering RTS / CTS frame are just 20 bytes / 16 bytes, the encrypted attachment is a considerable overhead. Karlof proposed the TinySec packet format for wireless sensor networks [6]. However, it has a weakness in a replay attack, because it does not use any kind of sequence number or time-stamp in the TinySec-Auth packet. So, TinySec is not suitable in WLAN. Bellardo has a different view point [7]. In contrast to the authentication scheme, this scheme places two different limits on the duration values accepted by stations. The low limit has a value equal to the amount of time required to send an ACK frame, plus media access back-off for that frame. The high limit has a value equal to the amount of time required to send the largest data frame, plus the media access back-off for that frame. However, this scheme still has a weakness in RTS/CTS flooding. Since each station cannot verify the adversary that frequently send false RTS/CTS frames, the stations defer their transmission, and consequently the wireless networks throughput is drastically lowered. Some research prevents Protocol Jamming Attacks using the creation of a series of protocol extension and replacements (e.g., WEP, WPA, 802.11i, 802.11w) [2]. However, these schemes need complicated key management, powerful computational process and frame size overhead. They are not easy to implement.

3 Authentication-Based Medium Access Control (A-MAC)

3.1 Malicious Adversary Models

Malicious adversary models are VCS jamming attacks and De-authentication / Disassociation jamming attacks. VCS jamming attacks exploit the weakness of the virtual carrier sensing mechanism, especially RTS/CTS process. De-authentication / Disassociation jamming attacks exploit the weakness of the state diagram in Fig. 2. When a station transfers from one state to another state, a vulnerable point occurs, especially in the De-authentication, Disassociation frames.

Fig. 3(a) shows RTS/CTS adversary and NAV adversary in VCS jamming attacks. If a malicious adversary frequently sends a RTS frame to A, A sets the NAV timer, as much as the duration field of RTS frame indicates. According to the VCS mechanism, A has to defer access to the medium until the NAV elapses, since the VCS mechanism indicates that the medium is busy. If a malicious adversary repeatedly sends a RTS frame to A, station A cannot have the opportunity to access the medium. This scenario will decrease the total wireless network throughput. We define this adversary as the RTS/CTS Adversary. In the case of Station B, a malicious adversary sets the NAV timer for a longer period. According to VCS mechanism, B simply thinks that the other station is about to send a large data set. As a result, B has to defer access to the medium during the time of the duration field. The duration field has 16 bits. An adversary can set the NAV duration for a 2^{16} time slot. Assume that one NAV time slot is $1 \mu s$; B has to wait approximately 0.06 second. This delay is considerable. If there are many stations around the malicious adversary, the total wireless network throughput is drastically lowered. We define this adversary as the NAV adversary. These two adversary models apply to the CTS transmission process too.

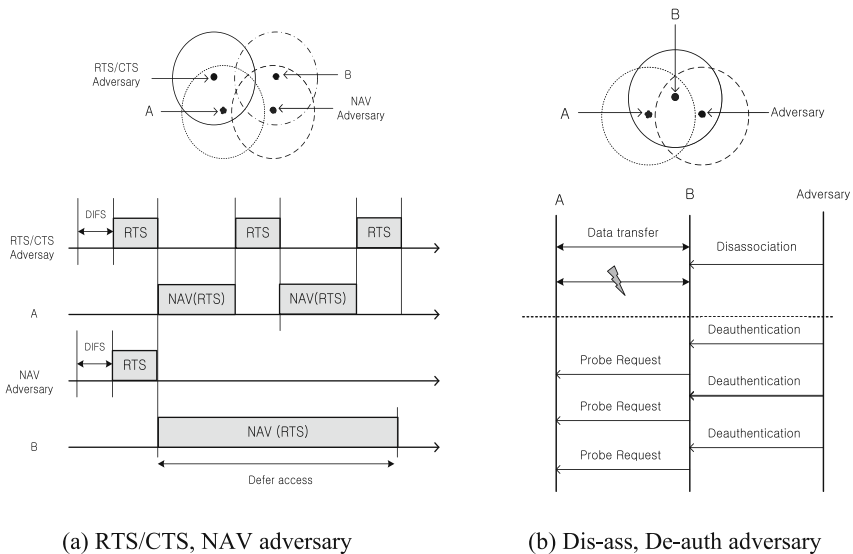


Fig. 3. Malicious adversary models

Fig. 3(b) shows Disassociation/De-authentication jamming attacks. These attacks exploit the weakness of the connection process. We have already learned that each station goes through States 1, 2, and 3. If a station gets a de-authentication frame when transferring from state 1 to State 2, the station cannot transfer to State 2. For the same reason, if a station gets a disassociation frame when transferring from State 2 to State 3 or sending data to the AP, the station reverts to State 2. According to the 802.11 standards, stations cannot reject the notification of disassociation and de-authentication frames. Unfortunately, a malicious adversary can easily fake disassociation and de-authentication frames. In Fig. 3(b), assume that station B is connected to Access Point A. If B gets a disassociation frame, B has to terminate the data communication. We define this attacker as the Disassociation Adversary. The adversary model is simple but, powerful. The total wireless throughput is to be a zero by only a few frames / second. After terminating the connection, B tries to reestablish a new connection. If a malicious adversary frequently sends a de-authentication frame at that time, B cannot transfer to State 3. We define this attacker as the De-authentication Adversary.

We define four adversary models. The characteristics of these models are easy implementation, use low energy, and decrease the total wireless network throughput. It is not easy for wireless IDS to detect the jammer, given that the jammer is hidden, uses low signal power, and sends fake frames not frequently but appropriately. Thus, the jammer can sustain a low network throughput that irritates network users and managers for a long time.

3.2 A-MAC Scheme

We propose the Authentication-based Medium Access Control (A-MAC) that achieves authentication, integrity and anti-replay attack without overhead to prevent Protocol Jamming Attacks that degrade throughput. A-MAC compresses the frames into a 32 bits digested message using Universal Hashing Message Authentication Codes (UMAC-32). Then, the digested message is masked by XOR with Hidden Sequence Numbers. Finally, we substitute the results for Cyclic Redundancy Check (CRC). In this scheme, we easily detect some abnormal frames that considerably decrease throughput. Thus, the total wireless network throughput does not decrease by discarding the instruction of abnormal frames. Fig. 4 details the proposed A-MAC process.

In Fig. 4, we assume that sender A transfers a message (M) to receiver B, and M is RTS, CTS, Disassociation, or De-authentication frame without CRC. A-MAC compresses the M into a 32 bits digested message by the UMAC-32 function and the secret key. Then, the compressed message is XORed with the Hidden Sequence Number. The Hidden SN is not seen in the transmitted frames. It is stored in a Sequence Table. The Hidden SN is very important in this scheme. It prevents a replay attack and complicates the digested message. Finally, the XORed message $\{UMAC(K, M) \oplus SN\}$ is concatenated with M. Sender A sends this frame to receiver B. B also compresses the received M into a 32 bits digested message by the UMAC-32 function and the secret key. Then, the compressed message is XORed with the received $\{UMAC(K, M) \oplus SN\}$. Receiver B derives SN from this process $[UMAC(K, M) \oplus \{UMAC(K, M) \oplus SN\} = SN]$. The derived SN is compared with SN of B's Sequence Table. If

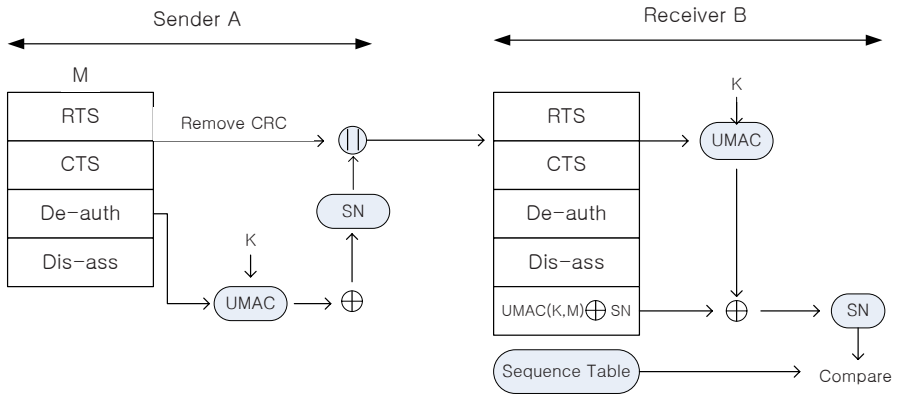


Fig. 4. Authentication-based Medium Access Control Process

the SNs agree with each other, the sender A is proved to be a legitimate station. If the SN does not agree with each other, there are either transmission errors or Protocol Jamming Attacks. Transmission errors are closely related to Signal to Noise Ratio (SNR). We can easily check the SNR in 802.11 systems by checking the Received Signal Strength Indication (RSSI). Therefore, although the RSSI is good, the SN does not agree with each other for several times. The probability of Protocol Jamming Attacks becomes high.

We can apply the A-MAC to an infrastructure WLAN. For example, an adversary sends de-authentication and disassociation frames to a certain station, the station can verify the frames using the A-MAC mechanism. If those frames are verified, they follow the instruction of each frame, otherwise they discard the frames. In case of sending lots of false RTS / CTS frames to a certain station, counter measures will be different. That is, key distribution schemes can make the method different. For example, an AP and stations share the same secret key, the stations can derive the SN from the frames using A-MAC. If the SNs of successive frames are incremented by one, the station will be legitimate, otherwise the stations will be suspected as an adversary. In contrast to the prior method, if the AP and each station share the 1:1 matching secret key, neighbor stations cannot derive the SN due to a different key. But the AP can derive the SN, and verify the frames. In that case, suppose the AP recognizes the adversary, it can notify the protocol jamming situation to the stations. The stations received the notification can discard the frames and cope with the situation.

A-MAC uses UMAC-32 for the digested message, instead of CRC, since UMAC-32 ensures the authenticity and the integrity of transmitted messages. Furthermore, it is the fastest message authentication code reported in cryptographic literature [8]. We can derive Tag from the UMAC-32 process. Tag's size is 32 bits, 64 bits, 96 bits, 128 bits. In this paper, A-MAC use a 32 bits Tag so that the digested message is XORed with a 32 bits Hidden SN. $Tag = H_{K1}(M) \oplus F_{K1}(nonce)$, M is an input message, H is a secret hash function, F is a pseudo random function, and K1 is secret random key shared by sender and receiver. Nonce is a value that changes with each digested message. In this paper, M is the content of the CRC removed frame, nonce is the Hidden SN. The proposed scheme can achieve authentication, integrity and anti-replay

attack without overhead. Therefore, A-MAC can overcome the weakness of Karlof's scheme [6] and Bellardo's scheme [7]. Especially, it does not expand the original frame size. This characteristic also overcomes the problem of Zhou's scheme [5].

3.3 Key Shift and SN Select

As using a short 32 bits digested message, A-MAC may be attacked by brute-force attack that is a technique to defeat the authentication mechanism by trying successively all the words in an exhaustive list. We propose two methods that complicate the attached 32 bits message of A-MAC to mitigate this vulnerability.

At first, A-MAC does not count the initial SN from zero. If the initial SN starts from zero, a malicious adversary can calculate $\{UMAC(K, M)\}$. Thus, the probability of detecting the key becomes high. A-MAC derives the initial SN from the 128 bits key to mitigate this probability. Table 1. is an example of Sequence Table. We make a 4×4 byte matrix using the 128 bits secret key. Each row and each column is 32 bits. The first row becomes the initial SN of the RTS frame, the second row becomes the initial SN of the CTS frame, the third row becomes the initial SN of the de-authentication frame, the fourth row becomes the initial SN of the disassociation frame. After being set, the initial SN is incremented by one every frame exchange. In this scheme, it is not an easy for a malicious adversary to estimate the key and the digested message.

Second, each correspondent (station or AP) can change the secret key and SN if necessary. We define this process as Key Shift & SN Select Notification (KSSN). Fig. 5 shows that each correspondent changes the secret key and SN simultaneously.

Table 1. Sequence Table

Sequence Table					
Node ID	Key	RTS	CTS	De-auth	Dis-ass
A	ABCD EFGH I JK L MNOP	ABCD	EFGH	I JK L	MNOP
B	BCDE FGH I JKLM NOPA	BCDE	FGH I	JKLM	NOPA
*	*	*	*	*	*
*	*	*	*	*	*

In Fig. 5, we assume that sender A wants to change the secret key and SN. A sets the subtype of Frame Control Header to 0111. In 802.11 standards, management subtype 0111 and control subtype 0111 are reserved and not currently used [3]. In this paper, we use this subtype 0111 to send the KSSN message. After setting the subtype, sender A transmits $\{UMAC(K, M) \oplus SN \oplus KSSN\}$ to receiver B. The received $\{UMAC(K, M) \oplus SN \oplus KSSN\}$ message is XORed with $\{UMAC(K, M) \oplus SN\}$ by receiver B. B derives the KSSN from this process. B compares the last 13 bits of KSSN with the last 13 bits of SN.

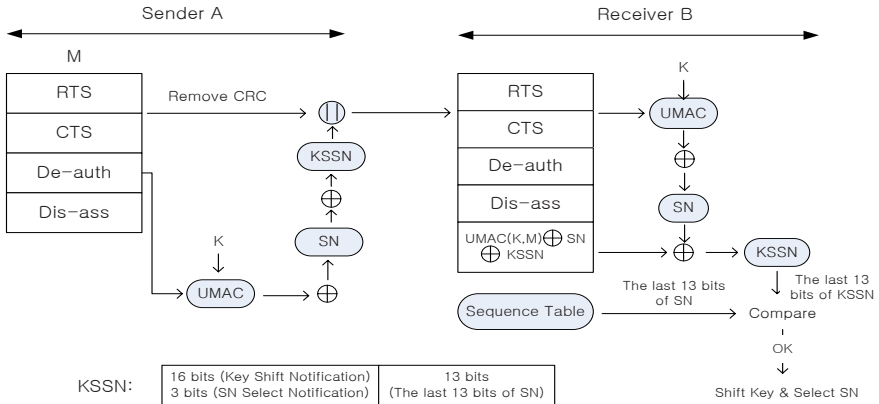


Fig. 5. Key Shift & SN Select Notification Process

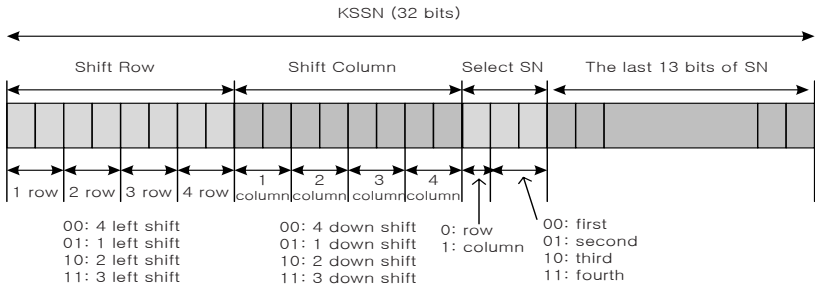


Fig. 6. KSSN format

Fig. 6 shows the 16 bits Key Shift Notification, 3 bits SN Select Notification, and the last 13 bits SN KSSN format. We form the 4 × 4 byte matrix using the 128 bits secret key, as in Fig. 7. The front 8 bits are related to Shift Row. The first 2 bits indicate the first row, 3 ~ 4 bits indicate the second row, 5 ~ 6 bits indicate the third row, 7 ~ 8 bits indicate the fourth row. 00 indicates four bytes left shift, 01 indicates one byte left shift, 10 indicates two bytes left shift, 11 indicates three bytes left shift. The 9 ~ 16 bits are related to the Shift Column. The allocation method is similar to Shift Row. The 17 ~ 19 bits indicate which row or column is selected as the new SN. If the 17th bit is 0, it selects the new SN from the row. If it is 1, it selects the new SN from the column. 00 indicates the first row or column, 01 indicates the second row or column, 10 indicates the third row or column, 11 indicates the fourth row or column. The last 13 bits are used for authentication. If the last 13 bits of KSSN and the last 13 bits of SN are verified, the receiver can change the secret key and SN. The sender transmits the KSSN several times to overcome transmission error. This depends on channel states.

Fig. 7 shows the Key Shift and SN Select operations. If the receiver derives the KSSN (KSSN: 00011011 00011011 010 + 13 bits) from the received message, the receiver verifies the last 13 bits of KSSN comparing them to the last 13 bits of SN.

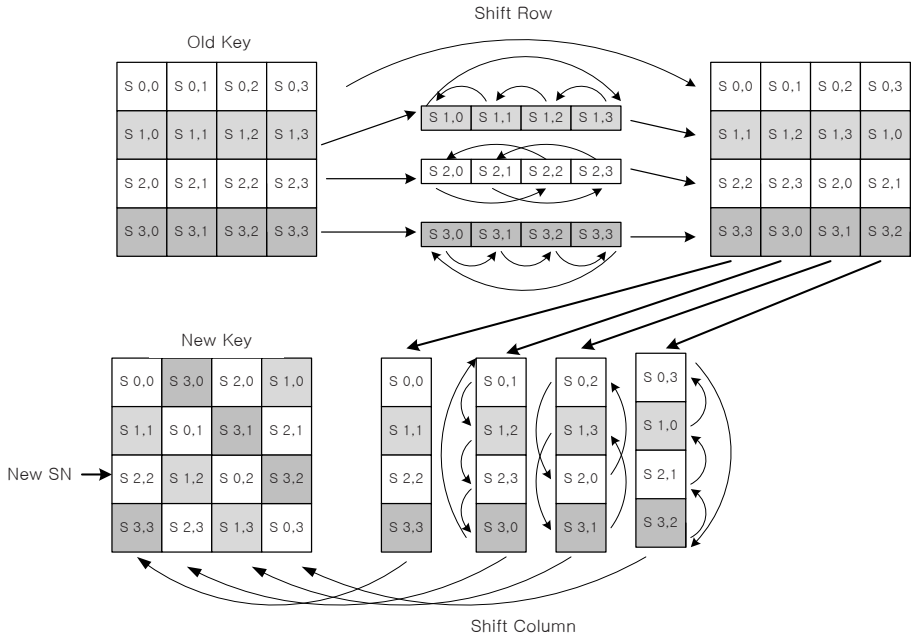


Fig. 7. Key Shift & SN Select Operations

If they are verified, the receiver changes the secret key and SN. The first 8 bits (00011011) indicate the first row is shifted by 4 bytes, the second row is shifted by 1 byte, the third row is shifted by 2 bytes, and the fourth row is shifted by 3 bytes. The next 8 bits (00011011) indicate that the first column is down-shifted by 4, the second column is down-shifted by 1 byte, the third column is down-shifted by 2 bytes, and the fourth column is down-shifted by 3bytes. As a result of the operations, the new key is generated. 17~19 bits (010) indicates that the third row is the new SN of the RTS frame. Then, other rows are sequentially allocated to other frames (e.g. the fourth row is the new SN of CTS frame, the first row is the new SN of de-authentication frame, and the second row is the new SN of association frame).

4 Security and Performance Analysis

4.1 Security Analysis

We designed the A-MAC to achieve authentication, integrity and anti-replay attack without overhead. The A-MAC compresses the frames into a 32 bits digested message using UMAC-32, and then the digested message is masked by XOR with Hidden Sequence Numbers. We obscure the digested message using KSSN operations to overcome the weakness of the short digested message (e.g. brute-force attack). In this process, we can obtain some security advantages. First, to prevent replay-attacks, we use a Hidden SN that makes it hard to estimate the next frame. If it were not for the

Hidden SN, a malicious adversary could perform Denial of Service (DoS) attacks. Assuming that A transmits a RTS control frame to B. $A \rightarrow B: \{FC \text{ (Frame Control Header)}, D \text{ (Duration)}, RA \text{ (Receiver Address)}, TA \text{ (Transmitter Address)}, SN, CRC\}$. An adversary can attempt to predict the SN, and then frequently sends a RTS control frame like this. Adversary $\rightarrow B: \{FC, D, RA, TA, SN + (i = i+1), CRC\}$. B will frequently send CTS or ACK frame. As a result, the total wireless network throughput will be lowered. However, the A-MAC can prevent these DoS attacks using the Hidden SN. Furthermore, it uses the 4 byte initial SN from the 128 bits secret key, instead of zero. This makes it more difficult to estimate the SN. Second, A-MAC can change the secret key and SN using the KSSN. This can prevent brute-force attacks. A malicious adversary can try to break the digested message with the probability of a 1 over 2^{32} . Assuming that an adversary repeatedly sends a forged frame, he could be accepted by AP after about 2^{31} . On 11Mbps 802.11b, an adversary can send about 27,000 forgery attempts per second (where $RTS = 15 \mu s$, $SIFS = 10 \mu s$) without considering transmission delay. Therefore, it would take over 44 hours. However, the A-MAC can change the secret key and SN frequently in a short time without a complex key exchange mechanism. In addition, A-MAC does not have any frame size overhead. These characteristics can overcome the weakness of Zhou's scheme [5] and Bellardo's scheme [7]. Therefore, the proposed A-MAC is useful for government, the military, and enterprises.

4.2 Performance Analysis

We analyze the performance and effectiveness of the proposed A-MAC. In the simulation environment, we do not consider complex key distribution systems (e.g. Remote Authentication Dial-In User Service: RADIUS), because the key distribution process does not affect our simulation, the secret keys are presumed to have been delivered to the AP and stations. The jammer and AP send RTS, CTS, de-authentication, and disassociation frames. When stations receive the frames, they process the A-MAC operation. If those frames are verified, they follow the instruction of each frame, otherwise they discard the frames.

Fig. 8 shows the simulation topology. There are one AP, four stations and one jammer. Especially, the jammer is hidden and near the AP. All stations can hear the RTS, CTS, de-authentication, and disassociation frame of the AP and Jammer. The jammer starts to attack at 6 seconds. Under transmission errors, although the station is a legitimate user, it cannot access the medium. In the real world, we can check the RSSI level and adjust the threshold. However, we just set the threshold value to one in this simulation. That is, we do not consider transmission error. Fig. 9 shows the throughput under the RTS/CTS Adversary model. The jammer transmits the faked RTS frame at the rate of 100, 200 Frames Per Second (FPS). As the faked RTS frame increases, the throughput lowers. Under a rate of 200 FPS, the throughput falls to 5.25 Mbps. However, we can sustain the throughput under the A-MAC. Fig. 10 shows the throughput under the NAV Adversary model. The Jammer sets the duration field to the length of 2^{11} , 2^{12} . This attack is more serious than the previous attack. Under a length of 2^{12} , the total throughput drastically falls to 50%, since the NAV at 2^{12} ($2^{12} \times 1 \mu s = 0.004 \text{ sec}$) has a considerable delay in 802.11 systems. However, we can also sustain the throughput under the A-MAC. Fig. 11 shows the throughput under the

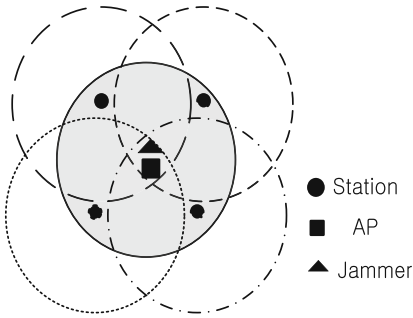


Fig. 8. Simulation Topology

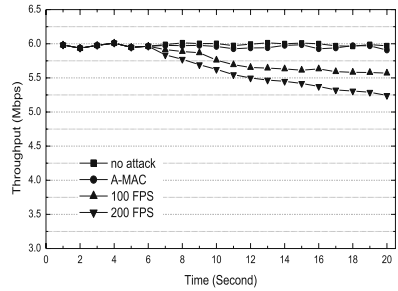


Fig. 9. RTS/CTS Adversary

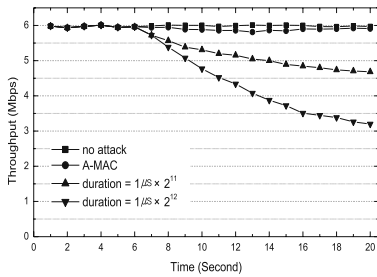


Fig. 10. NAV Adversary

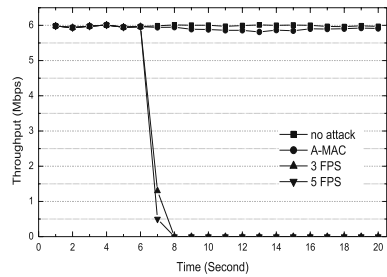


Fig. 11. De-auth / Dis-ass Adversary

De-authentication and Disassociation Adversary. The simulation results of the De-authentication Adversary model are similar to the results of the Disassociation Adversary model. This is due to the de-authentication frame-received station also terminating the connection during data transmission. So, we just show the results of the De-authentication Adversary model. The jammer transmits the faked frames at rates of 3, 5 FPS. These attacks are the most powerful and efficient attacks. Just a few frames can decrease the throughput to almost zero, since all stations undergo the state transition to State 1 or State 2, as soon as they get the frames. In the four cases, the A-MAC can prevent the attacks.

5 Conclusions

In this paper, we introduced the A-MAC to prevent VCS Attacks and De-authentication / Disassociation Attacks that are typical Protocol Jamming Attacks. A-MAC can authenticate frames using UMAC-32 and Hidden SN. A-MAC frequently changes the secret key and SN by using shift row and shift column operations to prevent the weakness of short 32 bits hashing codes. Therefore, A-MAC can achieve integrity, authentication, and anti-replay-attacks security features, and can prevent Protocol Jamming Attacks that degrade wireless network throughput. Our simulation

shows that A-MAC can sustain the network throughput under Protocol Jamming Attacks. In the near future, we will study a physical layer anti-jamming scheme in 802.11 based systems. If we mitigate Protocol Jamming Attacks and physical layer attacks in 802.11 based systems, security-sensitive groups (e.g. government, enterprises, and the military) have a preference to use WLAN. We hope that our research will help that situation come true.

References

1. <http://www.kismetwireless.net>
2. Liu, C., Yu, J.: Rogue Access Point Based DoS attacks against 802.11 WLANs. In: The Fourth AICT (2008)
3. Gast, M.S.: 802.11 Wireless Networks: The Definitive Guide, pp. 33–66. O'Reilly Publisher, Sebastopol (2002)
4. Malekzadeh, M.: Empirical Analysis of Virtual Carrier Sense Flooding Attacks over Wireless Local Area Network. *Journal of Computer Science* (2009)
5. Zhou, Y., Wu, D., Nettles, S.: Analyzing and Preventing MAC-Layer Denial of Service Attacks for Stock 802.11 Systems. In: Workshop on BWSA, Broadnets (2004)
6. Karlof, C., Sastry, N., Wagner, D.: TinySec: A link Layer Security Architecture for Wireless Sensor Networks. In: SenSys (2004)
7. Bellardo, J., Savage, S.: 802.11 Denial-of-Service Attacks: Real Vulnerabilities and Practical Solutions. In: USENIX Security Symposium (2003)
8. Krovetz, T.: RFC4418-UMAC: Message Authentication Code using Universal Hashing. In: IEEE Network Working Group (March 2006)

Femtocell Deployment to Minimize Performance Degradation in Mobile WiMAX Systems*

Chang Seup Kim, Bum-Gon Choi, Ju Yong Lee, Tae-Jin Lee,
Hyunseung Choo, and Min Young Chung**

School of Information and Communication Engineering
Sungkyunkwan University
300, Chunchun-dong, Jangan-gu, Suwon, Kyunggi-do, 440-746, Korea

Abstract. Femtocell is one of the promising technologies for improving service quality and data rate of indoor users. Femtocell is short-range, low-cost, and low-power base stations (BS) installed by the consumer indoors. Even though femtocell can provide improved home coverage and capacity for indoor users, it causes interference to macrocell users when femtocell uses the same frequency band of macrocell. To reduce the interference between macrocell and femtocell, it is needed to analyze the characteristic of macrocell and femtocell considering various interference scenarios. In this paper, we investigate the uplink and downlink capacity of macrocell and femtocell according to the strength of femtocell transmit power and the number of femtocells through simulations. Simulation results show that femtocell transmit power and the number of femtocells affect the performance of macrocell. From the results, we find the adequate femtocell transmit power which minimize the performance degradation of macrocell and femtocell. We also investigate capacities of macrocell and femtocell according to the locations of femtocell BS and macrocell UE.

Keywords: Femtocell, Macrocell, Co-Channel Interference, Power control.

1 Introduction

Mobile WiMAX (IEEE 802.16e) based on IEEE 802.16e standard has been commercialized in 2006 with rapid growth of various multimedia applications [1,2]. Recently, user demands for various multimedia services such as mobile IPTV are continuously increasing. On the other hand, it is expected that 60 percent of voice calls and 90 percent of data services will take place in indoors [3], and in current mobile communication systems, the signal strength transmitted from base station (BS) can be very low in indoor environments because the signal

* This work was supported by National Research Foundation of Korea Grant funded by the Korean Government (2009-0074466) and by the MKE (The Ministry of Knowledge Economy), Korea, under the ITRC (Information Technology Research Center) support program supervised by the NIPA (National IT Industry Promotion Agency) (NIPA-2009-(C1090-0902-0005) and NIPA-2009-(C1090-0902-0046)).

** Corresponding author.

strength may be severely attenuated when it penetrates the obstacles such as walls. Thus, providing high data rate services to indoor users is difficult only in mobile WiMAX systems. Therefore, the mobile WiMAX has performed femtocell standardization in two phases. The first phase of a WiMAX femtocell system, which requires no air interface or UE change, is expected to be available in the near future. More advanced and optimized femtocell features in phase 2 will be available upon completion of IEEE 802.16m and as part of WiMAX Release 2.0 with target deployments [4].

Femtocell is connected to IP based backhaul such as digital subscriber line (DSL) or cable modem, which provides lower cost than wireless mobile communication systems. Femtocells can achieve a high signal-to-interference-plus-noise ratio (SINR) with low transmit power because the distance between transmitter and receiver is short. In addition, the operators can save costs when femtocell accommodates traffic concentrated to macrocell of current mobile communication. On the other hand, users can receive improved data speeds and service quality with inexpensive costs. Although femtocell has many advantages in indoor environments, some technical issues, such as spectrum assignment of femtocell, access policy, network synchronization, handover, self-optimization, and self-configuration, should be solved to be effectively deployed in existing mobile WiMAX systems [5]. Especially, spectrum assignment for avoiding interference between macrocell and femtocell is an important issue since interference may cause severe degradation of throughput and service quality [6,7,8].

In mobile WiMAX system, throughput is determined by radius of macrocell, distribution of femtocell, and density of femtocell [9]. The interference characteristics under co-channel environments are presented in [10]. If femtocells are deployed in mobile WiMAX systems, interference from femtocells can affect the performance of mobile WiMAX systems. In this paper, we evaluate uplink and downlink capacity of macrocell and femtocell by simulation under various interference scenarios. We investigate the adequate femtocell transmit power level, which minimizes the performance degradation of mobile WiMAX system and can achieve enough femtocell throughput. Also, we examine the effect of the locations of femtocell BS and macrocell UE.

The rest of this paper is organized as follows. Section 2 provides interference scenarios in the environment where femtocells and macrocells coexist. Section 3 describes the system model and propagation models. In Section 4, we examine the uplink and downlink capacity of the macrocell and femtocell according to the strength of femtocell transmit power, number of femtocells, and locations of femtocell BS and macrocell UE. Finally, we conclude in Section 5.

2 Interference Scenarios of Femtocells in Mobile WiMAX Systems

When femtocells operate in the same spectrum in mobile WiMAX systems, the characteristic of interference depends on the method for channel assignment [11,12]. Channel assignment method for macrocell and femtocell is classified

into orthogonal and common channel assignment. Orthogonal channel assignment does not cause co-channel interference between the macrocell and femtocell because the femtocells use different channels with macrocell. However, orthogonal channel assignment has a low spectrum efficiency because frequency resources for the macrocell and femtocell are different. On the other hand, common channel assignment has a high spectrum efficiency because macrocell and femtocell can use all of the spectrum. However, there exists a problem of co-channel interference in common channel assignment method.

In common channel assignment, co-channel interference between macrocell and femtocell should be reduced to guarantee service quality of users. We analyze the co-channel interference between macrocell and femtocell according to link direction and the locations of femtocell BS and macrocell UE. Co-channel interference scenario between macrocell and femtocell is shown in Fig. 1. In case of uplink, interferences of macrocell and femtocell are caused by femtocell UE and macrocell UE, respectively. Similarly, in downlink, interferences of macrocell and femtocell are caused by femtocell BS and macrocell BS. If macro/femtocell UE is close to femto/macrocell BS, interference between macrocell and femtocell becomes large. Thus, the channel capacity of macro/femtocell may be significantly deteriorated.

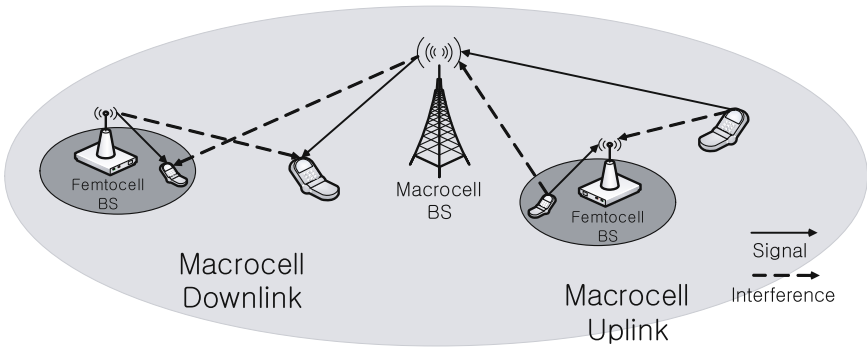


Fig. 1. Co-channel interference scenario between macrocell and femtocell

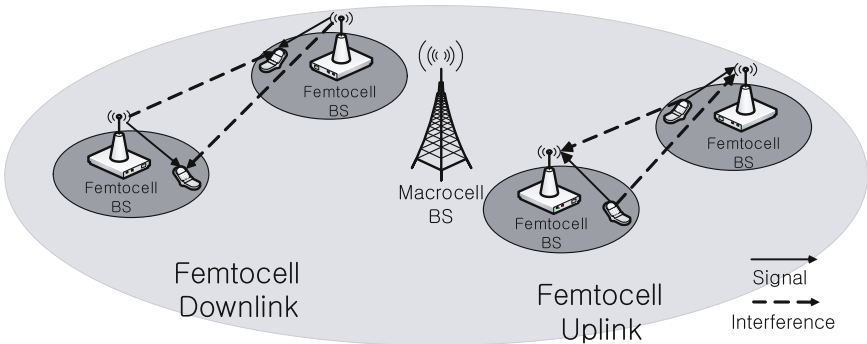


Fig. 2. Co-channel interference scenario between femtocells

When femtocells are deployed in a macrocell, we should also consider co-channel interference among femtocells as shown in Fig. 2. Interferences are caused by neighboring femtocell BS and femtocell UE. Especially, when the distance between femtocells is short, the impact of interference from neighboring femtocell may be high. That is, co-channel interference between femtocells can be a serious problem when a lot of femtocells are concentrated in a small area.

3 System Model

To analyze the performance when femtocells are deployed in mobile WiMAX systems, we consider a system model with 19 macrocells as shown in Fig. 3. In our model, femtocells are uniformly deployed in the centered macrocell 1. The attenuation of transmitted signal, called path loss, can differ according to channel environment (e.g., indoor or outdoor). As path loss models, we apply Wireless World Initiative New Radio (WINNER) II model to evaluate the throughput [13]. Simulation parameters are in Table 1 [14].

To calculate the path loss between macrocell BS and UE, we use non-line of sight (NLOS) outdoor propagation model. Also, we use NLOS indoor propagation model for calculating path loss between femtocell BS and UE. Path loss of the indoor and outdoor signal is given as [13]:

$$PL_{in}(d) = 36.8 \cdot \log_{10}(d) + 37.78,$$

$$PL_{out}(d) = 35.74 \cdot \log_{10}(d) + 35.68,$$

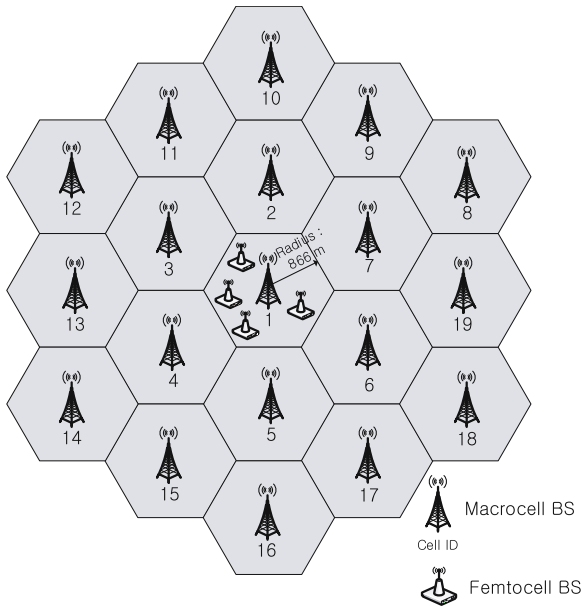


Fig. 3. Cell structure for simulation

Table 1. Simulation parameters

PARAMETER	VALUE	PARAMETER	VALUE
Carrier frequency	2.5 GHz	Channel bandwidth	10 MHz
Radius of macrocell	866 m	Radius of femtocell	10 m
Maximum power of macrocell BS	20 W	Maximum power of femtocell BS	0.1 W
Maximum power of UE	0.2 W	Antenna gain of macrocell BS	15 dBi
Antenna gain of femtocell BS	2 dBi	Antenna gain of UE	-1 dBi
Lognormal std. dev. (outdoor)	8 dB	Lognormal std. dev. (outdoor to indoor)	10 dB
Lognormal std. dev. (indoor)	4 dB	Lognormal std. dev. (indoor to outdoor)	7 dB

where d is the distance between the transmitter and the receiver in meter. In order to reflect the interference between macrocell and femtocell, indoor to outdoor and outdoor to indoor path loss model are considered as follows [13].

$$PL_{in_to_out}(d_{out}, d_{in}) = PL_{out}(d_{out} + d_{in}) + 14 + 15(1 - \cos(\theta))^2 + 0.5d_{in},$$

$$PL_{out_to_in}(d_{out}, d_{in}) = PL_{out}(d_{out} + d_{in}) + 13.8 + 0.5d_{in},$$

where d_{out} is the distance between the outdoor node and the wall nearest to the indoor node, d_{in} is the distance from the wall to the indoor node, and θ is the angle between the outdoor path and the vertical direction of the wall. In the simulation, θ is assumed to be 0 degree.

The received signal to interference ratio (SIR) is defined as the ratio of a signal power to the interference power. Let P_s^r and I_s^r be strengths of received signal and interference from sender s to receiver r , and M and F be the number of macrocells and femtocells.

SIR of user l in the downlink at macrocell BS i is:

$$SIR_{macroBS^i}^l = \frac{P_{macroBS^i}^l}{\sum_{j=1}^F I_{femtoBS^j}^l + \sum_{k=1, k \neq i}^M I_{macroBS^k}^l},$$

where $I_{femtoBS^j}^l$ is the interference which occurs due to the signal transmitted to the macrocell user l in the femtocell BS j and $I_{macroBS^k}^l$ is the interference which occurs due to the signal transmitted to the macrocell user l in the macrocell BS k . SIR of user l in the uplink at macrocell BS i is:

$$SIR_{macroUE^l}^i = \frac{P_{macroUE^l}^i}{\sum_{j=1}^F I_{femtoUE^j}^i + \sum_{k=1, k \neq l}^M I_{macroUE^k}^i},$$

where $I_{femtoUE^j}^i$ is the interference which occurs due to the signal scheduled to transmit from the femtocell UE j to the macrocell BS i , and $I_{macroUE^k}^i$ is the interference which occurs due to the signal scheduled to transmit from the

macrocell UE k to the macrocell BS i . SIR of user l in the downlink at femtocell BS i is:

$$SIR_{femtoBS^i}^l = \frac{P_{femtoBS^i}^l}{\sum_{j=1, j \neq i}^F I_{femtoBS^j}^l + \sum_{k=1}^M I_{macroBS^k}^l},$$

SIR of user l in the uplink at femtocell BS i is:

$$SIR_{femtoUE^l}^i = \frac{P_{femtoUE^l}^i}{\sum_{j=1, j \neq l}^F I_{femtoUE^j}^i + \sum_{k=1}^M I_{macroUE^k}^i}.$$

The received SIR is used to calculate channel capacity with the Shannon-Hartley theorem. Channel capacity C is calculated as follows.

$$C = BW \cdot \log_2(1 + SIR),$$

where BW is channel bandwidth in Hz.

4 Performance Evaluation

Fig. 4 shows uplink and downlink channel capacity of macrocell according to the number of femtocells. As the number of femtocell increases, the uplink channel capacity of macrocell steeply decreases, but downlink channel capacity of macrocell decreases more slowly than uplink channel capacity. The reason is that uplink transmit power of macrocell UE is lower than downlink transmit power of macrocell BS and strength of uplink interference is similar to strength of downlink interference. In order to maintain uplink and downlink channel capacity of macrocell high, the maximum transmit power of femtocell BS and UE is controlled. We set the transmitted power of femtocell BS and UE to 0.1 W, 0.01 W, and 0.001 W in the simulation. When the transmitted power of femtocell UE uses 0.001 W, the decreasing slope of uplink channel capacity of macrocell is reduced because co-channel interference from femtocell UE is decreased. On the other hand, macrocell UE in the downlink can get higher SIR than uplink of macrocell due to a high signal strength of macrocell BS. However, the downlink channel capacity of macrocell has smaller improvement than uplink of macrocell even if the transmitted power of femtocell BS uses 0.001 W.

Compared with the channel capacity of macrocell as shown in Fig. 4, uplink and downlink channel capacity of femtocell is high as shown in Fig. 5. Femtocell has a high received SIR because the distance between femtocell BS and UE is shorter than the distance between macrocell BS and UE. Since transmitted power of macrocell UE is low compared with macrocell BS, co-channel interference from macrocell UE is relatively low. Although the transmitted power of femtocell's UE is weakened, the uplink channel capacity of femtocell maintains high. Accordingly, we can increase the uplink channel capacity of macrocell. The downlink channel capacity of femtocell is constant because interference

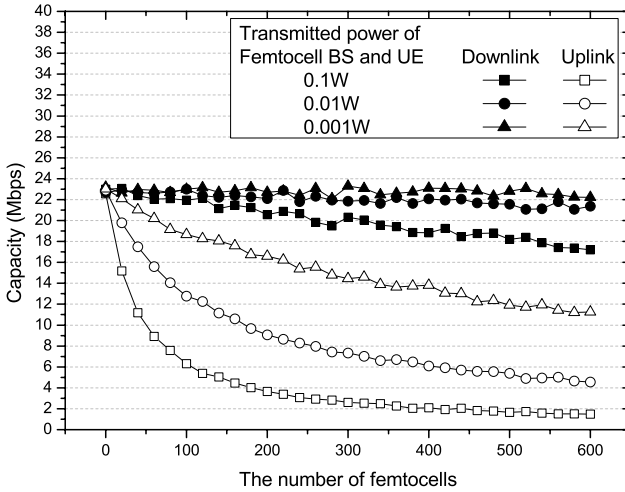


Fig. 4. Uplink and downlink channel capacity of macrocell

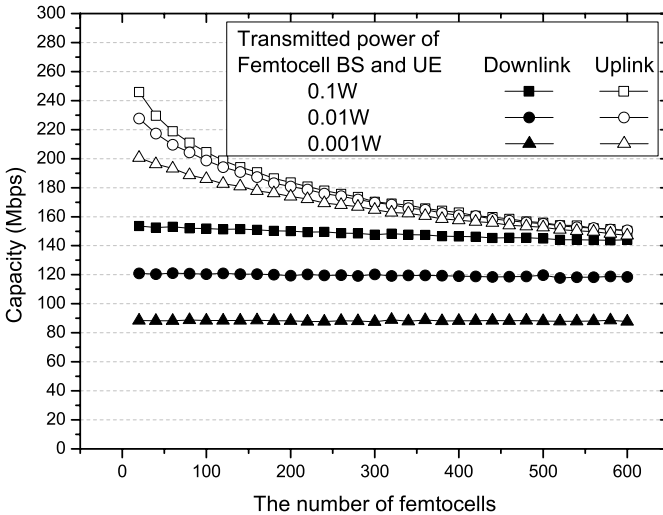
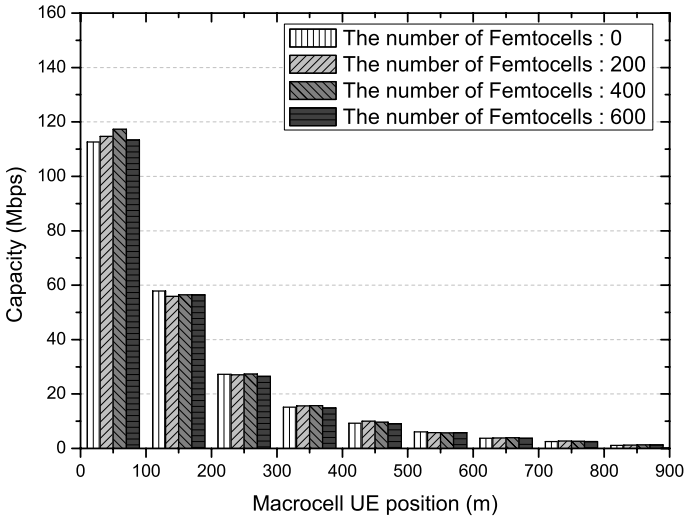


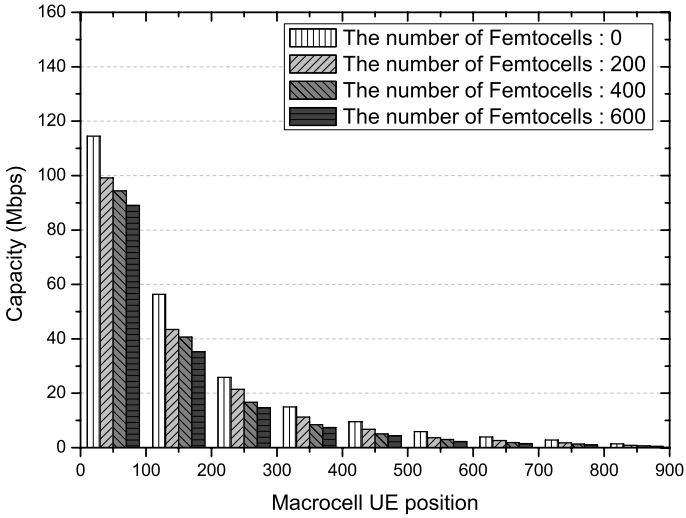
Fig. 5. Uplink and downlink channel capacity of femtocell

from macrocell BS is more dominant than interference of neighboring femtocells. When the transmitted power of femtocell BS uses 0.001 W, the downlink channel capacity of femtocell is about 90 Mbps. In this setting, femtocell can provide subscriber with enough service quality and data speed while minimizing the performance degradation of macrocell.

We consider the environment where transmitted power of femtocell BS and UE is fixed to 0.001 W. Fig. 6 shows downlink and uplink channel capacity

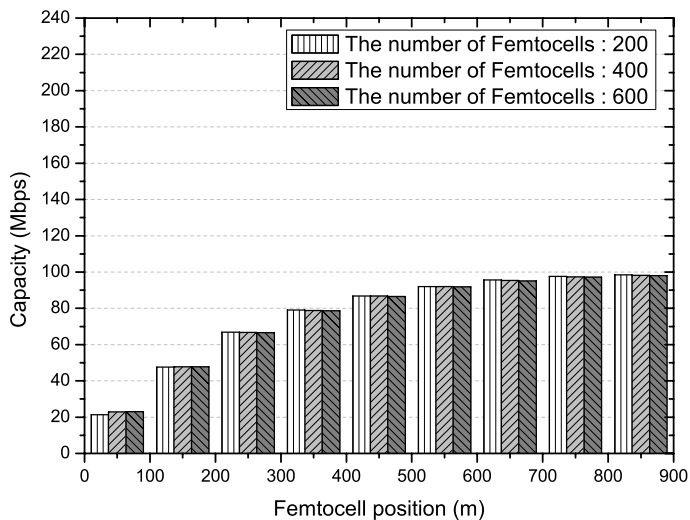


(a) Downlink

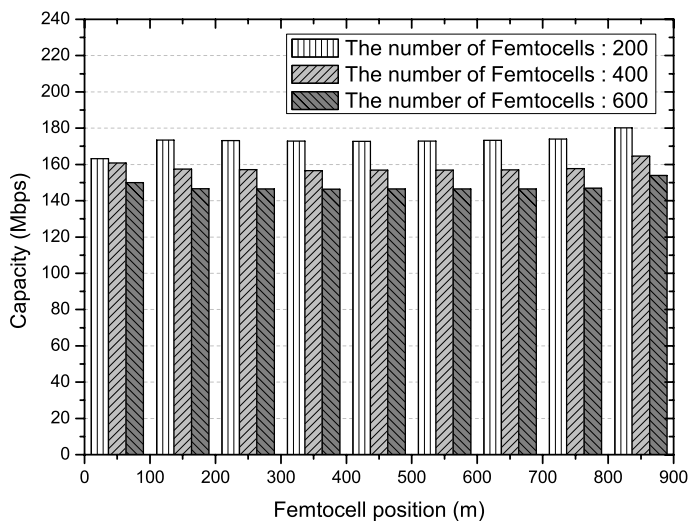


(b) Uplink

Fig. 6. Channel capacity of macrocell UE according to locations (The transmitted power of femtocell BS and UE is 0.001 W)



(a) Downlink



(b) Uplink

Fig. 7. Channel capacity of femtocell UE according to locations (The transmitted power of femtocell BS and UE is 0.001 W)

of macrocell according to the locations of macrocell UE. Macrocell UE has a high channel capacity when the location of macrocell UE is close to macrocell BS. The downlink channel capacity of macrocell is not influenced by co-channel interference of femtocell because the transmitted power of macrocell BS is high and the transmitted power of femtocell BS is low as shown in Fig. 6(a). On the other hand, co-channel interference according to the number of femtocell UEs highly affects the uplink channel capacity of macrocell UE since the signal strength of macrocell UE is low as shown in Fig. 6(b).

Fig. 7(a) shows the downlink channel capacity of femtocell according to the locations of femtocell. Femtocells close to macrocell BS has a low channel capacity because the transmitted power of macrocell BS becomes a strong interference to femtocell. If the location of femtocell is far away from macrocell BS, femtocell has a high channel capacity. Fig. 7(b) shows uplink channel capacity of femtocell according to the locations of femtocell. The uplink channel capacity of femtocell does not depend on the distance between macrocell BS and femtocell. However, the uplink channel capacity of femtocell decreases as the number of femtocells increases.

5 Conclusion

In this paper, we investigated channel capacity of macrocell and femtocell according to the number and power level of femtocells when femtocells are deployed in mobile WiMAX systems. When femtocell uses high transmit power, the channel capacity of macrocell severely decreases due to high interference from femtocell as the number of femtocells increases. In our simulation environment, the adequate transmit power of femtocell BS and UE is 0.001 W. This power level minimizes the performance degradation of macrocell while achieving enough femtocell throughput. With this femtocell transmit power, we also examined the effect on the locations of femtocell BS and macrocell UE.

References

1. IEEE 802.16e-2005: Local and Metropolitan Networks-Part 16: Air Interface for Fixed and Mobile Broadband Wireless Access System, Amendment 2: Physical and Medium Access Control Layers for Combined Fixed and Mobile Operation in Licensed Bands and Corrigendum 1 (2006)
2. WiMAX Forum: Mobile WiMAX-Part I: A Technical Overview and Performance Evaluation (2006)
3. Gordon Mansfield: Femto Cells in the Us Market-Business Drivers and Femto Cells in Us Market-Business Drivers and Consumer Propositions. FemtoCells Europe (2008)
4. Kim, R.Y., Kwak, S.K., Etemad, K.: WiMAX Femtocell: Requirements, Challenges, and Solutions. *IEEE Commun. Magazine* 47(9), 87–91 (2009)
5. Chandrasekhar, V., Andrews, J.G., Gatherer, A.: Femtocell Networks: A Survey. *IEEE Commun. Magazine* 46(9), 59–67 (2008)

6. 3GPP R4-071661, Ericsson: Impact of HNB with controlled output power on macro HSDPA capacity (2007)
7. 3GPP R4-080409, Qualcomm Europe: Simple Models for Home NodeB Interference Analysis (2008)
8. Claussen, H.: Performance of macro- and co-channel femtocells in a hierarchical cell structure. In: Proc. IEEE 18th International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC 2007), pp. 3–7 (2007)
9. Yeh, S.-p., Talwa, S., Lee, S.-C., Kim, H.: WiMAX femtocells: A Perspective on Network Architecture Capacity, and Coverage. *IEEE Commun. Magazine* 46(10), 58–64 (2008)
10. Sung, Y.-S., Jeon, N.-R., Woon, B.-W., Lee, J.-S., Lee, S.-C., Kim, S.-C.: Femtocell/Macrocell Interference Analysis for Mobile WiMAX System. In: IEEE VTS Asia Pacific Wireless Communications Symposium (2008)
11. López-Pérez, D., Valcarce, A., Roche, G.D.L., Zhang, J.: OFDMA Femtocells: A Roadmap on Interference Avoidance. *IEEE Commun. Magazine* 47(9), 41–48 (2009)
12. 3GPP TR25.820: 3G Home NodeB Study Item Technical Report. v8.2.0 (2008)
13. WINNER II WP1: WINNER II Part 1 Channel Models. IST-4-027756, D1.1.2, V1.1 (2007)
14. WiMAX Forum: WiMAX System Evaluation Methodology. v.2.01 (2007)

A Novel Frequency Planning for Femtocells in OFDMA-Based Cellular Networks Using Fractional Frequency Reuse*

Chang-Yeong Oh, Min Young Chung, Hyunseung Choo, and Tae-Jin Lee**

School of Information and Communication Engineering
Sungkyunkwan University
440-746, Suwon, South Korea, +82-31-290-7982
{ohchy, mychung, choo, tjlee}@ece.skku.ac.kr

Abstract. Femtocells are expected to be one of the emerging technologies for next generation communication systems. For successful deployment of femtocells in the pre-existing macrocell networks, there are some challenges such as the cell planning for interference management, handoff, and power control. In this paper, we focus on frequency planning which can provide interference avoidance for the co-existence of macrocells and femtocells. We propose a novel frequency planning for femtocells in cellular networks using fractional frequency reuse (FFR). We consider downlink performance of cellular systems based on Orthogonal Frequency Division Multiplexing Access (OFDMA), e.g., WiMAX and 3GPP Long Term Evaluation (LTE). Simulation results show that our scheme indeed reduces the effect of additional co-channel interference (CCI) between a given macrocell and deployed femtocells as well as neighboring macrocells.

1 Introduction

The concept of indoor cellular networks using home base stations (BSs), femtocells, with low transmission power is drawing much interest. Femtocells are one of the emerging cost-effective technologies for both operators and users to enhance coverage and to support higher data rate by bridging existing handsets and broadband wired networks. However, deployment of femtocells in pre-existing cellular networks causes some problems to be addressed. One of them is that it requires intelligent frequency allocation for femtocells and traditional macrocells when they operate simultaneously in the same network [3]. So it is important to efficiently allocate frequency resources for femtocells considering the effect of co-channel interference. Since frequency resources are usually limited, some methods to enhance the efficiency are desirable.

Traditionally, there exist frequency allocation mechanisms which allocate frequency resources according to frequency reuse factors (RFs) in multi-cell environment. One is

* This research was supported by the MKE(The Ministry of Knowledge Economy), Korea, under the ITRC(Information Technology Research Center) support program supervised by the NIPA(National IT Industry Promotion Agency) (NIPA-2009-C1090-0902-0005) and (NIPA-2009-C1090-0902-0046).

** Corresponding author.

the shared frequency allocation mechanism, in which frequency is shared with adjacent cells, i.e., universal frequency reuse or frequency RF of 1. The other is the orthogonal frequency allocation mechanism which allocates frequency according to the predetermined frequency patterns among multi cells with the frequency RF of 1 or more. The former may be viewed as more efficient because all frequency resource can be utilized in each of cells. However, it can suffer from performance degradation since the performance of users located at the edge of a cell is degraded due to co-channel interference (CCI) from neighboring cells. So the frequency RF of 3 or above is generally employed to mitigate CCI in 2G systems. The reduction of CCI is typically gained at the cost of efficiency in frequency resource allocation [6].

Fractional frequency reuse (FFR), which is the mixture of the shared and orthogonal frequency allocation, becomes one of the solutions to improve frequency efficiency in OFDMA-based cellular systems, e.g., WiMAX and 3GPP Long Term Evaluation (LTE). In most of FFR schemes, each cell is partitioned into two regions: *inner region and outer region*. The inner region around a base station (BS) in a cell can use the frequency RF of 1 due to low CCI from adjacent cells. The outer region far from a BS in a cell uses different frequency bands with those in adjacent cells with the frequency RF of 3 to reduce CCI [11].

FFR can be implemented easily in OFDMA-based cellular systems because the frequency band of a cell is divided into subchannels and can be handled by the unit of subchannel, a group of orthogonal subcarriers. Han et al. propose an FFR scheme to enhance the flexibility of frequency assignment and cell performance, in which mobile stations (MSs) in the inner region can use the subchannels assigned to both the inner and the outer region of a cell according to the frequency partitioning [4]. Another FFR scheme is introduced by noting that the inner region and the outer region are differently served by not only frequency bands but also time slots [5].

While the performance of the users located at the edge of a cell is important in the macrocell networks with FFR, interference management between macrocells and femtocells is a main issue in femto / macro co-deployment environment with different channel operation strategies: orthogonal channel allocation and shared channel allocation [1]. So the frequency assignment for femtocells is one of the key issues [2]. The authors in [2] propose that femtocells in the inner region of a macrocell use different subchannels with those for macrocell users to minimize interference, and femtocells in the outer region of a macrocell use the same subchannels as those for macrocell users. They suppose cell environment without FFR, in which the frequency RF of macrocells is greater than 1.

In this paper, we propose a frequency planning mechanism with FFR for macrocells and femtocells, in which co-channel operation is allowed in low CCI, and orthogonal frequency allocation is adopted in high CCI. To reduce the effect of CCI we utilize resource allocation in both frequency bands and time slots.

The remaining part of this paper is organized as follows. Section 2 describes the basic FFR scheme for macrocells and presents the proposed frequency planning

mechanism for femtocells and macrocells. Performance evaluation of the proposed scheme is provided in Section 3. We conclude in Section 4.

2 Proposed Frequency Planning

2.1 FFR for Macrocells

In our proposed FFR for macrocells, a macrocell is grouped into two parts: *inner region* and *outer region*. The decision on whether an MS belongs to inner or outer region is made by the reported signal-to-interference-plus-noise-ratio (SINR) of the reference signal at the initial configuration or at the periodic state update. If sufficient SINR is observed, i.e., the reported SINR of an MS is above the pre-determined SINR threshold, the BS considers that the MS belongs to the inner region. Otherwise, i.e., the reported SINR of an MS is below the pre-determined SINR threshold, the BS considers that the MS is in the outer region. The reported SINR of an MS is largely dependent on its location, i.e., the farther an MS is from the BS in a cell, the larger the path-loss of the signal from the BS becomes.

Fig. 1 shows an example, in which MSs in the inner region use the RF of 1 and those in the outer region use the RF of 3. It illustrates the feasible frequency bands for macrocells according to the region that an MS belongs to. So the whole subchannels can be allocated to the inner region while the outer region uses just 1/3 of all subchannels in a cell. Furthermore, the service times of the inner and the outer region are separated by time slots. The MSs in the inner region of macrocells can be served simultaneously because CCI from neighbor macro BSs is limited. They can utilize the whole subchannels

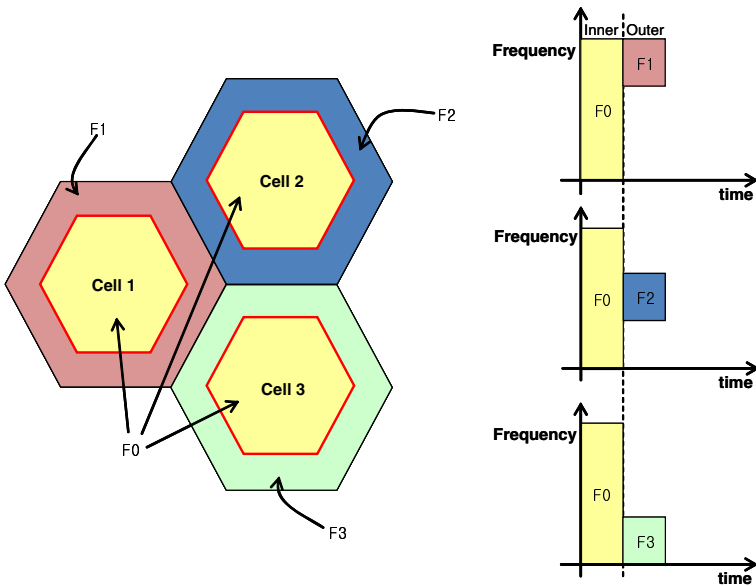


Fig. 1. Proposed frequency band/time slots allocation with FFR for macrocells

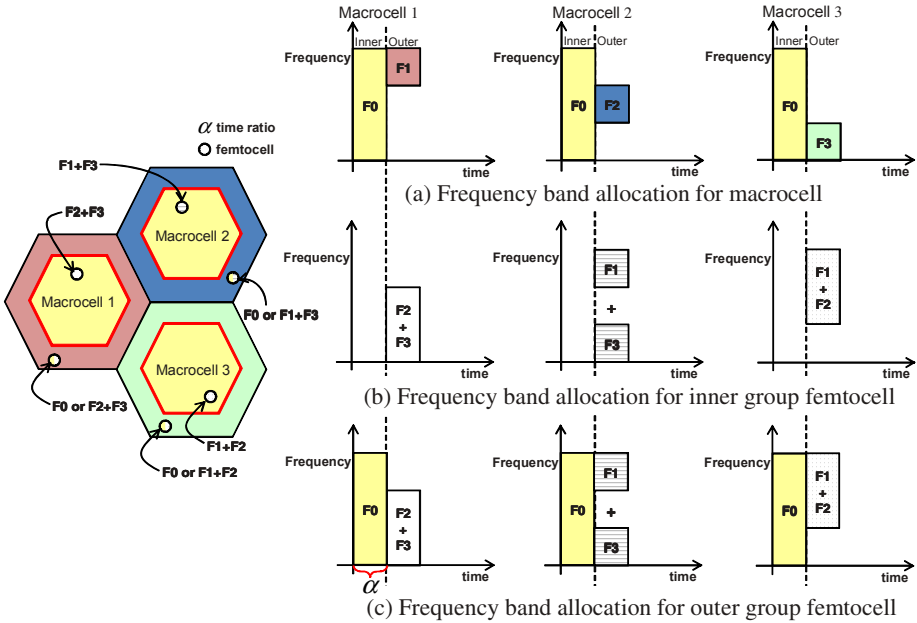


Fig. 2. Proposed frequency planning for femtocells

during their time slots. The MSs in the outer region of macrocells can also be served simultaneously but they must use the orthogonal subchannels in order to avoid high CCI from neighbor macro BSs. In summary, the inner region and the outer region are allocated shared or separate time slots and frequency bands to mitigate intra- and inter-cell interference among the inner and the outer regions as well as among cells.

2.2 Frequency Planning for Femtocells

In this section, we propose a novel frequency planning for femtocells under the macrocells with FFR. First of all, femtocells should be deployed harmoniously in order not to make minor impacts on the legacy macrocells but to maximize their performance. In our proposal, femtocells are divided into the inner group and the outer group depending on whether a femtocell is located in the inner region or the outer region of a macrocell. A femto BS is assumed to support auto-configuration because it has to be able to setup a user controlled hot-spot. A femto BS reports the SINR of the reference signal of a macrocell to the macro BS by pilot sensing, which is similar to the macro MSs when they do in the setup procedure. Then the macro BS determines whether the femtocell belongs to the inner or the outer group and notifies the appropriate frequency bands to the femto BS. We assume that all macro and femto BSs are synchronized.

We now consider three factors for the frequency planning of femtocells: frequency, time, and spatial state. The femtocells in the inner group, i.e., *inner femtocells*, are served during the service time of the macro MSs in the outer region, i.e., *outer service time*, and the subchannels orthogonal to the macrocell frequency bands are allocated

in order to avoid intra CCI from the macrocell that the femtocell belongs to. For the macro MSs in the outer region, only 1/3 of the overall frequency band is utilized to reduce CCI to adjacent macrocells. So these 2/3 unused subchannels can be allocated to the inner femtocells efficiently (see Fig. 2). The CCI between the inner femtocells and the neighboring macrocells is limited due to low transmission power of femtocells and relatively large path loss factors, e.g., high attenuation path-loss by long distance and wall-loss by the walls in buildings.

For the femtocells in the outer group, i.e., *outer femtocells*, we note that they use low transmission power and the CCI between the macrocell and the femtocells is small. So they are allowed to use the same frequency bands and time slots, i.e., *inner service time* as those of the macro MSs in the inner region. That is, the outer femtocells can utilize all the subchannels in the given time slots. As shown in Fig. 2(c), outer femtocells use the full frequency band during the inner service time and 2/3 of the band during the outer service time.

3 Performance Evaluation

We consider 2-tier cellular networks consisting of $M(=19)$ macrocells. The macro BSs are located at the center of each cell. Femto BSs which have uniform separation with one another are densely deployed within a macrocell (see Fig. 3). We use the modified COST-Walfish-Ikegami (WI) urban micro model for the non-line-of-sight (NLOS) outdoor path-loss model [10].

$$P_{loss}^{out}[dB] = 31.81 + 40.5 \log_{10}(d[m]) + \chi_{\sigma^{out}}, \tag{1}$$

where d is the distance between a sender and a receiver and $\chi_{\sigma^{out}}$ represents the outdoor shadowing (log-normal fading), which is characterized by the Gaussian distribution

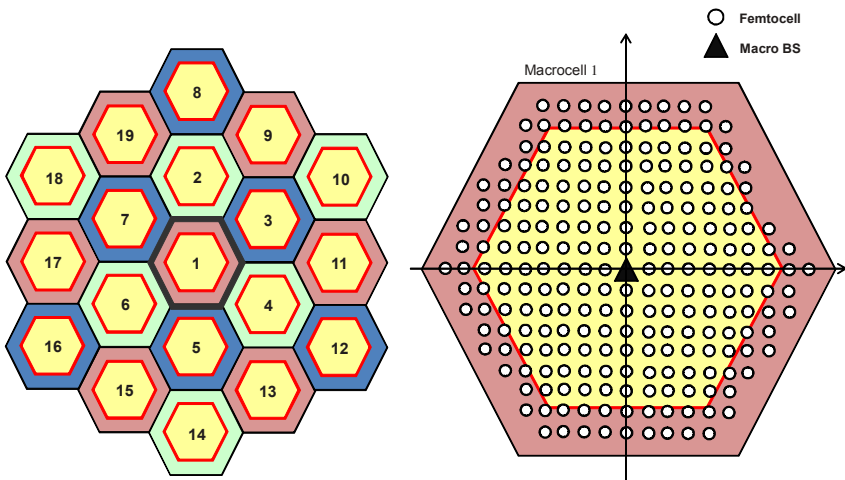


Fig. 3. 2-tier cellular networks and 234 uniformly separated femtocells in a macrocell

Table 1. Modulation and Coding Scheme (MCS) table

SINR	Code rate	Mod.	SINR	Code rate	Mod.
-4.34	1/12	QPSK	6.35	2/3	QPSK
-2.80	1/8	QPSK	9.50	1/2	16QAM
-1.65	1/6	QPSK	12.21	2/3	16QAM
0.31	1/4	QPSK	13.32	1/2	64QAM
1.51	1/3	QPSK	16.79	2/3	64QAM
4.12	1/2	QPSK	20.68	5/6	64QAM

Table 2. Simulation Parameters

Parameter	Value
Inter macro cell distance (ICD)	1000 m
Radius of a femtocell	20 m
FFT size	1024
Total number of data subcarriers in one cell	768
Down-link symbol rate	9.76 k symbols/sec
Total frequency bandwidth in one cell	10 MHz
AWGN power density (N_0)	-174 dBm/Hz
Macro BS power	20 W
Femto BS power	20 mW
Outdoor Log-Normal fading (σ^{out})	10 dB
Indoor Log-Normal fading (σ^{in})	4 dB

with zero mean and standard deviation (σ^{out}). And the modified COST 231-multi wall (MW) model one floor building [7] is employed for indoor propagation,

$$P_{loss}^{in} [dB] = 37 + 3.2 \cdot 10 \log_{10}(d[m]) + \sum_{t=1}^T L_w^t n_w^t + \chi_{\sigma^{in}}, \quad (2)$$

where d , T , L_w^t , n_w^t , and $\chi_{\sigma^{in}}$ are the distance between a sender and a receiver, the number of wall types, the wall loss according to wall type t , and the number of penetrated type t walls, and the shadowing, which is characterized by the Gaussian distribution with zero mean and standard deviation (σ^{in}), respectively. We assume each of the buildings with femto BSs has 1 outer (heavy) wall at 10m from the femto BS, $L_w^1=15$ dB, and 1 inner (light) wall at 5m from the femto BS, $L_w^2=3$ dB. We have K data subcarriers and F_i femtocells in macrocell i .

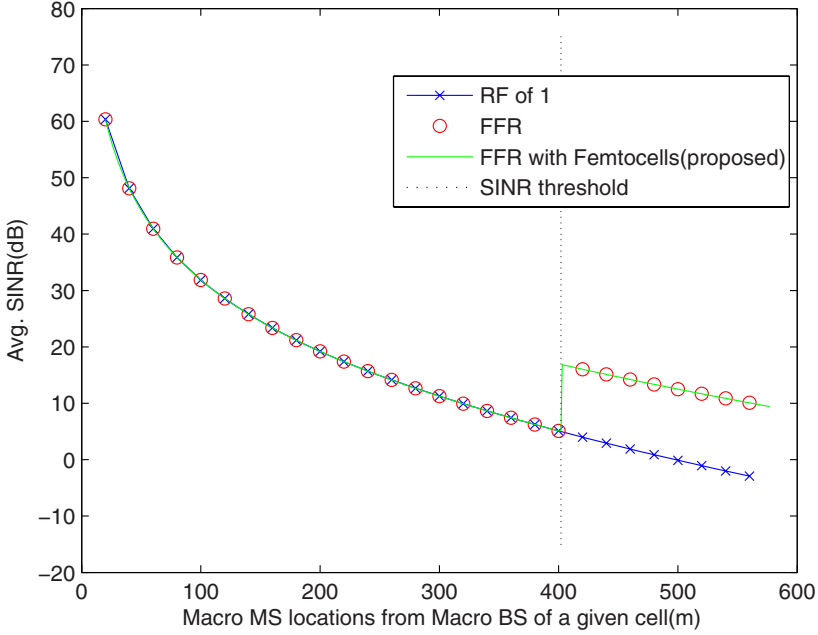


Fig. 4. Average (Avg.) SINR of a macro MS according to the distance from the macro BS in a given cell

The downlink SINR of user l in macrocell $i \in \{1, \dots, M\}$ on subcarrier $k \in \{1, \dots, K\}$ is

$$SINR_{l,k}^{(i)} = \frac{P_{l,k}^{(i)}}{\sum_{\eta=1, \eta \neq i}^M c_k^\eta \cdot I_k^\eta + \sum_{j=1}^{F_i} c_k^j \cdot I_k^j + N_0 \Delta f}, \quad (3)$$

where $P_{l,k}^{(i)}$, I_k^η , I_k^j , and $N_0 \Delta f$ are the received power of user l on subcarrier k in macrocell i , CCI from macrocell $\eta \in \{1, \dots, M\}$, CCI from femtocell $j \in \{1, \dots, F_i\}$ in macrocell i , and Additive White Gaussian Noise (AWGN) on a subcarrier, respectively. And $c_k \in \{0, 1\}$ is defined as the collision coefficient having the value 1 if subcarrier k is allocated to a specific cell so that it generates CCI to a reference cell, i.e., macrocell i , and 0 otherwise. The downlink SINR of user m on subcarrier k in femtocell j of macrocell i can be found as

$$SINR_{m,k}^{(i,j)} = \frac{P_{m,k}^{(i,j)}}{\sum_{\eta=1}^M c_k^\eta \cdot I_k^\eta + \sum_{\nu=1, \nu \neq j}^{F_i} c_k^\nu \cdot I_k^\nu + N_0 \Delta f}, \quad (4)$$

where $P_{m,k}^{(i,j)}$ and I_k^ν are the received power of user m on subcarrier k in femtocell j of macrocell i , CCI from femtocell $\nu \in \{1, \dots, F_i\}$ in macrocell i , respectively. We adopt

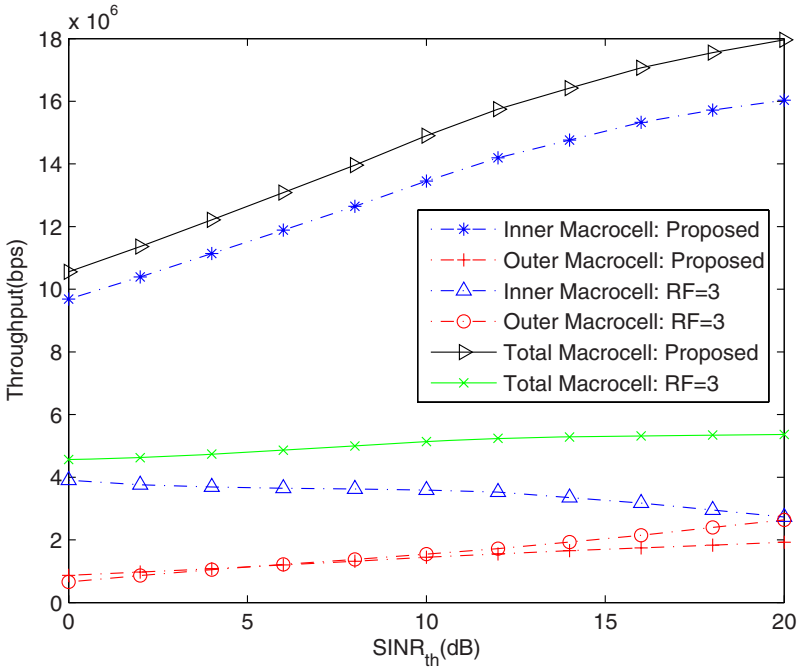


Fig. 5. Downlink throughput of a macrocell according to different SINR thresholds

Modulation and Coding Scheme (MCS) in Table 1 [8] and use simulation parameters in Table 2. The simplified proportional fairness allocation algorithm with the same power for subcarriers is considered in [9] to allocate subcarriers.

We present the average (Avg.) SINR variation of a macro MS in Fig. 4 as it moves on the x-axis from the center to the border when the macro BS is positioned at $(x, y)=(0,0)$. We consider the case where all macrocells use the same frequency band without femtocells, i.e., frequency RF of 1. In this case, we observe the Avg. SINR of a macro MS at the cell edge is relatively low as we expect. On the other hand, the basic FFR scheme without femtocells, shows that the Avg. SINR of a macro MS at the cell edge is greater than that of the RF of 1 due to CCI avoidance from tier-1 cells. We set the SINR threshold 5dB as the criterion to distinguish the inner and the outer region. We can observe the Avg. SINR of a macro MS with 234 femtocells is similar to that without femtocells. In the inner region there does not exist any effect of the deployment of femtocells at all because the inner femto BSs and the macro BS operate at different time slots. The outer femto BSs with co-channel operation do not interfere with the macro MSs much due to their low transmission power, high attenuation by long distance, and wall-losses. In the outer region, there is no additional CCI because all femto BSs and the macro BS are served at different frequency bands. We can conclude the performance of macro MSs with the numerous co-existing femtocells is not sacrificed in our proposed mechanism.

In Fig. 5, we evaluate the downlink throughput of a macrocell according to different SINR thresholds when there exist 36 active macro MSs in the cell. We compare ours

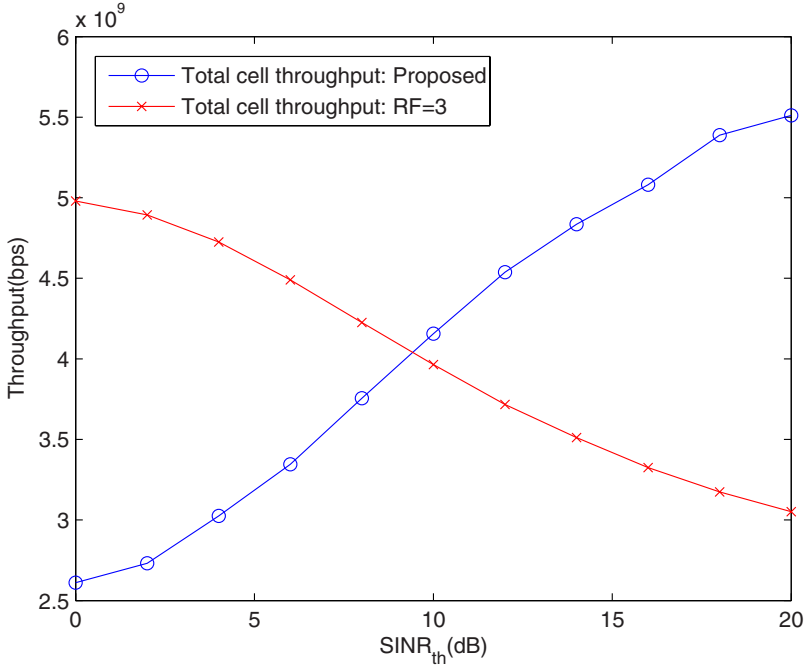


Fig. 6. Total throughput in a given cell area

with the scheme in [2] assuming that the RF of 3 is adopted. Note that the total bandwidth of FFR is much wider than that of the RF of 3 since the inner region adopt the RF of 1 and can use the full bandwidth for the macro MSs. So the throughput performance of macrocell by using FFR is highly enhanced. In our proposal, as the SINR threshold increases, Avg. SINR of the macro MSs in the inner region becomes higher. Thus the spectral efficiency is improved. Moreover, the throughput performance of the macro MSs in the outer region also increases thanks to the alleviation of CCI. However, if the SINR threshold is set to a high value, the frequency resource may be insufficient because many macro MSs may be included in the outer region but the bandwidth for the outer region is limited.

The scheme in [2] serves femtocells and macrocell simultaneously, especially, femtocells in the outer region and the macrocell share the frequency bands but the femtocells in the inner region and the macrocell use different frequency bands. If a macro MS is in an outer femtocell, it may experience high interference from the femto BS due to the co-channel operation. As the SINR threshold increases, the number of outer femtocells which may interfere with the macro MSs in the outer region also increases. This degradation is observed in the inner macrocell: RF=3 in Fig. 5. However, this interference scenario can be avoided in our proposed mechanism because macro MSs and outer femtocells are served by orthogonal frequency bands and time slots.

Fig. 6 shows the total cell throughput for downlink, i.e., the sum of throughput of all macro MSs and femto MSs in a given cell. Total cell throughput depends on the

throughput of femto MSs because femtocells are densely deployed and the spectral efficiency of femto MSs are much higher than that of macro MSs. We observe a trade-off in Fig. 5 as the SINR threshold varies. In our proposed mechanism, a femtocell in the outer region can use wider frequency band and longer time slots than a femtocell in the inner region while in the scheme in [2] with the RF of 3 a femtocell in the inner region use wider frequency bands than a femtocell in the outer region. The distribution of inner and outer femtocells is determined by the SINR threshold and the total cell throughput is shown to vary by the degree of distribution of femtocells.

4 Conclusion

In this paper, we have proposed a novel frequency planning for femtocells in OFDMA-based cellular networks using FFR. FFR scheme is an efficient method not only to enhance performance of edge users in a given cell but also to increase the spectral efficiency by adopting the mixture of reuse factors according to different region in a cell. Moreover, the basic rule of FFR, i.e., co-channel operation in low CCI or orthogonal channel operation in high CCI, is considered as a useful policy for frequency planning of femtocells when they are deployed in pre-existing macrocells. There may exist some difficulties and risks in sharing frequency among macrocells and femtocells, which cause performance degradation of macro MSs resulting from the effect of the deployment of femtocells. So we utilize unused frequency band in a given macrocell to inner femtocells for specific time slots in order to harmoniously serve macrocells and femtocells by maximizing diversity gains.

References

1. López-Pérez, D., Valcarce, A., Roche, G., Zhang, J.: OFDMA Femtocells: A Roadmap on Interference Avoidance. *IEEE Communications Magazine* 47(9), 41–48 (2009)
2. Guvenc, I., Jeong, M.-R., Watanabe, F., Inamura, H.: A hybrid frequency assignment for femtocells and coverage area analysis for co-channel operation. *IEEE Communications Letters* 12(12), 880–882 (2008)
3. Chandrasekhar, V., Andrews, J.G.: Femtocell Networks: A Survey. *IEEE Wireless Communications* 15(3), 1284–1536 (2008)
4. Han, S.S., Park, J., Lee, T.-J., Ahn, H.G., Jang, K.: A New Frequency Partitioning and Allocation of Subcarriers for Fractional Frequency Reuse in Mobile Communication Systems. *IEICE Transactions on Communications* E91-B(8), 2748–2751 (2008)
5. Giuliano, R., Monti, C., Loret, P.: WiMAX fractional frequency reuse for rural environments. *IEEE Communications Magazine* 46(9), 59–67 (2008)
6. Elayoubi, S.-E., Ben Haddada, O., Fourestie, B.: Performance evaluation of frequency planning schemes in OFDMA-based networks. *IEEE Transactions on Wireless Communications* 7(5), 1623–1633 (2008)
7. Nihtilä, T.: Increasing Femto Cell Throughput with HSDPA Using Higher Order Modulation. In: *Proc. INCC 2008, Lahore, Pakistan*, pp. 49–53 (2008)
8. Samsung Electronics, Experimental WiBro MCS Table (2007)

9. Nguyen, T.-D., Han, Y.: A Proportional Fair Algorithm with QoS Provision in Downlink OFDMA Systems. *IEEE Communications Letters* 10(11), 760–762 (2006)
10. Baum, D.S., Hansen, J., Salo, J.: An interim channel model for beyond-3G systems: extending the 3gpp spatial channel model (scm). In: *Proc. VTC 2005-Spring*, Stockholm, Sweden, pp. 3132–3136 (2005)
11. Huawei: Soft frequency reuse scheme for UTRAN LTE, 3GPP, R1-050507 (2005)
12. Damosso, E. (ed.): *Digital Mobile Radio Towards Future Generation Systems*, COST 231 final report (1998)

Association Analysis of Location Tracking Data for Various Telematics Services^{*}

In-Hye Shin and Gyung-Leen Park^{**}

Dept. of Computer Science and Statistics, Jeju National University
690-756, Jeju Do, Republic of Korea
{ihshin76,glpark}@jejunu.ac.kr

Abstract. This paper proposes an approach that extracts the association information from the location data obtained from the real fields but ignored so far. We provide and apply the approach to the real-life location tracking data collected from the Taxi Telematics system developed in Jeju, Korea. The analysis aims at obtaining taxis' meaningful moving patterns for the efficient operations of them. The proposed approach provides the flow chart which would not only take a glance around the overall analysis process but also help save temporal and economic costs required to employ the same or similar data mining analysis to similar services such as public transportations, distribution industries, and so on. Especially, we perform an association analysis on both of refined data and interesting factors extracted from the elementary analysis. The paper proposes the refined association rule mining process as follow: 1) obtaining the integrated dataset through the data cleaning process, 2) extracting the interesting factors from the integrated dataset using frequency and clustering method, 3) performing the association analysis, 4) extracting the meaningful and value-added information such as moving pattern, or 5) returning the feedback to adjust inappropriate factors. The result of the analysis shows that the association analysis makes it possible to detect the hidden moving patterns of vehicles that will greatly improve the quality of Telematics services considering the business requirements.

1 Introduction

The Taxi Telematics system has operated an efficient taxi dispatch service since 2006 in Jeju island, Korea [1]. Each taxi equipped with a GPS receiver reports its location to the central server every minute. The server is responsible for keeping such reports from each taxi and handles the call request from a customer by finding and dispatching a taxi closest to the call point. The tracking data obtained from the Taxi Telematics system can be utilized for various analysis on both research and business areas [2]. We have been doing our research to develop a data processing framework

^{*} This research was supported by the MKE (The Ministry of Knowledge Economy), Korea, under the ITRC (Information Technology Research Center) support program supervised by the IITA (Institute for Information Technology Advancement). (IITA-2009-C1090-0902-0040).

^{**} Corresponding author.

for an efficient analysis [3-9]. It is the exact fact that the empty taxi ratio can be reduced by means of guiding a taxi to the spot where many passengers are waiting for. In the project, the empty taxi ratio is around 80 % according to a survey in [6]. Thus, many taxi businesses are interested in taxis' moving patterns in order to decrease the consumption of fuels as well as increase the income. Today the fuels consumption is more important to protect the environments. The real pattern analysis using data mining is an interesting topic due to importance of the value-added information, as shown in [10]. Especially, [4] has proposed the framework and perform the clustering analysis for the location recommendation. The given analysis framework consists of both of the data processing framework and the analysis processing framework for the refined clustering. However, it does not provide association rules effectively so far. In this regard, this paper is to propose the association-specific analysis framework, perform the detailed analysis for association rule mining, and provide concrete pattern information. The proposed analysis processing framework consists of both of the given data processing framework [4] and a novel analysis processing framework for the refined association analysis aiming at detecting the frequent or meaningful moving patterns. This framework can help perform even a sophisticated and a quick analysis. Particularly, the paper is to develop the refined association analysis by means of taking into account interesting factors such as the driving type of taxi, boarding hour, pick-up/drop-off time, pick-up/drop-off area, boarding rate, and so on.

The result of analysis shows that the approach employed in the paper can be extended to similar areas such as developing public transportation systems, distribution systems, and so on.

The paper is organized as follows. Section 2 proposes the analysis process framework for the location tracking data collected from the Taxi Telematics system. Section 3 exhibits the summary of the results obtained from the refined association analysis to search the taxis' moving patterns, based on whole data, short driving data, and long driving data, respectively. Finally, Section 4 concludes the paper.

2 Proposed Framework

This section proposes both of the data processing framework and the analysis processing framework, in order to show effectively how to analyze the moving patterns of taxis using association analysis. The proposed frameworks provide the flow chart which would take a quick look around both of the data flow and the overall analysis process. Also, it helps quickly perform the same or similar analysis, while sparing the temporal and economic costs.

2.1 Proposed Data Processing Framework

Figure 1 outlines the overall procedure for location data processing. Taxis report their location records to the central call server every minute. Each record includes the basic GPS data such as timestamp, latitude, longitude, direction, speed taxi ID as well as status fields (related with pick-up and drop-off). First of all, an analyzer transforms the type of such records into data type readable in SAS then stores them with SAS dataset types. We can analyze the location dataset using SAS analysis engine including data mining analysis tool, Enterprise Miner (or E-Miner).

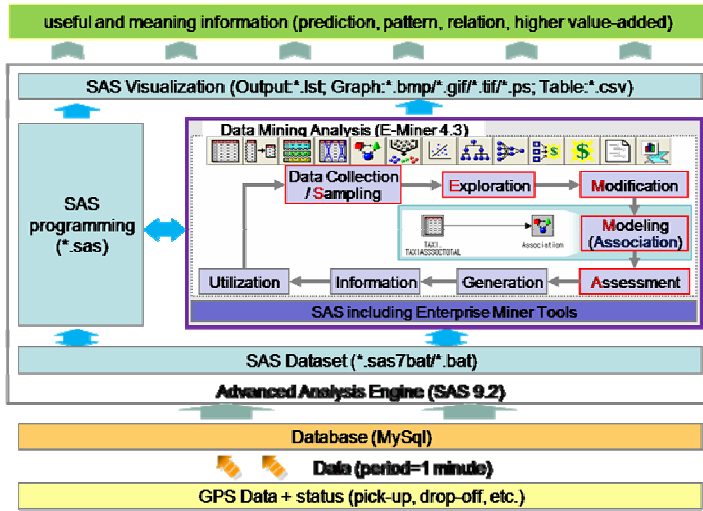


Fig. 1. Data processing framework including SAS Data Mining

To obtain more useful information, data mining analysis process consists of the data cleansing process, the modeling process which analyzes the refined data using a proper method, the assessment process, the utilization process, and feedback process, as shown in Fig. 1. As shown in Fig. 1, Modeling, fourth step of SEMMA [10], fits a predictive model; modeling a target variable using a cluster analysis, an association analysis, a regression model, a decision tree, a neural network, or a user-defined model. The paper exploits an association method for data modeling. Association searches interesting association or correlation relationships among a large set of data items.

2.2 Proposed Processing Analysis Processing for the Association Analysis

Fig. 2 shows the analysis processing framework for the refined association rule mining based on the taxi location records to inform their valuable moving patterns and moreover recommend the best location and time to the taxi drivers. The framework shows an overall process of the association analysis.

The processing framework could be expansively adapted to another transport industry, similar to a taxi business, such as delivery by changing some variables or generating new variables, and repairing the flow chart partially. As shown in Fig. 2, a data analyzer should load source records first, restore them as the SAS data type, modify the SAS dataset through the variable transformation and data filtering for data cleaning. We can create new data extracting the paired data by a pick-up and drop-off record using data split/merge. Then we select the interesting or significant variables from the paired dataset by the result of clustering. We adjust each variable's category by the result of frequency method, resulting in TaxiType (driving type of Taxi), SArea/EArea (pick-up area/drop-off area), SDay/EDay (pick-up day/drop-off day), STimeR/ETimeR (pick-up time/drop-off time). TaxiUseRatioR (boarding rate in terms of hour) variable, which means the boarding rate according to a driving hour, is obtained by dividing the number of pick-up report into the number of total report. The

analyzer completes the integrated dataset by merging original dataset and extracted interesting factors. It also needs new dedicated dataset for association analysis consisting of ID and Target variable. We divide association dataset into two by TaxiType variable with ShortDriveTaxi (short-driving) and the others (long-driving). The analyzer can return the feedback to the previous step (extraction of interesting factors), if not detect strong association rules. There is a recurrence of such a refined data mining analysis process, adjusting the feedback factors which mean modification of the given factor category and detection of a new factor until obtaining the good association rules.



Fig. 2. The analysis processing framework for the association rule mining of taxi records

2.3 Extracted Interesting Factors

Table 1 shows the extracted interesting factors, or categorical variables, obtained from clustering and frequency method using the refined and integrated SAS dataset. Note that each frequency should be alike within each variable’s category, or not biased. As shown in Table 1, we omit EArea, EDay, ETimeR variables (related with drop-off), which correspond to SArea, SDay, and STimeR (related with pick-up), due to space limitation of the paper as well as similar frequency distribution as pick-ups. Note that STimeR variable provides additionally the average frequency per hour in parenthesis, due to non-equalled time range.

Table 1. Summary of each interesting factor

Variable	Category	Description	Frequency	Percent (%)
TaxiType	ShortDriveTaxi	~ 30minutes	77447	95.51
	LongDriveTaxi	30minutes ~ 2hours	2309	2.85
	TourTaxi	2hours ~	1331	1.64
DrivingMinR	5MinDriving	1~5minutes	44507	54.89
	5-10MinDriving	5~10minutes	22254	27.44
	10-15MinDriving	10~15minutes	8787	10.84
	15~30MinDriving	15~30minutes	5539	6.83
TaxiUseRatioR	LowUse	10~15%	33082	40.80
	MiddleUse	15~20%	23389	28.84
	HighUse	20~25%	24616	30.36
SArea	SAirport	Ariport	3727	4.60
	SOldTown	Old town	43644	53.82
	SNewTown	New town	22929	28.28
	SEtc.	The others	10787	13.30
SDay	SMonday	Monday	9932	12.25
	STuesday	Tuesday	11405	14.07
	SWednesday	Wednesday	12335	15.21
	SThursday	Thursday	11582	14.28
	SFriday	Friday	12996	16.03
	SSaturday	Saturday	12033	14.84
	SSunday	Sunday	10804	13.32
STimeR	SOfficeStart	07~10o'clock	10469 (3490)	12.91
	SOfficeHour	10~18o'clock	30673 (3834)	37.83
	SOfficeEnd	18~21o'clock	17434 (5811)	21.50
	SEvening	21~23o'clock	8626 (4316)	10.64
	SNight	23~01o'clock	6906 (3453)	8.52
	SDawn	01~07o'clock	6979 (1163)	8.61

3 Association Analysis Results

This section shows the results obtained from the refined association analysis. In the previous section, we have split the taxi association datasets into three categories, whole (81,087 records), short driving (77,447 records), and long driving (3,640 records). Thus we present three association results according to each dataset. Association has a support level and a confidence level, two main measures of rule interestingness. In addition, lift is a measure of strength of the association rule. Typically, association rules are considered interesting if they satisfy both a minimum support threshold and a minimum confidence threshold. Such thresholds can be set by users or domain experts. Association rule mining is a two-step process: 1) finding all frequent itemsets; 2) generating strong association rules from the frequent itemsets. Note that we set minimum confidence to 10% by using the default association setting.

3.1 Analysis Results for Whole Taxis

Table 2 shows the meaningful and interesting association rules extracted from the association analysis based on the whole location data records. Stronger associations

Table 2. Association rules obtained from whole taxi records

Confidence	Support	Lift	Rule
100.00	10.64	3.29	SEvening ==> HighUse (rule 1)
70.04	15.06	2.31	SOfficeEnd ==> HighUse (rule 2)
83.53	31.60	2.05	SOfficeHour ==> LowUse (rule 3)
59.62	7.70	2.07	SOfficeStart ==> MiddleUse (rule 4)
100.00	5.69	3.29	SOLDTown & SEvening ==> HighUse (rule 5)
69.94	8.22	2.30	SOLDTown & SOfficeEnd ==> HighUse (rule 6)
74.85	8.79	1.46	SOLDTownTown SOfficeEnd & ==> EOldTown
83.40	17.61	2.04	SOLDTown & SOfficeHour ==> LowUse (rule 7)
72.40	8.50	1.46	SOLDTown & SOfficeEnd ==> ShortDriveTaxi & EOldTown
83.17	8.61	2.04	SOfficeHour & SNewTown ==> LowUse (rule 8)
33.80	5.42	1.11	SFriday ==> HighUse
45.03	5.52	1.10	SMonday ==> LowUse
39.38	5.84	0.97	SSaturday ==> LowUse
41.99	5.60	1.03	SSunday ==> LowUse
40.46	5.78	0.99	SThursday ==> LowUse
41.87	5.89	1.03	STuesday ==> LowUse
40.56	6.17	0.99	SWednesday ==> LowUse
82.41	4.80	2.02	SOfficeHour & SFriday ==> LowUse
84.48	5.03	2.07	SWednesday & SOfficeHour ==> LowUse
71.78	6.31	1.40	SFriday & SOLDTown ==> EOldTown
69.82	6.14	1.41	SFriday & SOLDTown ==> ShortDriveTaxi & EOldTown

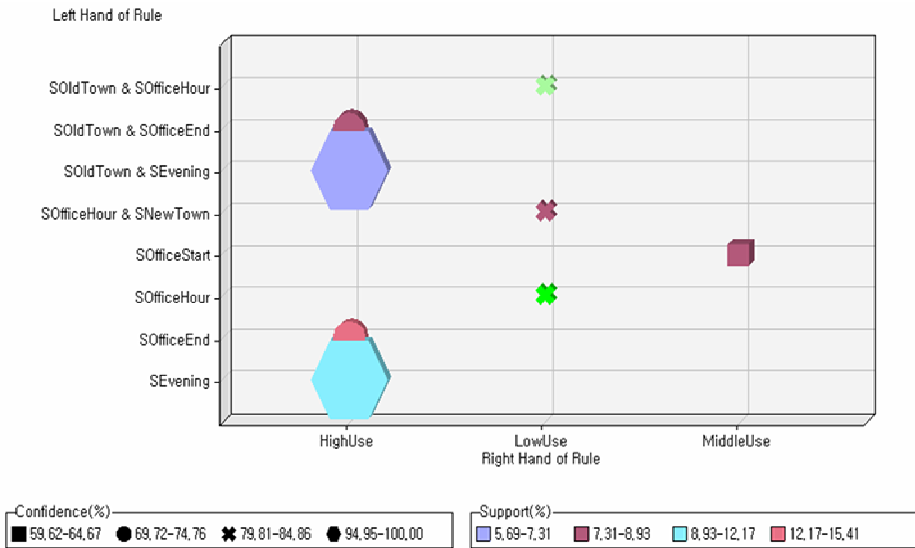


Fig. 3. The main association rules of whole taxi location records

are selected according the support level and the confidence level. The table shows that passengers use mostly taxi at the closing time of the office hours with 15.06% of the support level and 70.04% of the confidence level while the support level and the confidence level is 10.64% and 100%, respectively, for the evening hours (21~23 o'clock). Also it shows that only Friday is associated with HighUse of the taxies with 5.42% of the support level and 33.80% of the confidence level.

We depict graphs to compare those strong rules in Fig. 3 and Fig. 4. In Fig. 3, each polygon's shape, its color, and its size are regarded as confidence, support, lift, respectively. Fig. 4 shows the rule comparison in terms of two important factors, the confidence level and the support level.

3.2 Analysis Results for Short Driving Taxies

Table 3 shows the associations based on the short-driving taxi records. Eight strong associations are selected. Fig. 5 and Fig. 6 show the comparative graph among those

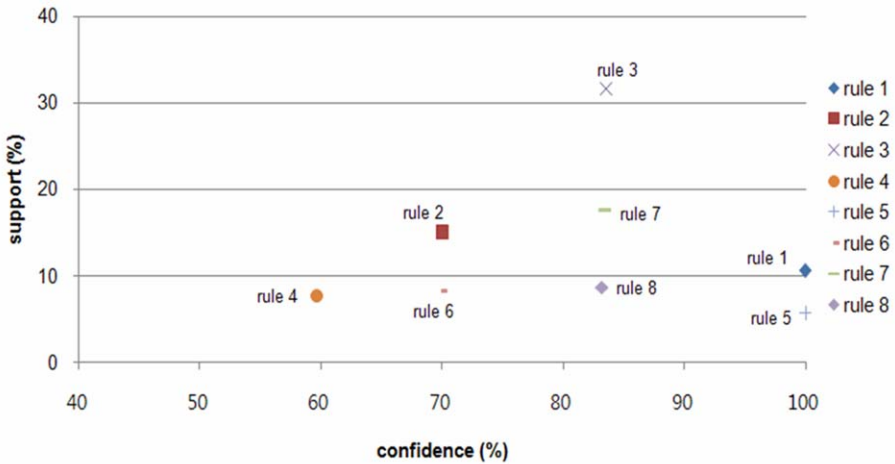


Fig. 4. Confidence and support of the main association rules

Table 3. Association rules obtained from short-driving taxi records

Confidence	Support	Lift	Rule
100.00	10.67	3.31	SEvening ==> HighUse (rule 1)
54.66	4.69	1.81	SNight ==> HighUse (rule 2)
45.34	3.89	1.57	SNight ==> MiddleUse (rule 3)
69.86	14.87	2.31	SOfficeEnd ==> HighUse (rule 4)
83.58	31.63	2.05	SOfficeHour ==> LowUse (rule 5)
69.85	8.18	2.31	SOldTown & SOfficeEnd ==> HighUse (rule 6)
70.45	4.46	2.33	SOfficeEnd & SNewTown ==> HighUse (rule 7)
100.00	3.11	3.31	SNewTown & SEvening ==> HighUse (rule 8)
75.89	13.61	1.44	SNewTown & ENewTown ==> 5MinDriving
68.69	27.33	1.30	SOldTown & EOldTown ==> 5MinDriving
76.03	8.90	1.47	SOldTown & SOfficeEnd ==> EOldTown
46.59	5.46	1.56	SOldTown & SOfficeEnd ==> EOldTown & 5MinDriving

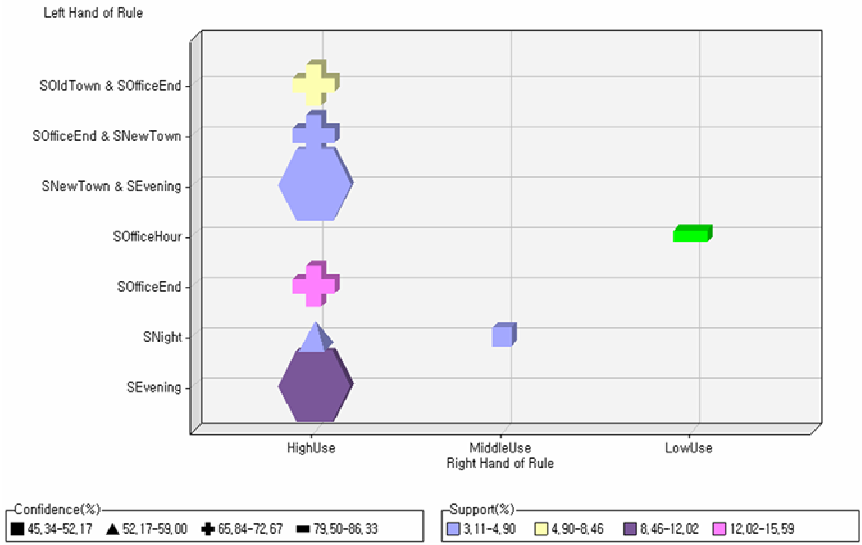


Fig. 5. The main association rules of short-driving taxis

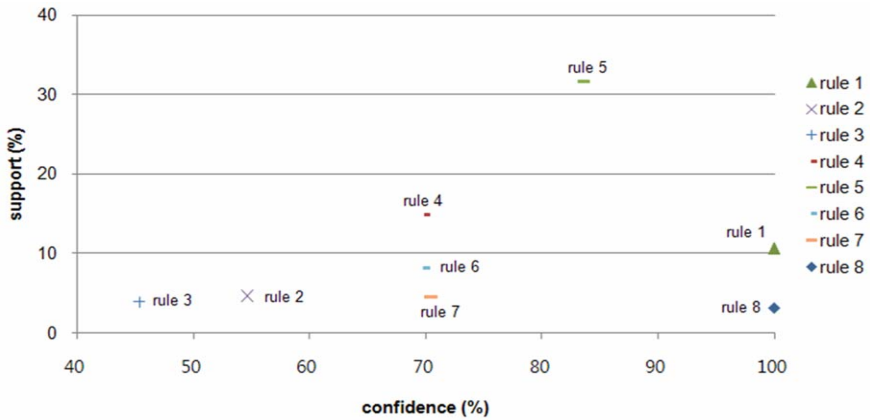


Fig. 6. Confidence and support of the main association rules for short-driving taxis

strong rules. The main association rules mean that short-driving taxis pick up many passengers at the close of office hours with support level 14.87% and confidence level 69.86%, respectively. Also mainly passengers move either within the old town or the new town within 5 minutes by a taxi. Especially, many taxis picking up a passenger move within the old town at the close of office hours.

3.3 Analysis Results for Long Driving Taxis

Taxi drivers are interested in the long driving pattern than the short due to the economic profit and the fuel cost, because the short driving consumes more fuel than

the long driving. Table 4 shows the meaningful and interesting associations based on the long-driving taxi records. Fig. 7 and Fig. 8 depict the main association rules. Finally, the main association indicates that long-driving taxis mainly move from old town to old town, especially at the close of office hours.

Table 4. Association rules obtained from long-driving taxi records

Confidence	Support	Lift	Rule
49.22	3.49	1.25	SDawn ==> LowUse (rule 1)
50.78	3.60	1.84	SDawn ==> MiddleUse (rule 2)
100.00	10.05	3.02	SEvening ==> HighUse (rule 3)
56.23	4.09	1.70	SNight ==> HighUse (rule 4)
43.77	3.19	1.59	SNight ==> MiddleUse (rule 5)
73.12	18.98	2.21	SOfficeEnd ==> HighUse (rule 6)
82.43	30.93	2.10	SOfficeHour ==> LowUse (rule 7)
40.00	4.84	1.02	SOfficeStart ==> LowUse (rule 8)
60.00	7.25	2.17	SOfficeStart ==> MiddleUse (rule 9)
100.00	4.56	3.02	SOldTown & SEvening ==> HighUse (rule 10)
100.00	3.05	3.02	SNewTown & SEvening ==> HighUse (rule 11)
71.87	8.98	2.17	SOldTown & SOfficeEnd ==> HighUse (rule 12)
82.27	14.53	2.10	SOldTown & SOfficeHour ==> LowUse (rule 13)
81.60	7.55	2.08	SOfficeHour & SNewTown ==> LowUse (rule 14)
61.18	2.86	2.22	SOldTown & SOfficeStart ==> MiddleUse (rule 15)
49.21	22.36	1.30	SOldTown ==> EOldTown
51.21	6.40	1.36	SOldTown & SOfficeEnd ==> EOldTown
28.62	6.40	1.10	SOldTown & EOldTown ==> SOfficeEnd

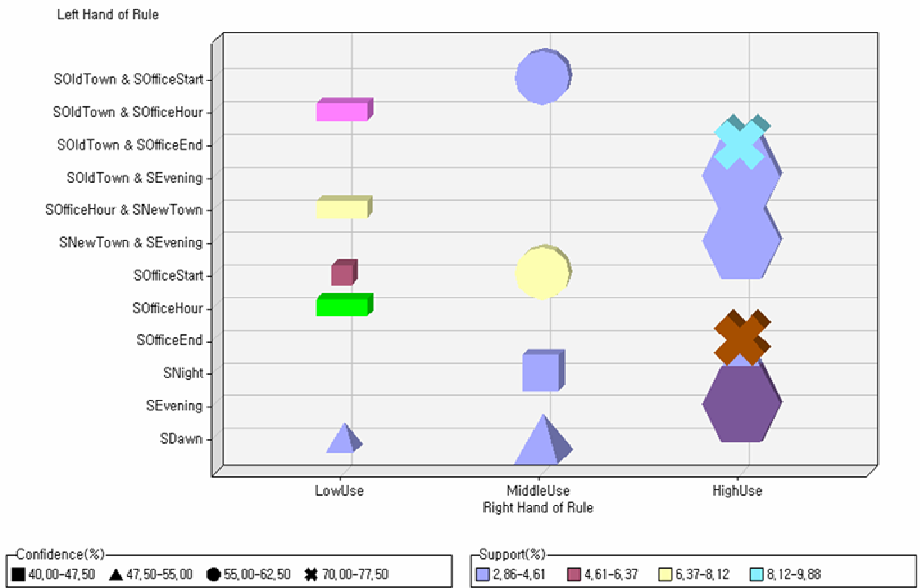


Fig. 7. The main association rules of long-driving taxis

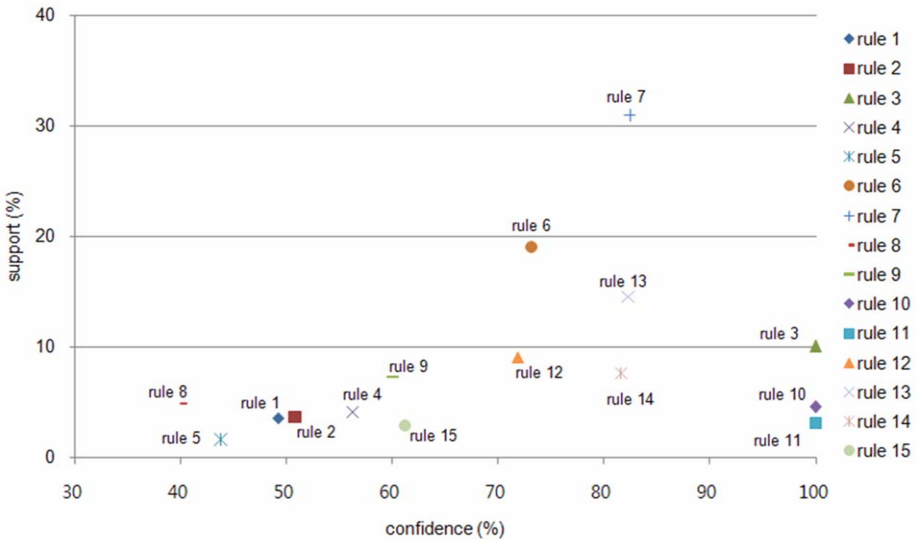


Fig. 8. Confidence and support of the main association rules for long-driving taxis

4 Conclusions

The paper has developed an analysis processing framework for detecting the hidden association rules of moving objects. In order to show the effectiveness of the framework, the framework has been applied to the real-life location history data accumulated from the Taxi Telematics system developed in Jeju Island. The proposed framework provides the flow chart which would have a glance around the overall analysis process as well as help save temporal and economic costs consumed to quickly address the same or similar mining analysis. The paper has extracted diverse interesting categorical factors such as taxi’s driving type, driving time, driving area, boarding rate, and so on. The approach enables us to obtain meaningful and value-added association rules of moving objects by performing repeatedly the refined association analysis as follow: 1) obtaining the integrated dataset through the data cleaning process such as filtering and split-merge, 2) extracting the various categorical variables from the integrated dataset after performing the primary analysis such as frequency and clustering, 3) performing the method based on another dedicated dataset created for association, and 4) deducing the meaningful and value-added pattern information such moving pattern, or 5) returning the feedback in order to appropriately adjust the categories of extracted factors, if desired. The approach can be applied not only to taxi business but also to similar areas such as public transportation services and distribution industries.

References

- [1] Lee, J., Park, G., Kim, H., Yang, Y., Kim, P., Kim, S.: A telematics service system based on the Linux cluster. In: Shi, Y., van Albada, G.D., Dongarra, J., Sloot, P.M.A. (eds.) ICCS 2007. LNCS, vol. 4490, pp. 660–667. Springer, Heidelberg (2007)

- [2] Hariharan, R., Toyama, K.: Project Lachesis: Parsing and modeling location histories. In: Egenhofer, M.J., Freksa, C., Miller, H.J. (eds.) *GIScience 2004*. LNCS, vol. 3234, pp. 106–124. Springer, Heidelberg (2004)
- [3] Lee, J., Park, G.: Design and implementation of a movement history analysis framework for the taxi telematics system. In: *Asia-Pacific Conference on Communications*, pp. 1–4 (2008)
- [4] Shin, I., Park, G., Saha, A., Kwak, Y., Kim, H.: Analysis of Moving Patterns of Moving Objects with the Proposed Framework. In: Gervasi, O., Taniar, D., Murgante, B., Laganà, A., Mun, Y., Gavrilova, M.L. (eds.) *ICCSA 2009*. LNCS, vol. 5593, pp. 443–452. Springer, Heidelberg (2009)
- [5] Lee, J., Hong, J.: Design and implementation of a spatial data processing engine for the telematics network. *Applied Computing and Computational Science* (2008)
- [6] Lee, J.: Traveling pattern analysis for the design of location-dependent contents based on the Taxi telematics system. In: *International Conference on Multimedia, Information Technology and its Applications* (2008)
- [7] Liao, Z.: Real-time taxi dispatching using global positioning systems. *Communication of the ACM*, 81–83 (2003)
- [8] He, H., Jin, H., Chen, J., McAullay, D., Li, J., Fallon, T.: Analysis of Breast Feeding Data Mining Methods. In: *Proc. Australasian Data Mining Conference*, pp. 47–52 (2006)
- [9] Madigan, E.A., Curet, O.L., Zrinyi, M.: Workforce analysis using data mining and linear regression to understand HIV/AIDS prevalence patterns. *Human Resource for Health* 6 (2008)
- [10] Matignon, R.: *Data Mining Using SAS Enterprise Miner*. Wiley, Chichester (2007)

An Efficient ICI Cancellation Method for Cooperative STBC-OFDM Systems

Kyunghoon Won, Jun-Hee Jang, Se-bin Im, and Hyung-Jin Choi

School of Information and Communication Engineering,
Sungkyunkwan University,
Suwon, Korea

{kairaess,hellojjh,yuner,hjchoi}@ece.skku.ac.kr

Abstract. In this paper, an efficient inter carrier interference (ICI) cancellation method for cooperative space time block coded orthogonal frequency division multiplexing (STBC-OFDM) system is presented. In cooperative STBC-OFDM system, ICI cancellation is needed because ICI due to the separated local oscillators always exists. To solve the complexity problem of matrix inversion operation in zero forcing method which shows the best performance, ICI cancellation method using sparse matrix decomposition (SMD) has been proposed. However, overall complexity is increased in proportion to the third order of required tap size which also increases in proportion to FFT size or carrier frequency offsets (CFOs). Considering implementation issue, the conventional method still has not sufficiently overcome the performance versus complexity trade-offs. Therefore, we propose an ICI cancellation method that focuses on solving practical complexity problem of conventional method. The proposed method adaptively decides the required tap size of each sparse matrix through signal to interference and noise ratio (SINR) measurement to reduce the complexity of conventional method, and we verified that the proposed method improves the performance versus complexity trade-offs compared with conventional method.

Keywords: Equalization, Frequency offset, ICI cancellation, STBC, OFDM.

1 Introduction

Spatial diversity offers prominent benefits in link reliability and spectral efficiency through the use of multiple antennas at the transmitter and/or receiver side [1], [2]. Unfortunately, the use of multiple-antenna techniques might not be practical for some reasons; especially due to the antenna size and power constraints. For this reason, recently, there has been significant interest in a class of techniques known as cooperative communication, which allows single-antenna mobiles to reap some of the benefits of multiple-antenna systems by cooperation of in-cell users [3], [4].

Since space time block codes (STBC) were developed originally for frequency flat channels [2], an effective way to use them on frequency selective channels is to use them along with orthogonal frequency division multiplexing (STBC-OFDM) [5]. By using space time coded with OFDM, frequency selective channel is converted into multiple frequency flat channels.

However, due to the distributed nature of cooperative STBC-OFDM system, each of the transmitters has a separate local oscillator, and synchronization becomes a difficult problem or requires feedback process to achieve. Furthermore, due to the physically separated local oscillators, different CFOs are received in one receiver antenna, and they cannot be easily assumed identical. Therefore, inter carrier interference (ICI) cancellation process is always needed in cooperative STBC-OFDM systems, because ICI caused by residual CFOs always exists at the receiver side.

Generally, zero forcing method has the best performance by multiplying the inversion of ICI matrix to the received signal. However it requires very high implementation complexity for matrix inversion. Therefore, in order to solve the complexity problem of matrix inversion operation, many proposals have been presented [6]-[9]. Some of the proposals achieve the complexity reduction, such as simplified zero forcing (ZF) method using sparse matrix decomposition (SMD), and iterative ICI cancellation method using the iterative operation of ICI estimation and cancellation [8], [9].

Especially, the ICI cancellation method which uses SMD can reduce the complexity of matrix inversion, because the major portion of ICI components are distributed in the diagonal or near diagonal position of ICI matrix [7], and the inversion operation of ICI matrix can be approximated by using SMD proposed in [10]. However, considering time-varying multi-path fading channel, it still have not sufficiently overcome the performance versus complexity trade-offs.

Therefore, in this paper, we propose an efficient ICI cancellation method which can efficiently reduce the practical complexity of conventional method without performance loss by using adaptive sparse matrix decomposition (ASMD).

This paper is organized as follows. In Section 2, we introduce the cooperative STBC-OFDM system model and ICI problem. The proposed ICI cancellation method using adaptive sparse matrix decomposition method is described in Section 3. In Section 4, the results of performance and complexity comparisons between the conventional method and the proposed method are presented. A brief conclusion is drawn in Section 5.

2 System Model

The cooperative communication system considered in this paper includes one source node, one destination node, and one of relay node. The basic premise in this work is that source and relay have information of their own to send, and would like to cooperate in order to send this information to the receiver according to the cooperative STBC scheme as shown in Fig. 1. That is, we do not consider signal distortion in the source-to-relay link. In Fig. 1, two consecutive symbols on subcarrier k , $X_{k,0}$ and $X_{k,1}$, are encoded and transmitted from the source and the relay antenna, respectively, to obtain the diversity order 2.

During even symbol period, $X_{k,0}$ and $X_{k,1}$ are transmitted from the source antenna and relay antenna, and during odd symbol period, $(-X_{k,1})^*$ and $(X_{k,0})^*$ are transmitted from the source antenna and the relay antenna, where $(\cdot)^*$ denotes the complex conjugate operation.

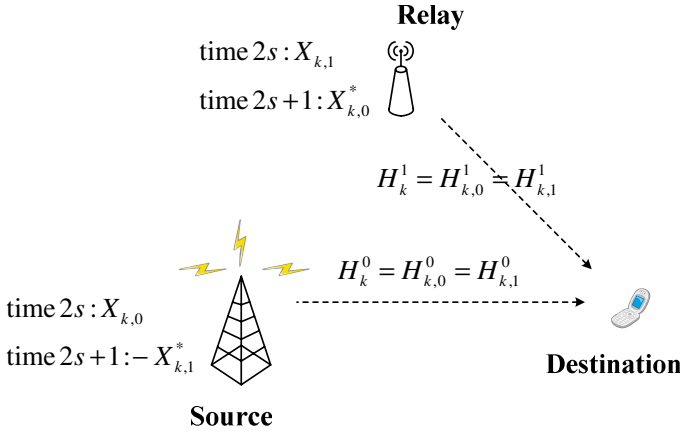


Fig. 1. Schematic representation of cooperative STBC-OFDM system

With the assumption of $H_k^j = H_{k,0}^j = H_{k,1}^j$, which is a case of quasi-static channel for two symbol period, the received signal $Y_{k,i}$ for symbol duration 0 and 1 ($i=0,1$) can be expressed as

$$\begin{aligned}
 Y_{k,0} &= H_k^0 \cdot X_{k,0} + H_k^1 \cdot X_{k,1} + W_k^0, \\
 Y_{k,1} &= H_k^0 \cdot (-X_{k,1})^* + H_k^1 \cdot (X_{k,0})^* + W_k^1, 0 \leq k \leq N-1
 \end{aligned}
 \tag{1}$$

where N is FFT size, $H_{k,i}^j$ is channel frequency response (CFR) of received signal on the j -th link, and W_k^j is zero-mean white Gaussian noise with variance σ^2 .

From (1), we can easily obtain diversity gain by using simple linear combination as described in [2]

$$\begin{aligned}
 \hat{X}_{k,0} &= (H_k^0)^* \cdot Y_{k,0} + H_k^1 \cdot (Y_{k,1})^* \\
 &= (|H_k^0|^2 + |H_k^1|^2) \cdot X_{k,0} + (H_k^0)^* \cdot W_k^0 + H_k^1 \cdot (W_k^1)^* \\
 \hat{X}_{k,1} &= (H_k^1)^* \cdot Y_{k,0} - H_k^0 \cdot (Y_{k,1})^* \\
 &= (|H_k^0|^2 + |H_k^1|^2) \cdot X_{k,1} + (H_k^1)^* \cdot W_k^0 - H_k^0 \cdot (W_k^1)^*.
 \end{aligned}
 \tag{2}$$

With the residual CFOs, the received signal for symbol duration 0 and 1 can be expressed as [6]

$$Y_{k,0} = \Gamma_0^0 \cdot H_k^0 \cdot X_{k,0} + \Gamma_0^1 \cdot H_k^1 \cdot X_{k,1} + \sum_{\substack{i=0 \\ i \neq k}}^{N-1} (\Gamma_{k-i}^0 \cdot H_i^0 \cdot X_{i,0} + \Gamma_{k-i}^1 \cdot H_i^1 \cdot X_{i,1}) + W_k^0$$

$$Y_{k,1} = -\Gamma_0^0 \cdot H_k^0 \cdot X_{k,1}^* + \Gamma_0^1 \cdot H_k^1 \cdot X_{k,1}^* + \sum_{\substack{i=0 \\ i \neq k}}^{N-1} (-\Gamma_{k-i}^0 \cdot H_i^0 \cdot X_{i,0}^* + \Gamma_{k-i}^1 \cdot H_i^1 \cdot X_{i,1}^*) + W_k^1 \quad (3)$$

$$\Gamma_k^j = \frac{\sin(\pi \cdot \varepsilon_j)}{N \cdot \sin(\pi(\varepsilon_j - k)/N)} \cdot \exp(j\pi\varepsilon_j(N-1/N) - j\pi(\varepsilon_j - k)/N)$$

where ε_j is the residual CFO from the j -th link which is normalized by the subcarrier spacing, and the second parts in (3) means the ICI components represented by a sum for $i \neq k$ on subcarrier k caused by CFOs.

Fig. 2 shows the performance degradations caused by ICI in cooperative STBC-OFDM systems. In spite of perfect channel estimation, from Fig. 2, we can see that there are severe performance degradations due to the loss of orthogonality between subcarriers.

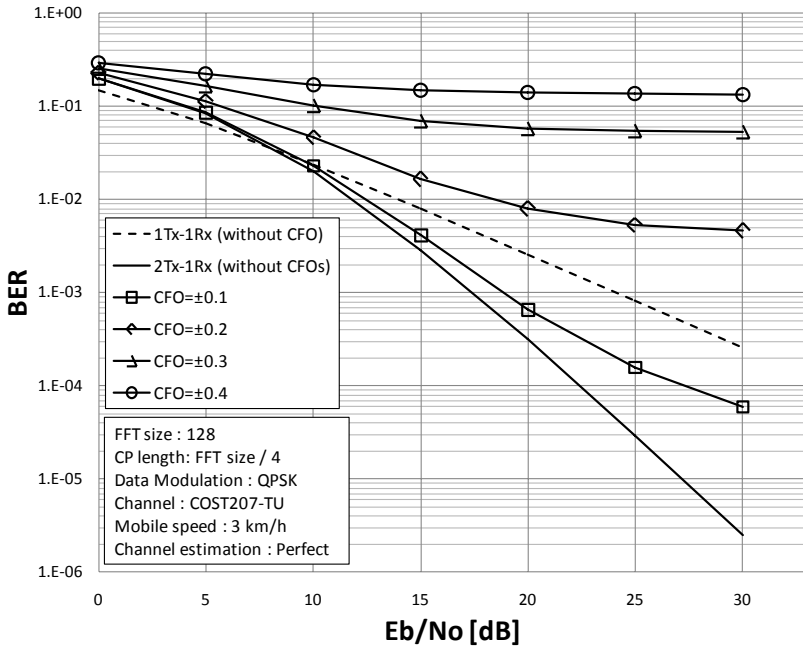


Fig. 2. Performance degradations caused by ICI in cooperative STBC-OFDM systems

3 An Efficient ICI Cancellation Method

Equation (3) can be expressed as matrix form in (4).

$$\begin{bmatrix} \mathbf{Y}_0 \\ \mathbf{Y}_1^* \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{X}_0 \\ \mathbf{X}_1 \end{bmatrix} + \begin{bmatrix} \mathbf{W}_0 \\ \mathbf{W}_1^* \end{bmatrix}$$

$$\mathbf{Y}_i = \begin{bmatrix} Y_{0,i} \\ Y_{1,i} \\ \vdots \\ Y_{N-1,i} \end{bmatrix}, \mathbf{X}_i = \begin{bmatrix} X_{0,i} \\ X_{1,i} \\ \vdots \\ X_{N-1,i} \end{bmatrix}, \mathbf{W}_i = \begin{bmatrix} W_{0,i} \\ W_{1,i} \\ \vdots \\ W_{N-1,i} \end{bmatrix} \quad (4)$$

where \mathbf{A} , \mathbf{B} , \mathbf{C} , and \mathbf{D} are $N \times N$ size sub-matrices.

$$\mathbf{A} = \begin{bmatrix} \Gamma_0^0 H_0^0 & \Gamma_{-1}^0 H_1^0 & \cdots & \Gamma_{1-N}^0 H_{N-1}^0 \\ \Gamma_1^0 H_0^0 & \Gamma_0^0 H_1^0 & \cdots & \Gamma_{2-N}^0 H_{N-1}^0 \\ \vdots & \vdots & \ddots & \vdots \\ \Gamma_{N-1}^0 H_0^0 & \Gamma_{N-2}^0 H_1^0 & \cdots & \Gamma_0^0 H_{N-1}^0 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} \Gamma_0^1 H_0^1 & \Gamma_{-1}^1 H_1^1 & \cdots & \Gamma_{1-N}^1 H_{N-1}^1 \\ \Gamma_1^1 H_0^1 & \Gamma_0^1 H_1^1 & \cdots & \Gamma_{2-N}^1 H_{N-1}^1 \\ \vdots & \vdots & \ddots & \vdots \\ \Gamma_{N-1}^1 H_0^1 & \Gamma_{N-2}^1 H_1^1 & \cdots & \Gamma_0^1 H_{N-1}^1 \end{bmatrix} \quad (5)$$

$$\mathbf{C} = \begin{bmatrix} \Gamma_0^{1*} H_0^{1*} & \Gamma_{-1}^{1*} H_1^{1*} & \cdots & \Gamma_{1-N}^{1*} H_{N-1}^{1*} \\ \Gamma_1^{1*} H_0^{1*} & \Gamma_0^{1*} H_1^{1*} & \cdots & \Gamma_{2-N}^{1*} H_{N-1}^{1*} \\ \vdots & \vdots & \ddots & \vdots \\ \Gamma_{N-1}^{1*} H_0^{1*} & \Gamma_{N-2}^{1*} H_1^{1*} & \cdots & \Gamma_0^{1*} H_{N-1}^{1*} \end{bmatrix} \quad \mathbf{D} = \begin{bmatrix} -\Gamma_0^{0*} H_0^{0*} & -\Gamma_{-1}^{0*} H_1^{0*} & \cdots & -\Gamma_{1-N}^{0*} H_{N-1}^{0*} \\ -\Gamma_1^{0*} H_0^{0*} & -\Gamma_0^{0*} H_1^{0*} & \cdots & -\Gamma_{2-N}^{0*} H_{N-1}^{0*} \\ \vdots & \vdots & \ddots & \vdots \\ -\Gamma_{N-1}^{0*} H_0^{0*} & -\Gamma_{N-2}^{0*} H_1^{0*} & \cdots & -\Gamma_0^{0*} H_{N-1}^{0*} \end{bmatrix}$$

Ignoring system complexity, for the perfect cancellation of ICI, we can simply multiply the inverse of ICI matrix to the received signal matrix (zero forcing method). However, as you can see in (4), the implementation complexity of the inversion operation of $2N \times 2N$ size ICI matrix is too high.

To reduce the complexity of matrix inversion operation in zero forcing method, [7] proposed the ICI cancellation method using SMD introduced in [10]. The major portion of ICI components in (5) are distributed in the diagonal or near diagonal position of ICI matrix. Therefore, by using SMD, inversion of ICI matrix can be transmitted into a series of inversions of its diagonal sub-blocks as illustrated in Fig. 3.

By using the SMD method, (4) can be approximated as

$$\begin{bmatrix} \hat{\mathbf{Y}}_{n,0} \\ \hat{\mathbf{Y}}_{n,1}^* \end{bmatrix} = \begin{bmatrix} \mathbf{A}_n & \mathbf{B}_n \\ \mathbf{C}_n & \mathbf{D}_n \end{bmatrix} \begin{bmatrix} \hat{\mathbf{X}}_{n,0} \\ \hat{\mathbf{X}}_{n,1} \end{bmatrix} + \begin{bmatrix} \hat{\mathbf{W}}_{n,0} \\ \hat{\mathbf{W}}_{n,1}^* \end{bmatrix}$$

$$\hat{\mathbf{X}}_{n,i} = \begin{bmatrix} X_{n,0} & X_{n+1,0} & \cdots & X_{n+2q,0} \end{bmatrix}^T \quad (6)$$

$$\text{where } \hat{\mathbf{Y}}_{n,i} = \begin{bmatrix} Y_{n,0} & Y_{n+1,0} & \cdots & Y_{n+2q,0} \end{bmatrix}^T, \quad 0 < n \leq N - (2q + 1)$$

$$\hat{\mathbf{W}}_{n,i} = \begin{bmatrix} W_{n,0} & W_{n+1,0} & \cdots & W_{n+2q,0} \end{bmatrix}^T$$

where n is sparse matrix index, $2q+1$ is the tap size of diagonal sub-blocks, and \mathbf{A}_n is shown in (7). $a_{n,n}$ is the elements of \mathbf{A} , \mathbf{B}_n , \mathbf{C}_n , and \mathbf{D}_n can be obtained similarly [7].

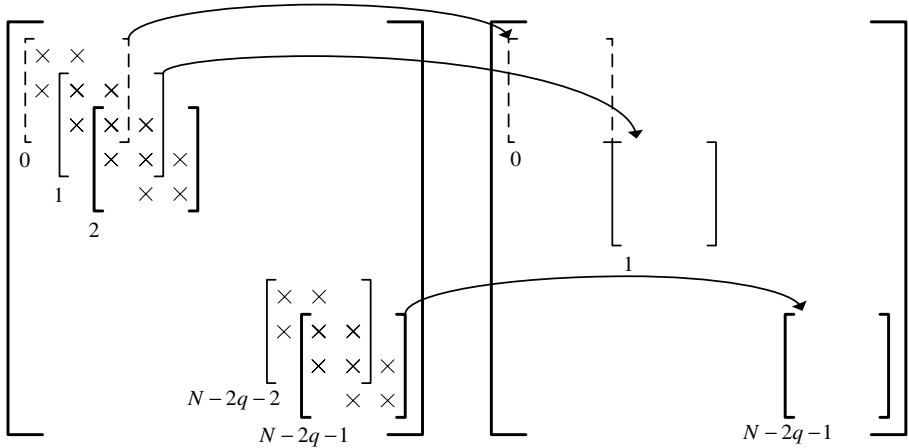


Fig. 3. Sparse matrix decomposition proposed in [10]

$$\mathbf{A}_n = \begin{bmatrix} a_{n,n} & a_{n,n+1} & \cdots & a_{n,n+q} & \cdots & 0 \\ a_{n+1,n} & a_{n+1,n+1} & & & & \vdots \\ \vdots & & & & & 0 \\ a_{n+q,n} & & & & & a_{n+q,n+2q} \\ 0 & & & & & \vdots \\ \vdots & & & & a_{n+2q-1,n+2q-1} & a_{n+2q-1,n+2q} \\ 0 & \cdots & a_{n+2q,n+q} & \cdots & a_{n+2q,n+2q-1} & a_{n+2q,n+2q} \end{bmatrix} \quad (7)$$

Therefore, mitigation of ICI is accomplished by multiplying the inversion of sparse matrix as follows

$$\begin{bmatrix} \hat{\mathbf{X}}_{n,0} \\ \hat{\mathbf{X}}_{n,1} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_n & \mathbf{B}_n \\ \mathbf{C}_n & \mathbf{D}_n \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{Y}_{n,0} \\ \mathbf{Y}_{n,1}^* \end{bmatrix} = \begin{bmatrix} \mathbf{X}_{n,0} \\ \mathbf{X}_{n,1} \end{bmatrix} + \begin{bmatrix} \mathbf{A}_n & \mathbf{B}_n \\ \mathbf{C}_n & \mathbf{D}_n \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{W}_{n,0} \\ \mathbf{W}_{n,1} \end{bmatrix}. \quad (8)$$

In the case of the above SMD method, although the total number of complex multiplications for the matrix inversion is reduced from $(2N)^3$ to $(N-2q) \cdot (2(2q+1))^3$, required complexity remains high because the tap size q increases as increasing CFOs and FFT size N .

However, considering time and frequency selectivity due to mobile movement and multi-path delay, respectively, each subcarriers are unequally interfered each other. Therefore, there is no necessity for the size of sparse matrix to be fixed. If each sparse matrix has adequately small tap size q^* within maximum tap size q depending on the amount of ICI from each neighboring subcarriers, the total amount of complex multiplications for the matrix inversion is reduced to less than $(N-2q) \cdot (2(2q+1))^3$.

From this point of view, we propose an adaptive sparse matrix decomposition (ASMD) method which is based on a kind of SINR measurement to adaptively decide

the efficient tap size q^* . Without loss of generality, we assume that all transmitted signals have average power 1, i.e., $E[\mathbf{X}^H\mathbf{X}]=\mathbf{I}$, where $(\)^H$ denotes the Hermitian transpose.

In the proposed ASMD method, the decision procedure of efficient tap size q_n^* of the n -th sparse matrix can be expressed as

$$\Psi_n^j = \begin{cases} \min \left\{ u \left| \frac{|\Gamma_{\pm(u+1)}^j \cdot H_{k\pm(u+1)}|^2}{|\Gamma_0^j \cdot H_k|^2} < \Psi_q^j \right. \right\}, u = 0, 1, \dots, q-1, \\ q, \text{ if } \frac{|\Gamma_{\pm(u+1)}^j \cdot H_{k\pm(u+1)}|^2}{|\Gamma_0^j \cdot H_k|^2} \geq \Psi_q^j \text{ for all } u = \{0, 1, \dots, q-1\} \end{cases} \quad (9)$$

$$\Psi_q^j = \left(|\Gamma_q^j|^2 + \sigma^2 \right) / |\Gamma_0^j|^2, k = n + q$$

where $\left| \Gamma_{\pm(u+1)}^j \cdot H_{k\pm(u+1)} \right|^2 / \left| \Gamma_0^j \cdot H_k \right|^2$ means signal to interference ratio (SIR) of the k -th subcarrier to the $\pm(u+1)$ -th subcarrier in frequency selective fading channel with CFO, Ψ_q^j is based on SINR of the k -th subcarrier to the $(k+q)$ -th subcarrier in frequency flat channel, and noise power σ^2 is inserted to the threshold Ψ_q^j to reduce unnecessary operation in low signal to noise ratio (SNR) environments.

From (8) and (9), with the assumption of quasi-static channel for two symbol period, efficient tap size q_n^* of n -th sparse matrix can be determined as follows

$$q_n^* = \max \{ \Psi_n^j \} \quad (10)$$

Therefore, by using the proposed ASMD method, the total amount of complex multiplications for ICI cancellation in cooperative STBC-OFDM systems can be approximately reduced from $(N-2q) \cdot (2(2q+1))^3$ to $(N-2q) \cdot (1+2 \cdot 4 \cdot (q^*+1) + (2(2q^*+1))^3)$, $q^* \leq q$. For each $(N-2q)$ times of sparse matrix, $1+2 \cdot 4 \cdot (q^*+1)$ times of complex multiplication is needed for the calculation of k -th subcarrier power $\left| \Gamma_0^j \cdot H_k \right|^2$ and SINR comparisons.

4 Simulation Results

In this section, we investigate the performance of the proposed method on multipath channel for cooperative STBC-OFDM system. The simulation parameters used in this paper are shown in Table 1. The multi-path fading channel is generated by COST 207 TU model and in order to avoid inter symbol interference (ISI) we set the length of cyclic prefix (CP) long enough.

Table 1. Simulation parameters

Parameters	Values
FFT size	128
Cyclic prefix (CP) length	FFT size / 4
System bandwidth	FFT size * 15 kHz
Subcarrier spacing	15 kHz
Center frequency	2 GHz
Channel model	COST 207 TU
Channel estimation	Perfect
Data modulation	QPSK
Normalized frequency offsets	± 0.1

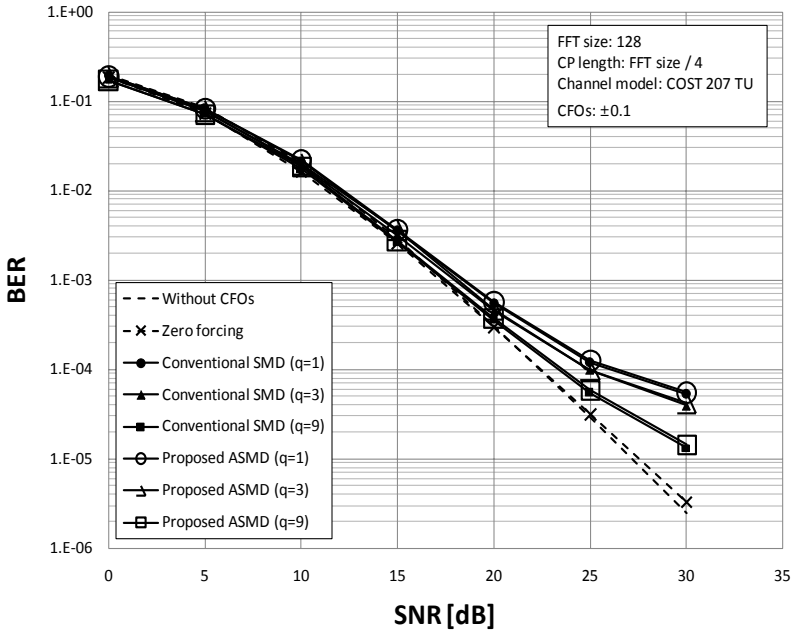
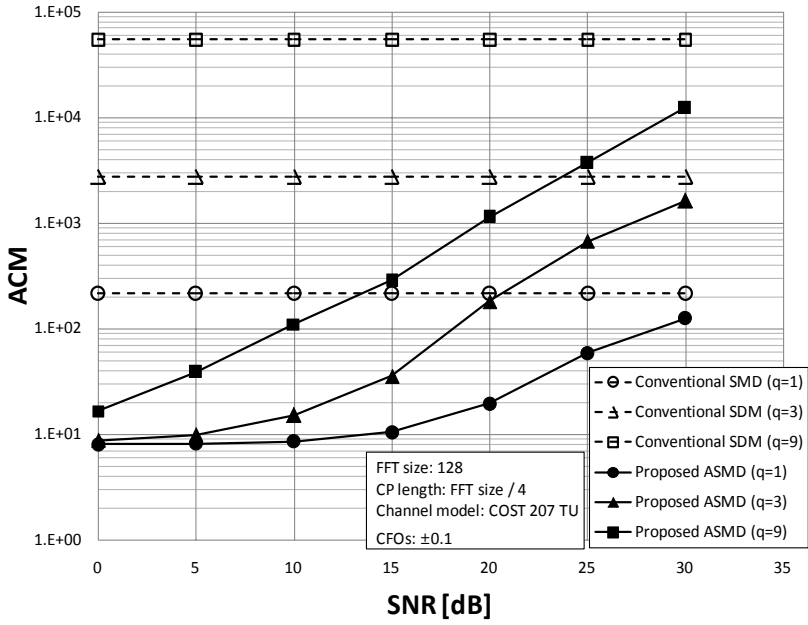
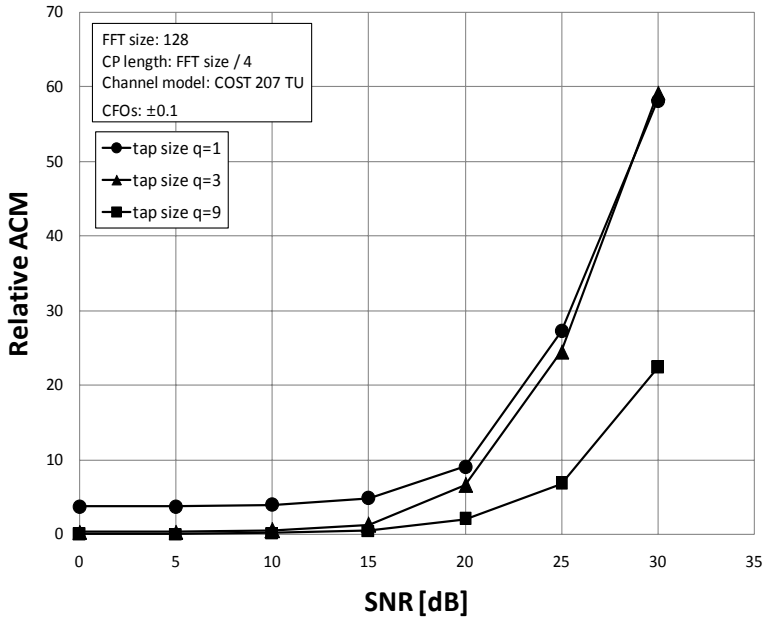
**Fig. 4.** BER performance of ICI cancellation methods versus SNR

Fig. 4 shows the bit error rate (BER) performances comparisons between conventional SMD and proposed ASMD for ICI cancellation. From Fig. 4, we can see that zero forcing method shows almost the same performance with and without CFOs. Also, we can see that each performance of the proposed ASMD method is very close to that of the conventional SMD method for the same number of taps.

Fig. 5 shows the average number of complex multiplications (ACM) and relative ACM for matrix inversion per each subcarrier. In Fig 5, in low SNR environment, overall ACM is dramatically reduced by using proposed ASMD, because threshold Ψ_q^j depends on the noise variance σ^2 ($\sigma^2 \approx 1$). And in high SNR environment, Ψ_q^j



(a) ACM versus SNR



(b) Relative ACM versus SNR

Fig. 5. Complexity comparisons between conventional SMD and proposed ASMD

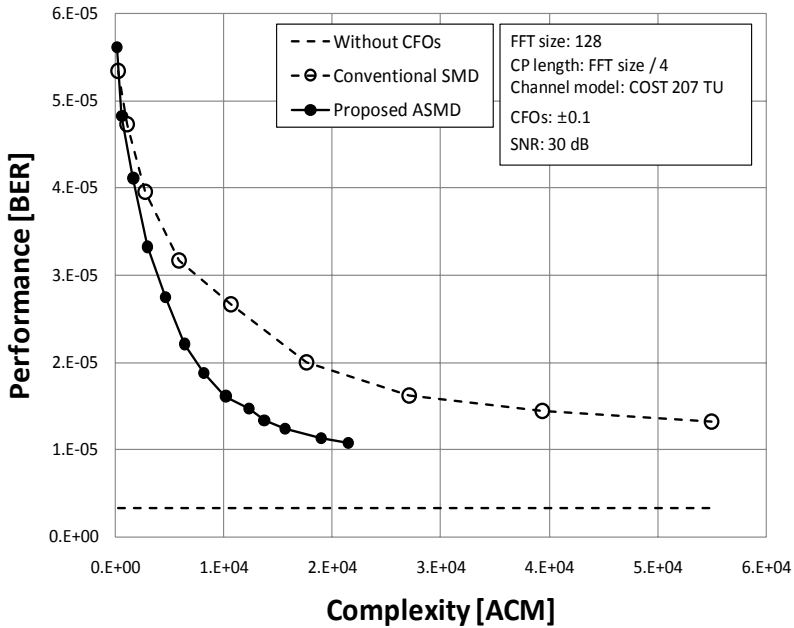


Fig. 6. BER performances versus AMC

depends on the SIR of at each subcarriers ($\sigma^2 \approx 0$). Therefore, although the ACM of the proposed ASMD method increases in proportion to SNR, it is always lower than the conventional SMD method, and the relative ACM ratio improves as the number of taps increases.

Fig. 6 compares performance versus complexity trade-offs of conventional SMD method and proposed ASMD method at SNR=30 dB environment. As shown in Fig 6, both conventional and proposed method, the more performance is improved, the more complexity is needed for the same amount of performance enhancement, because the amount of ICI on the k -th subcarrier from the $(k+q)$ -th subcarrier is described in the form of sinc function. However, we can see that the proposed ASMD method efficiently improves performance versus complexity trade-offs compared with the conventional SMD method by SINR measurement and adaptive tap allocation.

5 Conclusions

In this paper, we proposed an efficient ICI cancellation method for cooperative STBC-OFDM system. Through the frequency domain SINR measurement and adaptive tap allocation, we showed that the overall complexity for ICI cancellation can be efficiently reduced with approximately same performance compared with the conventional method. We expect that the proposed ASMD scheme can be applicable for various ICI suppression or cancellation methods.

Acknowledgments. This research was supported by the Ministry of Knowledge Economy, Korea, under the ITRC (Information Technology Research Center) support program supervised by the IITA (Institute of Information Technology Assessment) (IITA-2009-C1090-0902-0005).

References

1. Tarokh, V., Seshadri, N., Calderbank, A.R.: Space-Time Codes for High Data Rate Wireless Communication: Performance Criterion and Code Construction. *IEEE Trans. Inf. Theory* 44(2), 744–765 (1998)
2. Alamouti, S.M.: A Simple Transmit Diversity Technique for Wireless Communications. *IEEE J. Sel. Areas Commun.* 16(8), 1451–1458 (1998)
3. Sendonaris, A., Erkip, E., Aazhang, B.: User Cooperation Diversity-Part 1. *IEEE Trans. Commun.* 51(11), 1927–1937 (2003)
4. Laneman, J.N., Wornell, G.W.: Distributed Space Time Block Coded Protocols for Exploiting Cooperative Diversity in Wireless Networks. *IEEE Trans. Inf. Theory* 49, 2415–2425 (2003)
5. Lee, K.F., Williams, D.B.: A Space-Time Coded Transmitter Diversity Technique for Frequency Selective Fading Channels. In: *IEEE Sensor Array and Multichannel Signal Processing Workshop*, March 2000, pp. 149–152 (2000)
6. Mingqian, T.Z., Premkumar, A.B., Madhukmar, A.S.: Performance Investigation of STBC-OFDM Systems with Frequency Offset and a Semi-Blind Approach for the Correction. In: *IEEE Vehicular Technology Conference*, May 2004, pp. 1836–1839 (2004)
7. Li, Z., Qu, D., Zhu, G.: An Equalization Technique for Distributed STBC-OFDM System with Multiple Carrier Frequency Offsets. In: *IEEE Wireless Communications and Networking Conference*, April 2006, vol. 2, pp. 839–843 (2006)
8. Zhang, W., Qu, D., Zhu, G.: Performance Investigation of Distributed STBC-OFDM System with Multiple Carrier Frequency Offsets. In: *IEEE Personal, Indoor and Mobile Radio Communications Symposium*, September 2006, pp. 1–5 (2006)
9. Kim, Y.Z., Lee, H., Chung, H.K., Cho, W.S.: An Iterative Decoding Technique for Cooperative STBC-OFDM systems with Multiple Frequency Offsets. In: *IEEE Personal, Indoor and Mobile Radio Communications Symposium*, September 2007, pp. 1–5 (2007)
10. Jeon, W.G., Chang, K.H., Cho, Y.S.: An Equalization Technique for Orthogonal Frequency-Division Multiplexing Systems in Time-Varying Multipath Channels. *IEEE Trans. Commun.* 49, 27–32 (2001)

Low-Cost Two-Hop Anchor Node-Based Distributed Range-Free Localization in Wireless Sensor Networks

Taeyoung Kim¹, Minhan Shon¹, Wook Choi²,
MoonBae Song³, and Hyunseung Choo¹

¹ School of Information and Communication Engineering,
Sungkyunkwan Univ., Korea

tyenterprise@gmail.com, {minari95, choo}@skku.edu

² Department of Computer Science and Engineering,
Hankuk Univ. of Foreign Studies, Korea
twchoi@hufs.ac.kr

³ Digital Media and Communications Business,
Samsung Electronics, Korea
mbsong@gmail.com

Abstract. Due to the fact that most of the sensor network applications are based on the location information of sensor nodes, localization is an essential research area. Some localization schemes in the literature require sensor nodes to have additional devices to measure the distance or angle between two sensor nodes, but it is not suitable for low-cost sensor nodes. On the other hand, other localization schemes use only the connectivity information of sensor nodes so that localization is not much accurate enough. In this paper, we propose a range-free localization scheme, called *Low-cost Two-hop Anchor Node-based Distributed Range-free Localization (LADL)*, which offers a higher accuracy with lower cost than the previous works. LADL exploits a small portion of anchor nodes which know their own location beforehand. In LADL, sensor nodes collect the location information of the anchor nodes within two-hop distance and calculate their own location using a grid-scan algorithm. The simulation results show that LADL has a maximum of 12% lower delivery cost of location information messages and 25% higher accuracy than DRLS [8].

Keywords: Sensor Networks; Localization; Anchor node.

1 Introduction

Wireless sensor networks (WSNs) consist of a large number of low-cost sensor nodes with high resource constraints. A wide variety of applications and networking protocols in the WSNs such as environment, habitat monitoring, smart home/spaces, intrusion detection, object tracking, and geographic routing, are based on the position information of sensor nodes. The more accurate such position information is, the more reliable and efficient the service becomes. The

simplest method of accurately determining the position information of sensor nodes is to install the Global Positioning System (GPS) on each sensor node. However, it is not suitable for the low-cost sensor nodes. Therefore, the localization of such low-cost sensor nodes is a challenging problem to be tackled.

The localization schemes can be classified into two kinds – range-based and range-free. The range-based localization scheme is based on a specific technique such as Time of Arrival (ToA) [1], Time Difference of Arrival (TDoA) [2], Angle of Arrival (AoA) [3], and Received Signal Strength Indicator (RSSI) [4], in order to measure the distances or angles between transmitter and receiver. In this scheme, sensor nodes are required to have additional devices to run the distance or angle measurement technique. In addition, the measurement accuracy of this technique is affected by radio interferences like noise and fading and thus eventually having low localization accuracy. From this perspective, range-based localization is indeed not suitable for wireless sensor networks. Unlike the range-based schemes, in the range-free localization schemes, such as Centroid [5], Convex Position Estimation (CPE) [6], and Distance Vector-hop (DV-hop) [7], localization is accomplished using only the connectivity information between the sensor nodes. However, on the downside, the range-free schemes have higher cost of exchanging messages and calculating positions than the range-based schemes. Moreover, their positioning accuracy is relatively high.

Recently, Distributed Range-free Localization Scheme (DRLS) [8] has been proposed to enhance the localization accuracy. Assuming that there are a few anchor nodes in the sensor field that knows their own location information beforehand, DRLS allows the normal nodes to calculate their position. Each normal node calculates its initial position with the location information of one-hop anchor nodes using a grid scan algorithm, and enhances the localization accuracy with the location information of two-hop anchor nodes by means of vector-based refinement. To discover the anchor nodes, two-hop flooding is used but it increases the cost of message transmissions between the sensor nodes and may cause localization to be more inaccurate when the position balance between the anchor nodes is not matched (DRLS will be discussed in detail in section 2).

This paper proposes a range-free localization scheme, called Low-cost Two-hop Anchor-based Distributed Range-free Localization (LADL), which provides a higher accuracy with lower transmission cost than DRLS. All the normal nodes in LADL use the position information of a few anchor nodes to calculate their respective positions. LADL uses two-hop flooding as in DRLS but its message transmission cost is lower by a maximum of 12% than in DRLS as the broadcasting procedure is skipped after the two-hop flooding. In addition, the localization accuracy is higher by a maximum of 25% than in DRLS without even using the vector-based refinement process, and LADL eliminates inappropriate computations to the sensor nodes by applying a grid scan algorithm based on the position information of the anchor nodes within two hops.

The remainder of the paper is organized as follows. Section 2 presents DRLS and our motivations while section 3 introduces assumptions and basic definitions.

Section 4 describes the proposed scheme, LADL in detail. Section 5 discusses simulation results. Finally, section 6 concludes the paper.

2 Related Work

Centroid [5] is a fundamental scheme for range-free localization. In Centroid, the normal nodes find their position by collecting location information of neighboring anchor nodes and calculating average coordinate of them, but if the anchor node ratio is low, the accuracy of estimated positions becomes very low. To estimate the position of the normal nodes more accurately, CPE was proposed which is a centralized range-free localization in [6]. Anchor nodes in this scheme broadcast their position information, and normal nodes estimate their positions by constructing an Estimative Rectangle (*ER*). However, traffic load for message exchange is heavy, and thus it is not suitable for large scale sensor networks.

As mentioned above, DRLS [8] is a distributed range-free localization scheme. It calculates a position of a normal node based on the fact that both one-hop anchor nodes and two-hop anchor nodes can affect the position accuracy of the normal node. Fig. 1 shows how a two-hop anchor node, A_3 , can affect the position estimation of N . A normal node can exist anywhere within the data transmission radio range of two one-hop anchor nodes, A_1 and A_2 , but not within that of the two-hop anchor node, A_3 . The distribution of these anchor nodes enables the normal node to calculate its position more accurately. However, such a distribution of anchor nodes is not always possible.

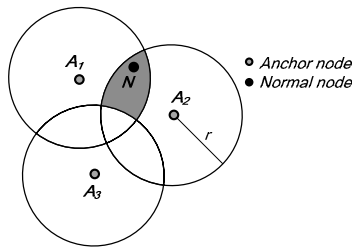


Fig. 1. The influence of an anchor node within two-hop distance

DRLS consists of three steps: beacon message exchange, grid scan, and refinement steps. In the beacon message exchange step, the normal node should know the position information of the anchor nodes existing within a two-hop distance from their neighboring anchor nodes by two-hop flooding and one additional broadcast. Thereafter, in the grid scan step, each normal node creates an estimative rectangle (*ER*) with the use of the transmission radio range of the one-hop anchor nodes and divides the *ER* into a set of small grid cells. Fig. 2 illustrates that how a normal node obtains a one-hop anchor-based initial estimated position using a grid scan algorithm. In the refinement step, the normal node that measured the initial position decides its own final estimated position via vector-based refinement. The direction of the vector, called virtual

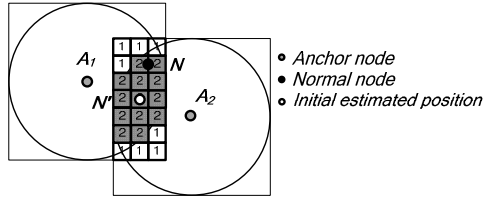


Fig. 2. Initial estimated position with grid scan algorithm in DRLS

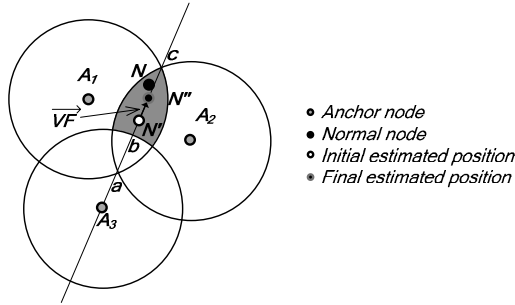


Fig. 3. Location refinement based on the vector in DRLS

force (VF), is from a two-hop anchor node to the initial estimated position, and its size becomes $|\vec{VF}| = \frac{ab \cdot N'c}{ac}$. Fig. 3 shows how a normal node refines its initial position based on a two-hop anchor node.

As DRLS conducts localization with the use of more information of the anchor nodes compared to that in the previous studies, it shows relatively high accuracy; however, it has several problems. First, as it uses two-hop flooding to collect the position information of two-hop anchor nodes, the message transmission cost becomes higher than that in the existing range-free methods. Second, its square-root computation for calculating the vector in the refinement step is not appropriate for low-cost sensor nodes [9]. Third, as it accomplishes the vector-based refinement, its location accuracy becomes lower if the distribution of the anchor nodes is not even. Fig.4 shows an example of the refinement which eventually lowers the position accuracy. Suppose that the initial estimated position of normal node N is N' . The sum of two vectors made by the anchor nodes A_3 and A_4 exists in an area where the normal node cannot be placed. Thus, the position accuracy becomes low. Followings are the motivations of this study to enhance DRLS by tackling the above-mentioned problems:

- Normal nodes in DRLS obtain the position information of two-hop anchor nodes by one-hop anchor nodes and ignore the position information of the anchor nodes that can be obtained without going through the one-hop anchor nodes. However, if they use such ignored information, however, and if additional broadcasting is skipped after two-hop flooding, the message transmission cost can be lowered.

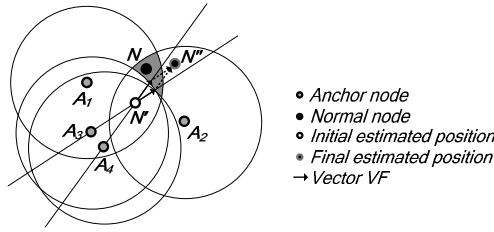


Fig. 4. The influence of farther anchors that are not distributed uniformly

- Square-root computation is not appropriate to low-cost sensor nodes [9]. Besides, vector-based refinement may cause lower position accuracy due to the anchor nodes that are placed in similar positions. Instead of skipping vector-based refinement, it can enhance the position accuracy and lessen the computational cost by utilizing all the one- and two-hop anchor nodes in grid scan step.

A summary of the differences between DRLS and LADL is shown in Table 1.

Table 1. Differences between DRLS and LADL

Head	DRLS	LADL
Selection of anchor nodes	A two-hop flooding and a broadcasting	A two-hop flooding
Construction of ER	Only one-hop anchor nodes	Both one- and two-hop anchor nodes
Grid Scan	Only one-hop anchor nodes	Both one- and two-hop anchor nodes
Refinement	Vector based refinement	Skipping the refinement step

3 Preliminaries

We assume the followings in this study:

- Normal nodes that do not have information of their own positions, and a few anchor nodes that have information of their positions, constitute the entire sensor field.
- All the sensor nodes are distributed randomly in a large sensor field, and cannot change their positions once these have been decided.
- All the sensor nodes can be identified through their unique IDs and have the same transmission radio and sensing ranges.
- The collision that can occur during data transmission is not considered, and there are no interferences such as fading and noise. In other words, a sensor node receives data without fail as it transmits data to its neighbor sensor nodes within the transmission radio range.

Anchor nodes within a one-hop distance from a normal node are called one-hop anchor nodes while anchor nodes within a two-hop distance from a normal node are called two-hop anchor nodes. Normal nodes can enhance localization accuracy by separating the one-hop anchor nodes from the two-hop anchor nodes among their neighboring anchor nodes. Fig. 5 shows where normal nodes can be located depending on the one-hop anchor nodes and the two-hop anchor nodes. Anchor node A is within a one-hop distance from normal node N_1 , thus becoming a one-hop anchor node of N_1 . On the other hand, it becomes a two-hop anchor node of N_2 since it is two hops away from normal node N_2 . Accordingly, the region where N_1 can exist is a circular area with a radius of transmission radio range r whereas the region where N_2 can exist is the circular area with a radius of $2r$ minus the circular area with a radius of r .

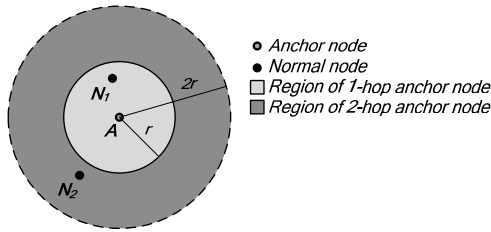


Fig. 5. The region that a normal node can be located from an anchor node

4 Proposed Scheme: LADL

LADL consists of three several steps - anchor node selection, ER construction, and grid scan. In the anchor node selection step, each normal node collects the position information of the anchor nodes which are located within a two-hop distance. In the step of the ER construction, each normal node constitutes an ER by using the collected position information of the anchor nodes. Note that as it creates an ER using one- and two-hop anchor nodes, the two-hop anchor node-based refinement process of DRLS is skipped. In the grid scan step, the positions of the normal nodes are identified by dividing the ER into a set of grid cells and then scanning each cell. Now, let us discuss the three steps of LADL in detail below.

4.1 Anchor Node Selection

In the anchor node selection step, each anchor node transmits its own ID and position information to its neighboring nodes. Each anchor node enables two-hop flooding by assigning 2 to the Time to Live (TTL) value of its own position information message. Thus, the normal node can identify either one-hop or two-hop anchor nodes based on the value of TTL.

Fig. 6 shows an example of two-hop flooding. In Fig. 6(a), anchor nodes A_1 , A_2 , and A_3 transmit their position information with TTL=2 to their neighboring

sensor nodes. In Fig. 6(b), every sensor node that has received the position information with TTL=2 transmits the position information of neighbor anchor nodes with TTL=1 to its neighboring sensor nodes. Via two-hop flooding, normal node N_1 selects A_1 and A_2 as one-hop anchor nodes, normal node N_2 selects A_2 as a one-hop anchor node while choosing A_3 as a two-hop anchor node.

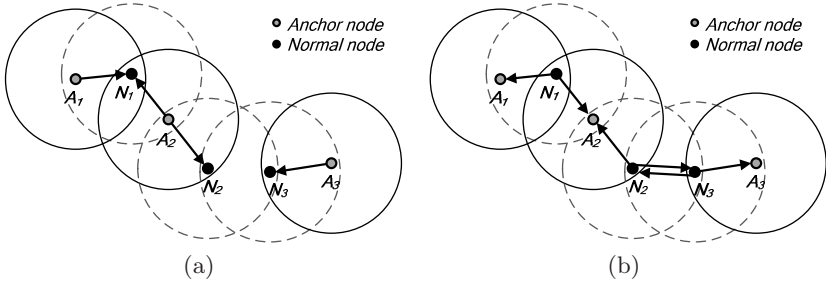


Fig. 6. Two-hop flooding: (a) Flooding with TTL=2, (b) Flooding with TTL=1

4.2 ER Construction

Every normal node can construct an *ER* with collected position information of both its one- and two-hop anchor nodes. To reduce the computation overhead for the localization, the normal node uses the rectangle that circumscribes the transmission radio range of each anchor node. Thereafter, an *ER* is constructed by calculating the overlapped region of these rectangles.

Fig. 7 shows an example of the *ER* construction. Normal node N keeps the position information of one-hop anchor node A_3 and two-hop anchor nodes A_1 and A_2 , and constructs an *ER* (i.e., dark shaded area in the figure) applying different rectangle drawing mechanisms to one- and two-hop anchor nodes, respectively. Then, the normal node should be located inside the *ER*.

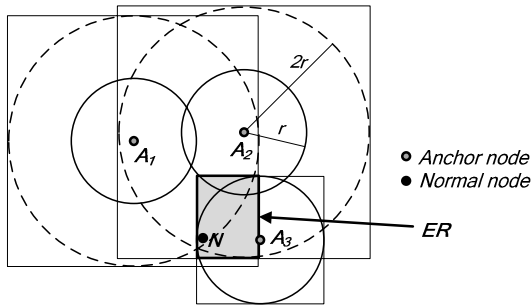


Fig. 7. ER made by the location information of one- and two-hop anchor nodes

4.3 Distributed Grid Scan Algorithm

The ER calculated by a normal node includes a region that is larger than the region in which the normal node can actually exist. Accordingly, the localization accuracy can be enhanced by excluding the region where a normal node cannot actually exist from the ER using the grid scan algorithm. After dividing the ER into a set of grid cells, the normal node decides if which grid cells should be excluded by scanning each cell. If the size of cell is too small, the scan computation cost increases, whereas the localization accuracy decreases if it is too large. In this paper, the length of one side of each cell is set at $0.1r$ (remember that r is the data transmission radio range), in consideration of both the scan computation cost and the localization accuracy. Additionally, the value of the cell is set at 1, and the coordinate of the center of gravity in each cell is selected as the coordinate value of the position of the cell. Here, 1, the value of each cell, means that a normal node could exist in the cell.

After removing the cells where each normal node cannot exist, the normal node applies a different criterion of scanning each cell for one- and two-hop anchor nodes. Fig. 8 shows an example of the grid scan algorithm. Note that the distribution of normal node N , and anchor nodes A_1 , A_2 , and A_3 is the same as in Fig. 7. Normal node N sets 1 to the value of each cell after dividing the ER into a set of $0.1rx0.1r$ grid cells. In Fig. 8, 0 is designated as the value of the cells whose distances from two-hop anchor nodes A_1 and A_2 are longer than $2r$ and shorter than r , and whose distances from one-hop anchor node A_3 are longer than r . That is, the cells with a value of 0 represent the excluded cells for the final scanning. The cells with a value of 1 (i.e., dark shaded area in Fig. 8) represents the possible cells that include the actual location of normal node N . Finally, normal node N obtains final estimated position N' by computing the average of the representative coordinate values in each shaded cell.

As all sensor nodes are distributed randomly in the sensor field, it may happen that an anchor node does not exist within the two-hop distance from the normal node. Therefore, the normal node waits until other neighbor normal nodes succeed in localization, and then accomplishes its localization by using neighbors' position information. For this purpose, the normal nodes that succeeded

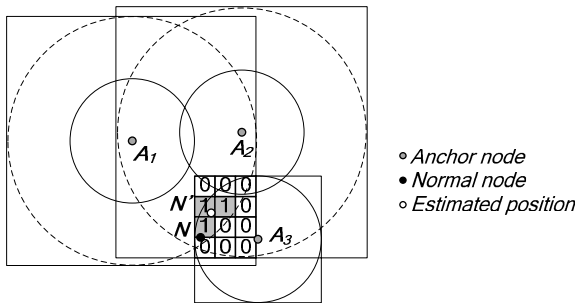


Fig. 8. The result of grid scan based on the location information of anchor nodes

in localization broadcast their own position information and that of the one-hop anchor node to its neighboring nodes.

5 Performance Evaluation

5.1 Methodology

In this section, DRLS, which has come about upon the motivation of LADL, will be compared with LADL to analyze the performance of LADL. We will use the following two metrics to measure the performance of LADL.

1. Localization accuracy: It is defined as the difference between the real position and the measured position of a normal node, and the data transmission radio range of sensor node r is the basic unit. It shows through a numerical value how accurately the normal node calculates the position compared with the real position.
2. Message transmission cost: It is defined by the number of position information messages of the anchor nodes used for each normal node to obtain information of the anchor nodes. It indicates the position information exchange cost of the anchor node to measure the positions of the normal nodes.

The simulator in this section is realized using JAVA. The data transmission radio range of the sensor nodes is r , and the simulation is performed in an environment where the sensor nodes were deployed randomly in the sensor field with $10r \times 10r$ sizes. A collision on the link, which occurs by two-hop flooding, is not considered, and when dividing the ER into specific grid sizes, it draws a grid with the size of each cell being $0.1r \times 0.1r$ in consideration of the position estimation accuracy, computational cost of grid scan, and memory capacity required. Moreover, it should be compared with the performance of DRLS by using the average value of the results derived from the 100 simulation repeats.

5.2 Localization Accuracy

Fig. 9 shows the position estimation accuracy of Centroid, CPE, DRLS, and LADL measured by changing the proportion of the anchor nodes from 5 to 47.5%, with the number of total sensor nodes fixed at 200. In addition, the result is about normal nodes which can get information of anchor nodes by two-hop flooding surely. LADL measures the position by using both one- and two-hop anchor nodes at the grid scan step; therefore, more accurate position measurement is possible by eliminating the inaccuracy that occurs in the vector-based refinement. Moreover, because the increase in the proportion of the anchor nodes causes an increase in the amount of information of the anchor nodes, position estimation accuracy increases. However, the localization without the normal nodes which does not have anchor node information is useless for applications in the WSNs. Therefore, the localization with all the normal nodes is necessary and its accuracy of DRLS and LADL is described in Fig. 10. The simulation result shows that LADL has a 4-25% higher accuracy compared to DRLS.

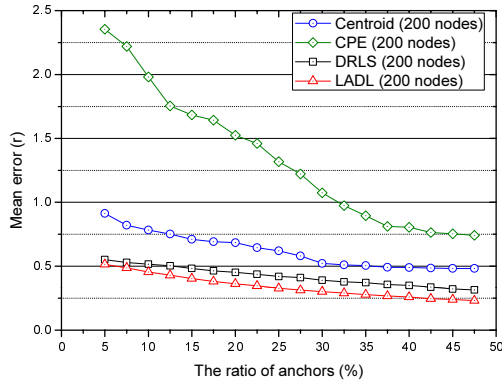


Fig. 9. Localization accuracy of normal nodes with different ratio of anchor nodes

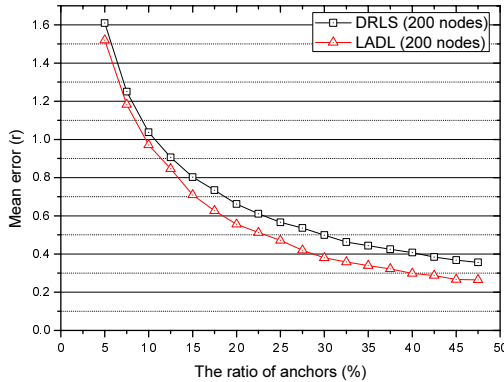


Fig. 10. Localization accuracy of all normal nodes with different ratio of anchor nodes

Fig. 11 shows how much the number of total sensor nodes affects the position estimation accuracy. The result of the accuracy of the position estimation that was done while changing the number of total anchor nodes to 200, 400, and 600 with fixing the proportion of the anchor nodes at 20% indicates that LADL has 16-49% better accuracy compared to DRLS. The increase in the number of total sensor nodes leads to more accurate position estimation because the higher density of the network causes that the amount of information of the anchor nodes that a normal node receives increases.

5.3 Message Transmission Cost

Fig. 12 shows the message transmission cost for the anchor nodes of DRLS and LADL measured by changing the proportion of the anchor node from 5 to 47.5% in the environment where the total number of sensor nodes is 200. DRLS transmits more messages compared to LADL because DRLS performs additional broadcasting in addition to two-hop flooding to provide information

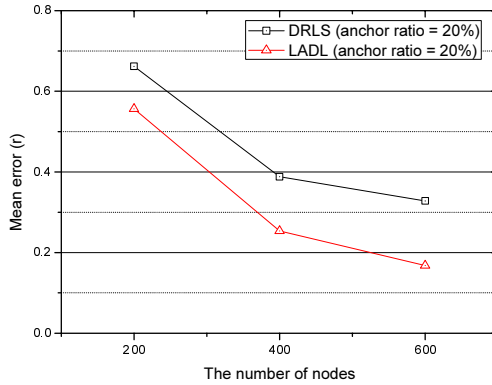


Fig. 11. Localization accuracy in terms of the total number of sensor nodes

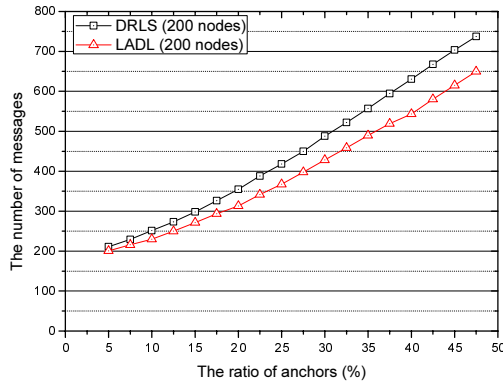


Fig. 12. Anchor node selection cost with different ratio of anchor nodes

of the anchor nodes to the normal nodes. Moreover, as the proportion of the anchor nodes increases, the number of messages increases significantly due to the fact that the number of anchor nodes which are the starting point of two-hop flooding increases. Fig. 9 shows that LADL provides the normal nodes with information of the anchor nodes with a 5-12% lower message transmission cost.

Fig. 13 shows how much the number of total sensor nodes affects the message transmission cost of two-hop flooding. The number of position information messages is measured while changing the number of total anchor nodes to 200, 400, and 600 with fixing the proportion of the anchor nodes at 20%. The result indicates that LADL has a 7-14% lower message transmission cost compared to DRLS. The reason for DRLS's poorer performance compared to LADL is that DRLS broadcasts additional position information more than LADL does. Further, the two-hop flooding cost increases rapidly as the number of total sensor nodes increases because the increase in the network density rapidly increases the number of sensor nodes existing within a one-hop distance from a sensor node.

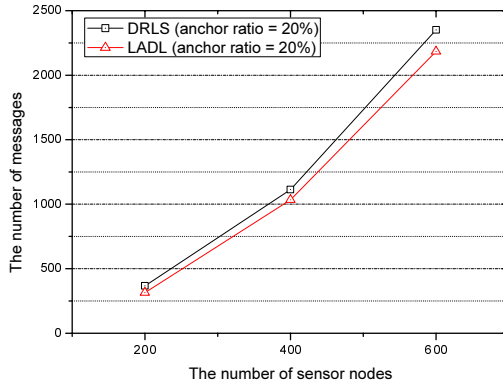


Fig. 13. Anchor node selection cost with changing of the total number of sensor nodes

6 Conclusions

In this paper, we proposed a range-free localization scheme, called low-cost two-hop anchor node-based distributed range-free localization (LADL), which can decrease the message transmission cost and can enhance the position estimation accuracy compared to DRLS. A normal node makes an ER by utilizing the position information of the anchor nodes within a two-hop distance from it, and points out where it will be located in the entire sensor field by scanning the interior of ER with a grid scan algorithm. Additionally, it decreases the cost of exchanging the position information of the anchor nodes via two-hop flooding, which is more effective than DRLS, to obtain the position information of the anchor node within a two-hop distance. The merit of LADL lies in the fact that it enables the normal node to conduct localization with a two-hop anchor node even if it does not have any one-hop anchor node. It was proven via simulation that LADL can conduct localization more precisely by as much as 25% with a message transmission cost that is up to 20% smaller than that of DRLS.

In the subsequent study, it should be proven that LADL is superior to DRLS in terms of computation cost when the square-root computation is removed from it, by analyzing the cost required by each normal node during the execution of the algorithm. Besides, we plan to study about identifying the most appropriate size of the transmission radio range and the cell at the grid scan step to obtain the higher localization accuracy.

Acknowledgments. This research was supported in part by MKE, Korea under ITRC NIPA-2009-(C1090-0902-0046), in part by MEST(Korea), under WCU Program supervised by the KOSEF (No.R31-2008-000-10062-0) and also supported partly by NIA (National Information Society Agency), KOREA under the KOREN program.

References

1. Hightower, J., Borriello, G.: Location Systems for Ubiquitous Computing. *Computer* 34(8), 57–66 (2001)
2. Pottie, G.J., Kaiser, W.J.: Wireless Integrated Network Sensors. *Communications of the ACM* 43(5), 51–58 (2000)
3. Priyantha, N.B., Chakraborty, A., Balakrishnan, H.: The Cricket Location-Support System. In: *Proceedings of ACM MobiCom 2000*, pp. 32–43 (2000)
4. Savvides, A., Han, C.-C., Strivastava, M.B.: Dynamic Fine-Grained Localization in Ad-Hoc Networks of Sensors. In: *Proceedings of ACM MobiCom 2001*, pp. 166–179 (2001)
5. Bulusu, N., Heidemann, J., Estrin, D.: GPS-Less Low Cost Outdoor Localization for Very Small Devices. *IEEE Wireless Communications* 7(5), 28–34 (2000)
6. Doherty, L., Pister, K.S.J., Ghaoui, L.E.: Convex Position Estimation in Wireless Sensor Networks. In: *Proceedings of IEEE INFOCOM 2001*, vol. 3, pp. 1655–1663 (2001)
7. Niculescu, D., Nath, B.: Ad Hoc Positioning System (APS). In: *Proceedings of IEEE Globecom 2001*, vol. 1, pp. 2926–2931 (2001)
8. Sheu, J.P., Chen, P.C., Hsu, C.S.: A Distributed Localization Scheme for Wireless Sensor Networks with Improved Grid-Scan and Vector-Based Refinement. *IEEE Transactions on Mobile Computing* 7(9), 1110–1123 (2008)
9. Acharya, M., Girao, J., Westhoff, D.: Secure Comparison of Encrypted Data in Wireless Sensor Networks. In: *Proceedings of IEEE Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks*, pp. 47–53 (2005)
10. He, T., Huang, C., Blum, B.M., Stankovic, J.A., Abdelzher, T.: Range-Free Localization Schemes for Large Scale Sensor Networks. In: *Proceedings of ACM MobiCom 2003*, pp. 81–95 (2003)
11. Boukerche, A., Oliveira, H.A.B.F., Nakamura, E.F., Loureiro, A.A.F.: Localization Systems for Wireless Sensor Networks. *IEEE Wireless Communications* 14(6), 6–12 (2007)
12. Akyildiz, I., Su, W., Sankarasubramaniam, Y., Cayirci, E.: A Survey on Sensor Networks. *IEEE Communications Magazine* 40(8), 102–114 (2002)
13. Lederer, S., Wang, Y., Gao, J.: Connectivity-based Localization of Large Scale Sensor Networks with Complex Shape. In: *Proceedings of IEEE INFOCOM 2008*, pp. 789–797 (2008)

SecDEACH: Secure and Resilient Dynamic Clustering Protocol Preserving Data Privacy in WSNs

Young-Ju Han, Min-Woo Park, and Tai-Myoung Chung

Internet Management Technology Laboratory,
Department of Computer Engineering,
School of Information and Communication Engineering,
Sungkyunkwan University,
300 Cheoncheon-dong, Jangan-gu,
Suwon-si, Gyeonggi-do, 440-746, Republic of Korea
Tel.: +82-31-290-7222; Fax: +82-31-299-6673
{yjhan,mwpark}@imtl.skku.ac.kr, tmchung@ece.skku.ac.kr

Abstract. Clustering scheme enabling the efficient utilization of the limited energy resources of the deployed sensor nodes can effectively prolong lifetime of wireless sensor networks. Like most of the protocols for wireless sensor networks, a clustering scheme is vulnerable to a number of security attacks. Especially, attacks involving cluster-heads are the most damaging since cluster-based schemes rely on their cluster-heads for routing and data aggregation. Several studies have proposed a secure clustering scheme and a resilient cluster leader election protocol in wireless sensor networks. However, most of resilient clustering schemes do not support dynamic clustering and most of secure clustering schemes are not robust and resilient against node(specially, cluster-heads) capture. In this paper, we present a secure and resilient clustering protocol on the basis of dynamic clustering protocol. Our scheme provides both a secure clustering and a resilient cluster-head election preserving dynamic clustering unlike the existing secure clustering schemes. In addition, our scheme provides data privacy through the concealed data aggregation. The analysis shows that our scheme offers the security services against both an inside attack and an outside attack while having the little memory, computation and communication overheads. The comparison with related works shows that our scheme has a great improvement in security and the reduction of the overheads compared to the existing schemes.

1 Introduction

WSNs(Wireless sensor networks) consist of a number of sensor nodes deployed over a geographical area for the purpose of monitoring certain phenomena of interest. Wireless sensor networks enable the reliable surveillance of a variety of environment for health, military, home and civil applications [7][8]. Each sensor node has wireless communication capability, data aggregating ability, and self-organizing capability [8].

When embedded in critical application, WSNs are likely to be attacked [3]. Aside from the well-known vulnerabilities due to wireless communication, WSNs lack physical protection and are usually deployed in open and unattended environments, which

makes them more vulnerable to attacks [39]. It is therefore crucial to devise security solutions to these networks [3].

Sensor nodes are often organized into several clusters for the purpose of efficiency and scalability. Clustering scheme enabling the efficient utilization of the limited energy resources of the deployed sensor nodes can effectively prolong lifetime of wireless sensor networks [7,12]. Adding security to clustering scheme is challenging, as its dynamic and periodic rearrangement of clustering makes key distribution solutions inadequate [39]. Like most of the protocols for WSNs, a clustering scheme is vulnerable to a number of security attacks including passive eavesdropping, jamming, spoofing and replay attacks [3]. Especially, attacks involving CHs(cluster-heads) are the most damaging since cluster scheme fundamentally relies on the CHs for data aggregation and routing. If the adversary manages to become a CH, he can perform various attacks to disrupt the function of the network such as routing and data aggregation. Therefore, to guarantee for WSNs to function properly, a secure clustering protocol is required.

A number of secure routing protocols have been proposed recently in WSNs [15,9,10,14]. These can be generally categorized into a secure cluster formation protocol [19] and a secure CH election protocol [5,14,10]. Secure cluster formation protocols significantly improve the security of clustering in the presence of malicious attacks. However, they only focus on the cluster formation and do not consider the security against disrupting CH election by the compromised node. Resilient cluster leader election schemes in [5,14,10] are more resilient against a node capture than the above mentioned secure cluster protocols. However, they do not support dynamic clustering which is a heuristic method for an optimal clustering since they have assumption that a cluster sector is pre-defined and not changed. This feature is not suitable for WSNs.

In this paper, we present a secure and resilient clustering scheme considering energy efficiency and a concealed data aggregation in WSNs. Our scheme provides both the resilient CH election preserving a dynamic clustering and the secure cluster formation unlike the existing secure clustering schemes. In addition, our scheme provides data privacy through a concealed data aggregation. The comparison with related works shows that our scheme offers the security services against both an inside attack and an outside attack. Our scheme has a little memory and computation overhead but it is endurable in sensor nodes.

This paper is organized as follows. Section 2 introduces the backgrounds of the proposed scheme. Section 3 presents our goals and some assumptions required for describing the proposed scheme. Section 4 describes the procedures of the proposed scheme in detail. Section 5 explains the results of analysis comparing the our proposed scheme with the other schemes. At last, section 6 summarizes our work and discusses the future works.

2 Backgrounds

2.1 DEACH and Its Vulnerabilities

DEACH. The basic operation of a clustering routing protocol in WSNs is to partition network into several clusters. After clustering, the CH of each cluster aggregates the sensed data from its member nodes for energy efficiency of data processing and then

sends the aggregated data to the BS(base station) [11][15]. The most common technique in famous clustering routing protocols is the probabilistic clustering algorithm based on the randomized CH rotation for distributing the energy consumption among sensor nodes in each cluster [11][15].

As our previous works, we have proposed DEACH(Distance based low-Energy Adaptive Clustering Hierarchy Protocol) [16] which is an energy-efficient distance based clustering scheme. In clustering scheme, how to control the balance of the energy consumption among CHs can greatly affect the network lifespan which is generally determined by the first dead-node [12]. Among the sources of energy consumption of a CH, wireless data transmission is the most critical [6]. Because the energy consumption of wireless data transmission increases generally in proportion with the distance from sending node to receiving node [6], the energy consumption of CH depends on its location from the BS. For example, in single-hop inter-cluster communication, the CHs farther to the BS will die much faster than the other CHs, while in multi-hop inter-cluster communication, the CH closest to the BS is burdened with a heavy relay traffic load and die first(i.e., the hot spot problem). Most of clustering schemes such as LEACH [15] and HEED [11] utilize mainly CH frequency or residual energy of each node as criterion of CH election and do not consider the distance from the BS to the CH. So, let those schemes cause the unbalanced energy consumption among CHs in spite those schemes are able to prolong the network lifetime.

To solve this problem, DEACH considers the distance from the BS to CHs as well as the residual energy as the criterion of CH election for the balanced energy consumption among CHs. For a single-hop inter-cluster environment, we use Eq. 1 as probability formula for electing a CH. That is, each sensor node has the probability, P_i , of becoming a CH which is determined by the distance to the BS and its residual energy.

$$P_i = \begin{cases} \frac{P'_{opt}(i)}{1 - P'_{opt}(i) \times (r \bmod \frac{1}{P'_{opt}(i)})} \times \frac{E_{res}(i)}{E_{init}(i)} & C_i(t) = 1 \\ 0 & C_i(t) = 0 \end{cases} \quad (1)$$

where P_{opt} is the optimal probability being a CH which is a system parameter, $P'_{opt}(i)$ is the optimal probability being a CH of sensor node i based on the distance from itself to the BS, r is the current round of the CH election, $E_{init}(i)$ and $E_{res}(i)$ are the initial energy and the residual energy of the sensor node, $d(i, j)$ is the distance between node i and j , d_{max} represents the distance of the farthest sensor node from the BS and d_{min} represents the distance of the closest sensor node. $C_i(t)$ presents whether or not sensor node i has been a CH in the most recent $r \bmod 1/P'_{opt}$ rounds. That is, $C_i(t) = 0$ if node i has been a CH and one otherwise). Equation 1 shows that the closer the sensor node is to the BS, the lower probability of becoming a CH it could have. In other words, the CHs closer to the BS can have more members than the CHs farther to the BS since the few CHs are elected in the area close to BS. Besides, among the nodes having the same distance from the BS, the more the sensor node has the residual energy, the higher probability of becoming a CH it could have. According to Eq. 1, every sensor node

elects itself as a CH only once during $1/P'_{opt}$ rounds on average. In this way, the CHs are selected randomly among the sensor nodes and the energy concentration on CHs is distributed.

DEACH provides fully distributed manner by utilizing local information and good energy-efficiency by performing the load balanced clustering. Through simulation experiments, we have showed that the proposed scheme is more effective than the other clustering protocols in prolonging the lifespan of WSNs [16].

Its Vulnerabilities. Adding security to clustering scheme is challenging, as its dynamic and periodic rearrangement of clustering makes key distribution solutions inadequate [39]. Like most of the protocols for WSNs, a clustering scheme is vulnerable to a number of security attacks including passive eavesdropping, jamming, spoofing and replay attacks [3]. Especially, attacks involving CHs are the most damaging since a cluster scheme fully relies on the CHs for data aggregation and routing. If the adversary manages to become a CH, he can perform various attacks to disrupt the function of the network such as routing and data aggregation.

Firstly, in existing CH election protocols, nodes rely on their local status or on reported status of peer nodes for electing their CH. For example, in LEACH [15], HEED [11] and DEACH [16], a node can independently decide to become CH based on its energy-level or a distance. However, such approaches have substantial limitations with the respect to security. A fundamental observation is that any election mechanism based on a concrete election metrics such as the residual energy level or distance, can be manipulated or forged. So, the adversary can become the CH or disturb a legitimate node not to become a CH for the upcoming epoch by forging or replaying CH announcement.

Secondly, if the CH is compromised by the adversary, the adversary can learn the sensed data and forge the data because the adversary knows key material even though security primitive such as encryption is used to provide confidentiality of the sensed data. This feature makes the application in WSNs very more critical and vulnerable.

2.2 Secure Routing Protocols in WSNs

A number of secure routing protocols for sensor clustering in hostile situations have been proposed recently in WSNs [15,9,10,14]. These can be generally categorized into a secure cluster formation protocol and a secure CH election protocol.

F-LEACH [1] is proposed to protect the CH election in LEACH. A sensor itself elects as a CH by using common keys shared with the BS and the BS authenticates the CH declaration using the same keys. Then, the BS securely broadcasts array of the authenticated CHs using μ TESLA [2]. Sensors join only one authenticated CH. However, this scheme cannot authenticate the sensors which join the service of CH. To resolve this problem, Oliveiral et al. have proposed SecLEACH [9] in which the BS authenticates the CH declaration from sensors and CHs also authenticate the joining sensors.

However, they have some limitations. Firstly, they do not consider the security of cluster leader election since these focus on the security issues during the formation of initial clusters. Secondly, a malicious insider can be a CH because these do not control the frequency of CH declaration. That is, they can prevent only external attackers not internal attackers from declaring themselves as CHs. Thirdly, the limitation related

with a compromised CH is not to guarantee confidentiality of the sensed data from sensor nodes. Even though each sensor node encrypts its sensed data for confidentiality, the compromised CHs can decrypt the data and finally they know which information is sensed from sensor nodes. So, a concealed data aggregation is required against a node capture. Finally, they have a single point of failure problem because they are the centralized mechanism which fully depends on the BS for authenticating the CH.

A number of secure cluster leader election protocols have been proposed in [5][10][14]. All of them work in a decentralized way and use the light-weighted cryptographic algorithms. The main technique of their schemes is that cluster member securely elects its CH through the agreement of a elected CH within a pre-defined cluster. [5] and [10] use a commitment based scheme. Each sensor sends its commitment to other sensors in the peer-to-peer manner. A commitment is an encrypted random value using a shared key and the random value is created by each sensor. Then, each sensor sends the fulfillment value(that is, its random value) to other sensors within its sector. The sensors receiving the fulfillment values verify them using the shared key and sum them to make an agreed random value. They divide the real sum of agreed random values by the number of sensors and get the remainder which indicates the position of CH node in the cluster. Like [5][10], Dong’s scheme [14] also has the assumption that cluster sector is pre-defined and not changed. In Dong’s scheme, the new CH is selected based on the remaining energy on cluster members for the purpose of the load balancing. Their idea for a resilient clustering is to let all member within a cluster have the identical list of the candidate CHs. To control the list of each sensor node, they use two one-way hash chains, YES chain and NO chain. In initialization phase, every sensor node sends its YES commitment value to the other nodes within its sector. Each sensor node receiving all YES commitment from the other nodes creates the list of the candidate CHs through the shuffle algorithm. The shuffle algorithm makes the list of the candidate CHs of all sensor nodes within the same sector identical. In each round of CH election, each sensor node sends YES key or NO key according to its residual energy. If sensor node does not have energy for serving as the CH, it releases NO key, otherwise it releases YES key. If sensor node receiving NO key, it removes identifier(ID) of the sending node from the list of the candidate CHs. Through this mechanism, all cluster members in the

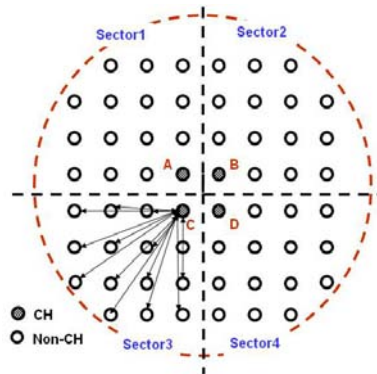


Fig. 1. The problem of pre-defined clustering

list of the candidate CHs will server as the CH in a round-robin manner unless they release their NO keys during the CH election. When the benign cluster members are well-connected, all benign nodes that are eligible for the CH role will be included in the list of the candidate CHs.

Resilient cluster leader election schemes are more resilient against a node capture than the above mentioned secure clustering protocol. However, they do not support a dynamic clustering which is a heuristic method for an optimal clustering since they have the assumption that cluster sector is pre-defined and not changed. Figure 1 shows the problem of a pre-defined clustering. In Fig. 1 WSNs are divided into 4 sectors. Nodes A,B,C and D simultaneously can become the CHs in any round. In this case, even though they are very close to each other, they are not clustered together because they belong to different sectors. As a result, intra-cluster communications of them cause the heavy overhead. This feature is not suitable for WSNs.

3 System Models and The Objective

3.1 Threats Model

We assume an adversary may launch arbitrary attacks against the cluster formation protocol and cluster based routing protocol except for completely jamming the communication channel. In this section, we classify adversaries in clustering routing protocol as follows: external adversary and internal adversary.

- **External adversary.** An external adversary can randomly select some nodes and extract key generation information from them. Then, he tries to eavesdrop the communication between other ordinary nodes. Besides, he tries to inject and replay packet to disrupt the cluster formation and cluster routing.
- **Internal adversary.** We assume that he already has knowledge on the algorithms or protocols. Using this knowledge, he can selectively capture the nodes and then he tries to eavesdrop the communication between other ordinary nodes. Then, he can use any method to interrupt/corrupt cluster formation and to forge the sensed data or the aggregated data, e.g., compromising some nodes, modifying the communication message, replay attack etc.

Ensuring a good level of security for such types of networks is not a trivial task. As WSNs use wireless communications, the threats and attacks against them are more diverse and often large-scale.

3.2 The Sensor Network Model and Notation

We consider a sensor network composed of N resource-constrained sensor nodes distributed uniformly on the square area of size $A = a \times a$ in a homogenous spatial Poisson process with density $\lambda = \frac{N}{A}$. For a cluster formation, each sensor node must embed DEACH algorithm. For DEACH, the optimal probability of being CH, P_{opt} , is pre-defined, which is CH_{opt}/N where CH_{opt} is the optimal number of CH. The sensor network model of the proposed scheme has the following assumptions.

- The sensor nodes are stationary, identical, and energy-constrained.
- All sensor nodes can communicate directly with the BS and are equipped with power control capabilities to vary their transmitted power.
- The BS has no energy constraint and is stationarily located far from sensor field.
- The communication is symmetric and based on a single-hop.
- A sensor node can compute the approximate distance based on the received signal strength.

We denote some common notations for describing the our scheme as follows.

- S_i is i^{th} ($i \in N$) sensor node.
- $S_i(id)$ is the identifier of the sensor node, S_i .
- $S_i(r)$ is the maximum transmission range of the sensor node, S_i .
- P_i is the probability of becoming a CH of the sensor node, S_i .
- $P_{opt}(i)$ is the optimal probability of becoming a CH of the node S_i referred in Eq. [11](#).
- CH_i is i^{th} cluster-head.
- $S_{CH}(i)$ is the set of the candidate CHs of the sensor node, S_i .
- $E_{res}(i)$ is the residual energy of the sensor node, S_i .
- E_{th} is the minimum energy required for serving as CH.
- $S_{neighbor}(i)$ is the set of neighborhood nodes of the node, S_i , within $S_i(r)$.
- $Q_i(j)$ is the qualification information of the node, S_j , as a CH of the node, S_i .
- $S_q(i)$ is the set of the qualified nodes of the node, S_i .
- d_i is the sensed data from the sensor node, S_i .
- er_{max} is the maximum number of rounds for CH election.

3.3 The Objectives

The objective of SecDEACH is to provide for both a secure clustering protocol considering a concealed data aggregation and a resilient CH election against an internal attacker and an external attacker. To prevent from security threats, WSNs must satisfy the security requirements and design requirement as following:

Security Goal 1: Resilience against fabrication attacks, impersonation attacks and replay attacks. To prevent an external attacker, a clustering protocol ensures that unauthorized sensor nodes should not join a cluster routing protocol. For doing this, it should guarantee the basic security primitives such as authentication, confidentiality, freshness and integrity.

Security Goal 2: Fault-tolerance. It is impossible to detect a compromised node acting like a ordinary node. So, if a compromised node acts as a CH, the application is damaged since a CH is responsible for data aggregation and data transmission schedule of its member nodes. To prevent this attack, a cluster routing scheme guarantees that a compromised node does not become a CH [\[14\]](#). That is, the attacker cannot arbitrarily increase or decrease the chance of any ordinary node being elected as a CH.

Security Goal 3: Concealed data aggregation in-networking. If a CH is compromised, the compromised CH node can learn all of the sensed data from its member nodes by decrypting them, even though its member nodes encrypt their sensed data. To

prevent this attack, a clustering routing protocol guarantees confidentiality and privacy of the sensed data which does not depend on CH.

Design Goal: Low communication, computation overheads and low storage requirement. Adding to the security primitives in a clustering protocol makes communication and computation overheads larger. The stronger robustness of security is, the overhead is larger. So, we consider the energy-efficient mechanism preserving the robustness of security. In addition, we consider the limitation of storage resulted from the resource constraints of the sensor node.

4 SecDEACH

In this section, we will briefly review the security primitives used in the proposed scheme and then describe its technical detail.

4.1 The Security Primitives

Our protocol takes advantage of the following three security primitives, the *one-way key chain*, the *Blundo's pairwise key pre-distribution scheme* and the *symmetric double additive homomorphic encryption*.

One-way key chain [14]. An one-way key chain $\{K_0, K_1, \dots, K_R\}$ is generated by iteratively performing the one-way hash function $H(\cdot)$ on the last key K_R in the chain. This key chain has the following property: given any K_j , all former keys can be derived by computing $K_i = H^{j-i}(K_j)$, $0 \leq i < j$, while none of the later keys can be computed due to the one-wayness of function H . Therefore, with the knowledge of $K_0 = H(K_1)$ (called the key chain commitment), anybody can verify the authenticity of any later key by only performing hash operations. The reason is that the attacker cannot change any hash value because the hash function in use is collision resistant. In our protocol, the one-way key chain will be used to protect the announcement reported by the self-elected CHs.

Blundo's pairwise key pre-distribution scheme [4,14]. Blundo [4] proposed a polynomial-based key pre-distribution scheme that was aimed to pre-distribute group key. In this scheme, the BS generates a t -degree bivariate polynomial $f(x, y)$ that has a property $f(x, y) = f(y, x)$ over a finite field F_q . The q is a prime number and is enough large to generate cryptographic keys. For each node having a unique ID i , the BS distributes the polynomial key $f(x, i)$. Based on this polynomial key, any two nodes sharing the key can generate a pairwise security key. Let us have an example. When the node i wants to communicate with a node j , it generates a key $f(j, i)$ and the node j can also generate a key $f(i, j)$. Due to the property of $f(x, y) = f(y, x)$, the two keys are same and this key becomes a unique pairwise key between node i and node j . This scheme can tolerate up to t compromised nodes, i.e., if an adversary wants to obtain one polynomial key $f(x, y)$, he must capture at least $t + 1$ nodes having the polynomial key [4]. We do not consider the case where a large number of nodes are compromised since most applications will fail anyway in such situation [14]. Therefore, we always assume no more than t compromised nodes in the networks [14]. We will use this for resilient CH election in our proposed scheme.

Symmetric double additive homomorphic encryption scheme. [13] A homomorphic encryption scheme allows arithmetic to be performed on ciphertext. One example is a multiplicatively homomorphic scheme, whereby the multiplication of two ciphertexts followed by a decryption operation yields the same results which is the multiplication of the two corresponding plaintext values. Homomorphic encryption scheme is especially useful in scenarios where someone who does not have decryption keys needs to perform arithmetic operations on a set of ciphertext. A more formal description of additive homomorphic encryption scheme is as follows [13].

Let be $Enc: K \times P \rightarrow C$ and $Dec: K \times C \rightarrow P$, where Enc is a symmetric encryption function, Dec is the corresponding decryption function, P is the plaintext space, C is the ciphertext space, and K is a set of symmetric secret keys. Enc is additively homomorphic, if and only if for given messages $a_1, a_2, \dots, a_n \in P$ and for secret keys $k_1, k_2, \dots, k_n \in K$, there exists a key $k \in K$ such that the following equation holds

$$(a_1 + a_2 + \dots + a_n) = Dec_k(Enc_{k_1}(a_1) + Enc_{k_2}(a_2) + \dots + Enc_{k_n}(a_n)) \quad (2)$$

A double homomorphic encryption scheme is an encryption that is homomorphic in both the key and the plaintext. That is, it holds that

$$(a_1 + a_2 + \dots + a_n) = Dec_{k_1+k_2+\dots+k_n}(Enc_{k_1}(a_1) + Enc_{k_2}(a_2) + \dots + Enc_{k_n}(a_n)) \quad (3)$$

In our proposed scheme, we will use this for concealed data aggregation.

4.2 Protocol Description

Initialization. Firstly, to use μ TESLA [2] for the BS's authenticated broadcast, the BS creates BROADCAST key chain by using one-way hash key chain method in subsection 4.1, $\{K^0, K^1, \dots, K^{br_{max}}\}$, where br_{max} is the maximum number of broadcasts from the BS. The BS creates a pairwise symmetric keys, $\{K_1, K_2, \dots, K_N\}$, for all sensor nodes to protect the node-to-BS communication. The BS then maintains a symmetric bivariate polynomial $f(\cdot, \cdot)$ which is a key material for establishing a pairwise key between sensor nodes after deploying and pre-distributes a polynomial share $f(i, \cdot)$ to every node S_i . We will describe later how to generate sensor node id, $S_i(id)$.

Next, prior to deployment, each sensor node S_i is assigned with node ID, $S_i(id)$, two random one-way key chains, the A - key chain [1] and the G - key chain [2], a pairwise key shared with the BS, K_i , and K^0 which is shared by all members of the network for BS's authenticated broadcast. The G-key chain is used by this node to tell other nodes that it has no energy left for serving as CH. Therefore, it only includes a key chain commitment, $G_{i,0}$, and G-key, $G_{i,1}$. The A-key chain is used to inform other nodes that it is a self-elected CH. Therefore, it includes a key chain commitment, $A_{i,0}$, and a number of A-keys, $\{A_{i,1}, A_{i,2}, A_{i,3}, \dots, A_{i,\delta * P_{opt} * er_{max}}\}$. We will discuss in section 5.4 why we determine the number of A-keys as $\delta * P_{opt} * er_{max}$. Lastly, we define that the ID of sensor node is the hash of the two key chain commitments, i.e. $S_i(id) = H(A_{i,0} || G_{i,0})$ like in [14].

¹ ANNOUNCEMENT key chain.

² GIVEUP key chain.

After deployment, the sensor node S_i firstly establishes a pairwise key with neighborhood sensor nodes within their communication range $S_i(r)$ by using Blundo's pairwise key pre-distribution scheme and initializes their neighborhood sensor nodes list, $S_{neighbor}(i)$. A pairwise key between two sensor nodes S_i and S_j is denoted as $K_{i,j}$. Next, the every sensor node calculates its $P'_{opt}(i)$ for the CH election. Finally, the every sensor node needs to exchange their key chain commitments (A-key commitment, $A_{i,0}$ and G-key commitment, $G_{i,0}$) between each other within its transmission range. For doing so, each sensor node sends $initMsg(S_i) = \{S_i(id) \| E_{k_{i,j}}(A_{i,0} \| G_{i,0} \| P'_{opt}(i))\}$, within its transmission range. The sensor node receiving this message can easily verify these commitments according to the node ID. Whenever a sensor node S_i receives the verified A-key commitment from a node S_j , S_i creates $Q_i(j)$. $Q_i(j)$ means that S_j has a qualification being a CH of S_i . It consists of node ID, $1/P'_{opt}(j)$ rounds, the last round that the node S_j was a CH, the number of announcements within last $1/P'_{opt}(j)$ rounds and last A-key. After collecting all those authenticated A-chain commitments from other nodes, every sensor node S_i creates a qualification list, $S_q(i)$. $S_q(i)$ consists of array of $Q_i(j)$. This is meaningful because a dynamic clustering formation is only performed within its communication range.

Secure Cluster Head Election. We use DEACH in section 2.1 as a fundamental clustering method. DEACH consists of a cluster-setup phase and a data transmission phase. The secure CH election and secure cluster formation phases belong to the cluster-setup phase. For simplicity, we assume that no collision occurred in the MAC(Media Access Control) layers of sensors during the clustering routing.

At the beginning of this phase, if the residual energy $E_{res}(i)$ of the sensor node S_i is less than the threshold E_{th} , the sensor node S_i broadcasts the key of its G-key chain to announce its decision. E_{th} is the minimum energy required for serving as CH and pre-determined system value. If $E_{res}(i)$ is greater than E_{th} , this node chooses a random number between 0 and 1. If the random number is less than its P_i in Eq. 11 this node becomes the CH and releases its A-key from A-key chain and broadcasts its announcement message, $advHeadMsg(CH_i) = \{S_i(id) \| C_{i,c}\}$, β times to tolerate the channel loss within transmission range. c is the number of the cumulative announcements sent by S_i . If a sensor node S_j receives a key in the A-key chain of S_i , node S_j can verify it by computing $A_{i,c-1} = H(A_{i,c})$ and comparing with the key $A_{i,c-1}$ received earlier. If the verification succeeds, node S_j extracts $Q_j(i)$ from its $S_q(j)$. And then S_j determines whether S_i has a qualification being a CH by using Eq. 4.

$$V(i) = \begin{cases} 1 & , \text{if } (r - r') \geq r \bmod (1/P'_{opt}(i)) \\ 0 & , \text{otherwise} \end{cases} \quad (4)$$

where r is the current round, r' is the last round in $Q_j(i)$. If $V(i)$ is one, S_j adds S_i into $S_{CH}(i)$ since it means that S_i is a non-compromised node. Otherwise, this announcement is ignored. If this announcement continues over α times within $1/P'_{opt}(i)$, S_i is removed permanently from $S_q(j)$. If node S_j receives the correct G-key $G_{i,1}$ from S_i , S_j will remove S_i from $S_q(j)$. Finally, after $S_{CH}(j)$ is created, S_j chooses the closest CH_i among the CHs in $S_{CH}(j)$ and joins to the CH_i as the member. After joining, it updates the A-key and the last round in $Q_j(i)$. Table 1 shows the CH election process in a sensor node having 5-node neighbors within its transmission range when $\alpha = 3$.

Table 1. Illustration of the resilient CH election in a sensor node having 5-node neighbors($\alpha = 3$)

	Released Information					$S_i(Q)$ =(list of ($\text{id}^a, \underline{r}^b, \underline{p}^c$))	$S_i(CH)$	CH Node	Notes
	1	2	3	4	5				
initial Stage	$A_{1,0}$ $G_{1,0}$	$A_{2,0}$ $G_{2,0}$	$A_{3,0}$ $G_{3,0}$	$A_{4,0}$ $G_{4,0}$	$A_{5,0}$ $G_{5,0}$				Exchange Commit. and generate $S_i(Q)$
Round 1		$A_{2,1}$		$A_{4,1}$		(1,0,3),(2,1,4),(3,0,6) (4,0,7),(5,0,4)	2,4	2 or 4	select 2
Round 2			$A_{3,1}$	$G_{4,1}$		(1,0,3),(2,1,4),(3,2,6) ,(5,0,4)	3	3	remove 4 , α for 2 =1
Round 3		$A_{2,2}$			$A_{5,1}$	(1,0,3),(2,1,4),(3,2,6) ,(5,3,4)	5	5	select 5 , α for 2 =2
Round 4	$A_{1,1}$	$A_{2,3}$	$A_{3,2}$		$A_{5,2}$	(1,4,3),(3,2,6),(5,3,4)	1	1	select 1 and remove 2 because α for 2 =3

^a $S_i(id)$.

^b Last round which it was a CH.

^c $1/P'_{opt}(i)$.

Secure Cluster Formation. After electing the CH_j , S_i sends the $joinMsg$ to the CH_j . The $joinMsg(S_i, CH_j)$ is as follows.

$$joinMsg(S_i, CH_j) = \{S_i(id) \| S_j(id) \| nonce \| MAC_{k_{i,j}}(S_i(id) \| S_j(id) \| nonce)\}$$

The $nonce$ is used to prevent the replay attack. After receiving $joinMsg$, the CHs send the time slot schedule to their cluster member nodes by using a pairwise key.

Secure Data Aggregation. In this step, S_i encrypts the sensed data using a shared key K_i with the BS for the concealed data aggregation of its CH and sends $reportMsg$. The $reportMsg(S_i, CH_j)$ is as follows.

$$reportMsg(S_i, CH_j) = \{S_i(id) \| S_j(id) \| Enc_{k_i}(d_i) \| nonce + l \| MAC_{k_{i,j}}(S_i(id) \| S_j(id) \| Enc_{k_i}(d_i) \| nonce + l)\}$$

The node-to-CH communication is protected by using the same key used to protect the $joinMsg$ and guarantees integrity through MAC³. A value computed from the $nonce$ and the reporting cycle(l) [9] is also included to prevent replay. The CHs can now check the authenticity of sensing reports they receive and then perform a concealed data aggregation $d_{sum} = \sum Enc_{k_i}(d_i)$ for all encrypted sensing data using symmetric double homomorphic encryption scheme in Eq. [3]. That is, not only the CHs do not need to decrypt the encrypted sensing data but also can not the CH decrypt the encrypted data since the encryption is performed by using the shared key with the BS. After performing the concealed data aggregation, the CH sends $reportMsg$ to the BS. The $reportMsg(CH_i, BS)$ is as follows where $L_{members}$ is the IDs list of reporting member nodes of CH_i . It is required for the BS to decrypt the sum of encrypted data by the CH_i .

$$reportMsg(CH_i, BS) = \{S_i(id) \| L_{members}(i) \| d_{sum} \| counter \| MAC_{k_i}(S_i(id) \| L_{members}(i) \| d_{sum} \| counter)\}$$

³ Message authentication code as keyed hash function.

The aggregate result is protected by using the symmetric key shared between the CH and the BS. For freshness, a *counter* (shared between the CH and the BS) is included in the reporting message. The BS receiving the aggregated data can decrypt the data using shared keys between sensor nodes using Eq. 2.

5 Analysis

5.1 Security Goal 1: Resilience against Fabrication Attacks, Impersonation Attacks and Replay Attacks

This security goal requires authenticity, integrity, confidentiality and freshness to node-to-node communications.

For authenticity, it has to guarantee that only legitimate cluster can participate in cluster routing operation. In other words, we have to ensure that the unauthorized nodes do not have to join the CH election and cannot impersonate any ordinary sensor node. These two properties are actually achieved by using the Blundo's key pre-distribution protocol [4]. The ID of a sensor node is tied to its cryptographic key. It is thus not possible for an unauthorized node to join the cluster routing procedure. That is, the messages, *initMsg*, *joinMsg* and *reportMsg*, are encrypted with a pairwise key and a successful decryption of this message allows the sensor node to conclude that the message is originated from a legitimate node in the network. In addition, even though *advHeadMsg* is not encrypted, the adversary cannot forge the announcements from a legitimate node. The announcement contains the ID of the sensor node and the key in A-key or G-key chains. To forge these keys, the attacker needs to either invert the one-way hash function H or fool benign cluster members into accepting incorrect key chain commitment. However, because of our ID assignment scheme, the ID of a sensor node, $S_i(id)$, is the hash of its two key chain commitments, i.e., $S_i(id) = H(A_{i,0} || G_{i,0})$. Therefore, to forge the key chain commitment, the adversary has to forge the ID as well. This requires the knowledge of keying materials related to the forged ID, which is not available to the adversary. Therefore, the adversary has no way to make its forged commitments accepted by any benign cluster member. Without the proper chained keys, the adversary may try to replay the keys released earlier by the legitimate member. However, this will not impact our scheme because our scheme uses either the fresh A-key or the G-key in each round of CH election.

Because the messages, *joinMsg* and *reportMsg*, are encrypted by MAC using a pairwise key and includes the *nonce*, the CH can also conclude that it is not changed and it is not a stale message being replayed. The freshness of all subsequent sensor reports from the ordinary nodes to their CH is guaranteed by *nonce* values that are incremented each time. For CH-to-BS communication, the freshness is guaranteed by the *counter* value shared between the CH and the BS. The *counter* value is also incremented each time the CH sends a new report to the BS.

5.2 Security Goal 2: Fault-Tolerance

A compromised node can always behave normally like benign node and be elected as the new CH at some point during the field operation. There are no effective ways

to identify those passive malicious node since there are no evidences of them being malicious. As a result, this goal focuses on making sure that an adversary cannot significantly impact the chance of any benign members elected as a new CH. Thus even if the CH is compromised, it cannot continue to control properly the cluster. The CH election protocol should ensure the fairness of the sensor nodes being elected as CHs. For doing so, our idea is to ensure that each sensor node becomes a CH one per $1/P'_{opt}$ rounds on average.

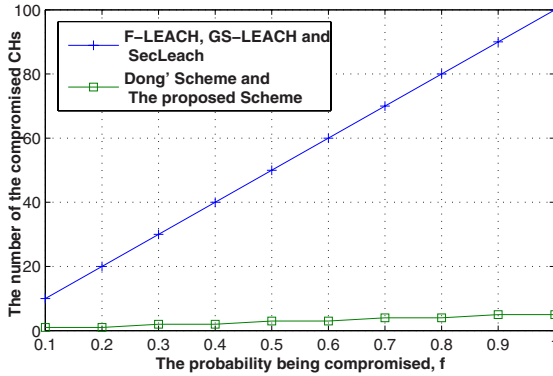


Fig. 2. The number of compromised CHs

In our scheme, the CH election is based on the neighbors list $S_{neighbor}(i)$, which is initialized at the beginning of the CH election. All nodes in the neighbors list has the qualification serving as the CH. In our scheme, the sensor node can receive several announcements being a CH from the neighborhood in each round of CH election. The sensor node receiving them has to verify whether they have a qualification serving as the CH and chooses its CH among the verified candidate CH nodes. To verify the qualification, our scheme uses Eq. 4. If the announcement is advertised more than α times within $1/P'_{opt}$, the verification is failed. If the verification is failed for a node, the node cannot become a CH. It ensures that each sensor node becomes a CH one per its $1/P'_{opt}$ rounds on average. Figure 2 shows the comparison our scheme with the existing secure cluster routing schemes [19] in the aspect of the compromised CH election. Figure 2 shows the number of the compromised node being elected CH according to the probability being compromised of the sensor node, f , where $EX(1/P'_{opt}) = 5$ and $N = 100$. The number of compromised CH in SLEACH, F-LEACH and SecLEACH is very larger than that of our scheme and Dong's scheme because the secure clustering schemes like F-LEACH cannot know if the a self-selected CH is compromised. As we can see in Fig. 2 our scheme provides resilient CH election like Dong's scheme [14].

5.3 Security Goal 3: Concealed Data Aggregation In-Networking

For the concealed data aggregation, our scheme ensures that only BS can decrypt the aggregated data from all sensor node. This property is actually achieved by using symmetric double additive homomorphic encryption [13]. Each sensor node encrypts the

sensed data with a shared key with only BS. After receiving all *reportMsg* from all member node in its cluster, the CH just only sum all the encrypted data and then sends the summed data to the BS. The BS can extract the sensed data from the summed encrypted data since the BS knows all of the shared keys with sensor nodes. It offers a concealed data aggregation as well as data privacy.

5.4 Overhead

Before deployment, every sensor node is assigned with an ID and a t -degree polynomial share $f(i, \cdot)$. It also needs to store two key chains along with their commitments, one contains a single key, and the other consists of $\delta * P_{opt} * er_{max}$ keys. These chains are similarly used in Dong's scheme [14]. However, in Dong's scheme, if a node has capability being a CH, sensor node must inform its YES key per every cluster leader election round since a YES chain is used to announce being a candidate CH. So, the number of keys of a YES chain in Dong's scheme is the same as the maximum number of rounds for CH election, er_{max} . On the other hands, in our scheme, the sensor node being a self-elected CH only broadcasts its A-key to announce the fact that it becomes a CH. Without loss of generality, the ordinary sensor node is elected itself as CH one within last $1/p'_{opt}$ rounds on average because the sensor node elects itself as the CH by Eq. 1 which is optimized into P_{opt} [16]. Thus, in our scheme, each node S_i only needs $P'_{opt}(i) * er_{max}$ keys for A-key chain unlike Dong's scheme. So, the number of our A-key chain is very smaller than that of Dong's scheme. However, there is one limitation to consider. To calculate P'_{opt} , each node has to know its location information which will be deployed. However, we cannot calculate $P'_{opt}(i)$ in the initialization phase since each sensor node will be deployed randomly. So, we determine the length of A-key chain as $\delta * P_{opt} * er_{max}$ based on P_{opt} . According to Eq. 1, $P'_{opt}(i)$ is less than or equal to $\frac{3}{2}P_{opt}$. As a result, the adequate value of δ can be 2.

Figure 3 shows the comparison our scheme with Dong's scheme in the aspects of the memory and communication overheads. Figure 3(a) and Fig. 3(b) show the memory overhead for saving A-key chain and the communication overhead for sending announcement by the number of CH where *key size* = 8bytes, $er_{max} = 50$ and $N = 100$. As we can see in Fig. 3, our scheme is more light-weighted and energy-efficient than Dong's scheme. In addition, if there are n sensor nodes within a transmission range, a sensor node will also need to buffer the neighbors list of n IDs, $(n - 1)$ pairwise keys, and $2(n - 1)$ key chain commitments. These space requirements are usually not a problem since the transmission range is limited.

The computational overhead mainly comes from three parts. First, with the Blundo's key pre-distribution protocol, a sensor node needs to compute a t -degree polynomial to establish a pairwise key and verify the ID of another sensor node. Second, a sensor node has to perform the hash function H one time to verify the key chain commitment. Finally, every sensor node will need to perform symmetric key operations such as encryption and authentication to protect the communication between sensor nodes. Obviously, all these operations only involve light-weighted computation.

The proposed scheme only causes small communication overhead like Dong's scheme [14]. The Blundo's key pre-distribution technique does not introduce any additional communication overhead at all for two nodes to establish a pairwise key [14]. In the

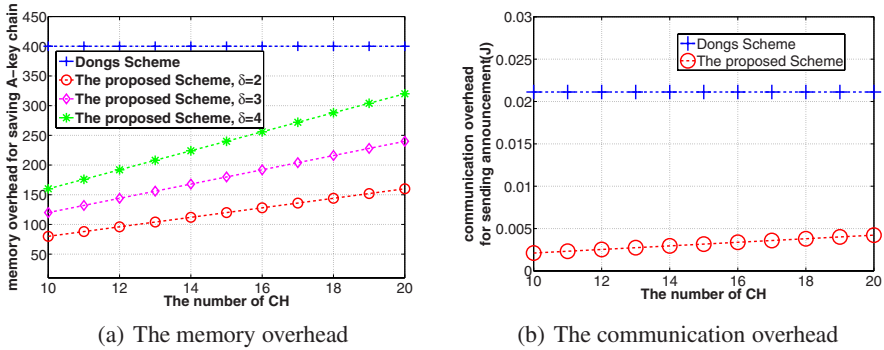


Fig. 3. The comparison of the overhead

initialization step, each node only needs to disclose its ID and two key chain commitments. In each round of CH election, every sensor node is only required to broadcast a single chained key, A-key. Different from centralized approaches [119] where sensor nodes exchange messages with the BS that is far away, our scheme only involves the communication between the sensor nodes that are physically close to each other. Therefore, the communication cost of our protocol is not a big problem [114].

6 Conclusion and Future Works

In this paper, we have presented a secure and resilient clustering protocol preserving data privacy. Our scheme provides both a secure clustering and a resilient CH election preserving dynamic clustering unlike the existing secure clustering schemes. In addition, our scheme provides data privacy through the concealed data aggregation. The analysis shows that our scheme offers the security services against both inside attack and outside attack while having the little memory, computation and communication overhead. The comparison with related works shows that our scheme has the improvement in security and the reduction of the overheads compared to the existing scheme.

A further direction of this study will perform experiment on using our proposed schemes and perform the network wide security analysis. In addition, we will apply public key material into our scheme to support more strong authentication and electronic signature.

References

1. Ferreira, A.C., Vilaca, M.A., et al.: On the security of cluster-based communication protocols for wireless sensor networks. In: Lorenz, P., Dini, P. (eds.) ICN 2005. LNCS, vol. 3420, pp. 449–458. Springer, Heidelberg (2005)
2. Perring, A., Szewczyk, R., Well, V., Culler, D., Tygar, J.D.: SPINS: Security Protocol for Sensor Networks. *Wireless Networks* 8(5), 521–534 (2002)
3. Karlof, C., Wagner, D., et al.: Secure routing in wireless sensor networks: attacks and countermeasures. *Elsevier’s AdHoc Networks J.* 1(2-3), 293–315 (2003)

4. Liu, D., Ning, P.: Establishing pairwise keys in distributed sensor networks. In: CCS 2003, Proceeding of the 10th ACM conference on Computer and communications security, pp. 52–61. ACM Press, New York (2003), Network
5. Wang, G., Cho, G.: Secure Cluster Head Sensor Elections using Signal Strength Estimation and Ordered Transmissions, *Sensors*. MDPI (Molecular Diversity Preservation International) 9(2), 4709–4727 (2009)
6. Karl, H., Willig, A.: *Protocols and Architectures for Wireless sensor Networks*. John Wiley & Sons, Ltd., Chichester (2005)
7. Al-Karaki, J.N., Kamal, A.E.: Routing techniques in wireless sensor networks: a survey. *IEEE Wireless Communications* 11(6), 6–28 (2004)
8. Akyildiz, L.F., et al.: A Survey on Sensor Networks. *IEEE Communications Magazine* 40(8), 102–114 (2002)
9. Oliveira, L.B., Ferreira, A., et al.: Secleach-on the security of clustered sensor networks. *Signal Process* 87(12), 2882–2895 (2007)
10. Sirivianos, M., et al.: Non-manipulable aggregator node election for wireless sensor networks. In: ICST wiOpt (January 2007)
11. Younis, O., Fahmy, S.: HEED: A Hybrid, Energy-Efficient, Distributed Clustering Approach for Ad Hoc Sensor Networks. *IEEE Transactions on Mobile Computing* 3(4), 660–669 (2004)
12. Younis, O., et al.: Node Clustering in Wireless Snsor Networks: Recent Developments and Deployment Challenges. *IEEE Network* 20(3), 20–25 (2006)
13. Tiwari, P.: *Data Aggregation in Cluster-based Wireless Sensor Networks*, A master thesis, Deemed University (July 2008)
14. Dong, Q., Liu, D.: Resilient Cluster Leader Election for Wireless Sensor Networks. In: Proceedings of The IEEE Communications Society Conference on Sensor, Mesh and Ad Hoc Communications and Networks (SECON), June 2009, pp. 1–9 (2009)
15. heinzelman, W., et al.: An application-Specific Protocol Architecture for Wireless Microsensor Networks. *IEEE Transactions on Wireless Communications* 1(4) (October 2002)
16. Han, Y.-J., Park, S.-H., Eom, J.-H., Chung, T.-M.: Energy-Efficient Distance based Clustering Routing Scheme for Wireless Sensor Networks. In: Gervasi, O., Gavrilova, M.L. (eds.) ICCSA 2007, Part II. LNCS, vol. 4706, pp. 195–206. Springer, Heidelberg (2007)

Avoidance of Co-channel Interference Using Switched Parasitic Array Antenna in Femtocell Networks

Yeonjune Jeong, Hyunduk Kim, Byung-Sung Kim, and Hyunseung Choo*

School of Information and Communication Engineering
Sungkyunkwan University, Korea
{yjjeong83, yuby82, choo}@skku.edu, bskim@ece.skku.ac.kr

Abstract. Femtocells are low-cost and low-power cellular home base stations that connect mobile users to a operator's network by means of a broadband backhaul. Through this approach, the network operators can extend indoor coverage at a low-cost and reduce the operating cost. However, there exists a co-channel interference between macrocell and femtocells when they use the same frequency band. This co-channel interference causes severe performance degradation of femtocells and macrocell. Therefore, we propose a novel scheme to reduce the co-channel interference to macrocell without performance degradation of femtocells. The proposed scheme avoids the co-channel interference to macrocell users effectively using the SPA antenna which offers 8-directions to control the antenna beam pattern and transmitting power. Performance evaluation results have confirmed that the proposed scheme reduces the co-channel interference by 80% compared to the case using omni-directional antenna.

Keywords: We would like to encourage you to list your keywords within the abstract section.

1 Introduction

According to recent surveys, 50 percent of phone calls and 70 percent of data services have taken place indoors [1]. Therefore, network operators need to provide high data rate and quality of service (QoS) indoors. Femtocell has been introduced to solve indoor coverage problem and QoS requirement. One definition of femtocell is a low-power wireless access point that operates in licensed spectrum to connect standard mobile devices to the operators network using a digital subscriber line (DSL) connection [2]. Femtocell provides benefits for both subscribers and cellular poerators. Subscribers can achieve high data rate indoors with low-cost with deployment of femtocell. Wireless operator can extend the indoor coverage with low-cost, increase cell capacity, and gain the high frequency reuse efficiency. Moreover, cullular operators can reduce traffic for macrocell, thus reducing infrastructure and operating cost.

* Corresponding author.

One of the important issues of introducing femtocells is solving the interference problem with macrocell [3-7]. If macrocell and femtocells use the same frequency in orthogonal frequency division multiple access (OFDMA) system, there is a co-channel interference between macrocell and femtocells. This interference causes performance degradation of macrocell and femtocell. Cell capacity decreases as the number of femtocells increases [3]. Since the femtocell is deployed by its users, it is difficult for cellular operators to estimate the location of femtocells. Therefore, cellular operator cannot solve the interference problem between macrocell and femtocells with existing network planning and optimization methods.

To solve the above problem, a method to control femtocell coverage has been proposed to reduce the interference between macrocell and femtocells. It defines the area where the signal strength of femtocell is equal to that of macrocell as interference-limited coverage area (ILCA) and control the ILCA to reduce the amount of interference [8]. However, it reduces the femtocell coverage as the distance between femtocell base station (FBS) and macrocell base station (MBS) is getting closer. Another method has been introduced which controls the transmitting power and steers the beam pattern of the femtocell to reduce the interference to macrocell [9]. In this scheme, femtocell adjusts its transmitting power equal to the macrocell power at the border of femtocell considering the path loss of the macrocell and femtocell. However, it needs complicated algorithm to reduce the interference between macrocell and femtocell.

In this paper, we propose a new physical approach adopting dynamic antenna beam planning to avoid the interference using a switched parasitic array (SPA) antenna. The SPA is an antenna capable of steering its main beam direction dynamically [10-12]. Compared to the multi-element antennas used in [9], the SPA provides the evenly distributed directional antenna patterns, which makes it easier to estimate the receiving power and optimize the control algorithm. FBS controls the beam pattern of antenna to support its user only with high front-back ratio. Therefore, it can mitigate the interference to macrocell users effectively. Additionally, it can reduce the power consumption of FBS without any degradation of the performance of femtocell because SPA may achieve higher antenna gain compared to omni-directional antenna.

The paper is structured as follows. In section 2, previous approaches to solve the interference problem between macrocell and femtocells are discussed and the structure and characteristics of SPA antenna are described. Then we propose a method of data transmission and power control in femtocells using SPA in section 3. In section 4, the performance of the proposed method using SPA is evaluated compared to that of femtocell using omni-directional antenna. Finally, conclusions are given in section 5.

2 Related Work

If the macrocell and femtocell use the same frequency band, it causes co-channel interference between them. Figure 1 shows that the co-channel interference scenario between the macrocell and femtocell in downlink. The macrocell users are

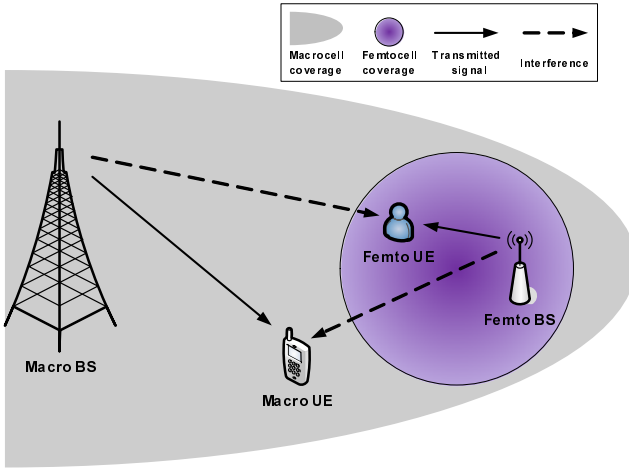


Fig. 1. Scenario of macro and femtocell Interference in downlink

interfered from FBS, and the femtocell users are interfered from MBS as shown in Figure 1. When a macrocell user approaches the femtocell coverage, the effect of interference from FBS may grow up, and when the femtocell user approaches MBS, it may suffer heavy interference from MBS. Such cases, the performance of macrocell and femtocell are decreased. If the macrocell user gets into the femtocell coverage, communication with MBS becomes impossible, because the interference to macrocell user in femtocell coverage is much bigger than received signal strength from MBS.

In the previous study, FBS uses omni-directional antenna to transmit the signal to all direction irrespective of its user’s position [7, 8]. Therefore, macrocell users are interfered from FBS when they are located nearby the femtocell coverage. Additionally, co-channel interference is getting bigger as transmitting power of FBS increases. If FBS makes beam pattern to its user’s location, the signal level of FBS is very low in the rest of its user’s location. If the macrocell user is located outside of femtocell coverage, the co-channel interference from FBS can be avoided based on macrocell user’s location.

In [9], a multi-element antenna solution is proposed to decrease the interference in macrocell users. It can control the antenna beam direction by switching of two patch antennas and IFA antennas. The effect of interference in macrocell users is reduced by controlling the transmission power P_{femto} of FBS according to the equation given below:

$$P_{femto} = \min (P_{macro} + G(\theta) - L_{macro}(d) + L_{femto}(r), P_{max}) \quad (1)$$

where P_{macro} is the transmitting power of MBS, $G(\theta)$ is the antenna gain of MBS, $L_{(macro)}(d)$ is average outdoor path loss when the distance between macrocell and femtocell is d . $L_{femto}(r)$ is average indoor path loss when the distance between FBS and its user is r . P_{max} is the maximum transmitting

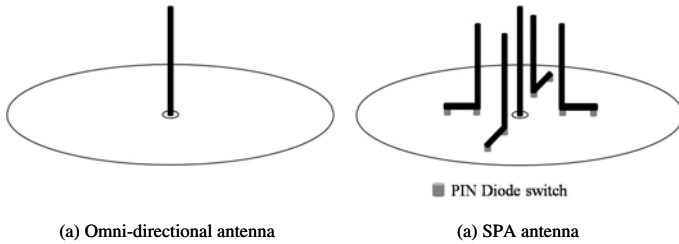


Fig. 2. Structure of antennas

power of FBS. The method of (1) means that FSB increases the transmitting power to overcome the interference by macrocell when it is closely located to MSB and conversely the FBS far away from MBS lowers its power.

Figure 2 shows omni-directional antenna and SPA antenna. In this paper, FBS uses the SPA antenna in Figure 2 (b). A monopole antenna located at the center of the ground plane, and four identical parasitic array elements are placed surrounding the monopole symmetrically. Switches are inserted between array elements and ground plane to control the main beam direction of the monopole antenna by alternating the ground connection of the parasitic elements. A bias applied to the knee in the element provides an RF shorting path to the ground plane, and the shorted parasitic elements acts as a “director” to guide the main beam direction to its direction. A bias applied to the end of element provides a long RF path to the ground plane when the switch at the knee is turned-off. In this case, the parasitic element works as a “reflector.”

If two switches in the parasitic element are turned-off, the element is “invisible.” Four near orthogonal beam patterns can be generated by sequentially switching two element pair placed in a diagonal position as front director and back reflector. In addition, eight directional patterns can be also obtained by utilizing two adjacent reflectors and two directors. All directional beam patterns are distinctly identified from each other according to their direction with the same shape and easily synthesized by controlling the switches. Even radiation pattern and reasonable front-back ratio for each direction makes it possible to simplify the power control algorithm to reduce the interference.

3 Proposed Scheme

In this section, we explain the procedure of data transmission to its user in FBS using SPA antenna. We assume that FBS are configured to limit access to only a few authorized users to protect limited resources. If unauthorized users request to connect to the closest femtocell, FBS restricts connection to femtocell for them. In this condition, our proposed scheme is as follows. The FBS monitors the user location by measuring the signal strength of surrounding users. Then, FBS performs 1-directional power control of transmitting signal with antenna gain based on user’s position. Then, check the users’ ID to avoid unwanted

data transmission to unauthorized users. Finally, FBS make the beam pattern to transfer the signal to authorized user. These processes guarantee the same performance of femtocell with FBS using omni-directional antenna. Additionally, it reduces the co-channel interference to macrocell users nearby femtocell coverage. The procedure of data transmission of FBS is shown in Figure 3.

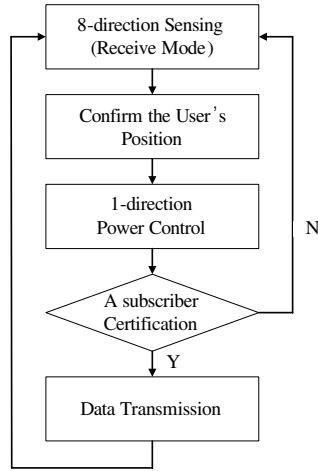


Fig. 3. Procedure of data transmission of FBS.

First, FBS measures the signal level of nearby users. Second, FBS confirms the user's location based on received signal strength. Third, according to the measured distance and location, FBS controls the transmitting power as below:

$$P_{femto} = \min(P_{omni} - G(\theta) + L_{femto}(r), P_{max}) \quad (2)$$

where P_{omni} is transmitting power when FBS uses omni-directional antenna. $G(\theta)$ is antenna gain of SPA in direction of the user, where θ is the angle to FBS with respect to the user's location. $L_{femto}(r)$ is the average path loss at the distance r between FBS and the user. Fourth, FBS checks the user ID. In this paper, assume that the FBS acts as closed access mode which only allows the wireless connection to registered users. Thus, FBS can reduce the unwanted data transmission for unregistered users. If unregistered user request the wireless connection to femtocell, FBS goes back to receive mode. By restricting wireless connection for unregistered users, femtocell can avoid the unnecessary data transmission, and reduces co-channel interference to macrocell user effectively.

It has advantages controlling the transmitting power of FBS based on its user's location. First, femtocell can reduce interference to adjacent macrocell users. For example, macrocell users get interfered from FBS without their locations when FBS uses omni-directional antenna. However, if our proposed scheme is applied, FBS makes the beam pattern to the certain limited area where its serving user

is located. If the angle of the location between macrocell user and femtocell user is quite different, macrocell user can avoid the co-channel interference from the FBS. Otherwise, macrocell user may get interfered from FBS as the difference of angle with femtocell user's location.

Second, by using SPA antenna in FBS, we can obtain higher antenna gain than FBS using omni-directional antenna. Thus, FBS using SPA antenna has the same performance with omni-directional antenna while transmitting power of FBS is lower than that of omni-directional antenna. Additionally it can avoid unnecessary data transmission due to transmit the signal only for registered users. When FBS uses SPA antenna, power consumption can be lower than that of omni-directional antenna.

Table 1. Simulation Parameters

Parameters	Values
System bandwidth	10MHz
Inter-macrocell distance	1000m
Inter-femtocell distance	40m
Radius of femtocell	10m
Macrocell BS power	20W
Femtocell BS power	20mW
Number of subcarriers	768
Frame duration	5ms
Number of OFDM symbols	48
Thermal noise(AWGN) density	-174dBm/Hz
System frequency(λ)	2GHz

4 Performance Evaluation

In this paper, we performe the system level simulation in three cases: conventional OFDMA macrocell with no femtocells; macrocell with femtocells using omni-directional antenna; and macrocell with femtocells using SPA antenna. We consider the reuse factor 1 for macrocell as 2-tier network that has 19 cells. We assume that MBS is located at the center (0m) of the macrocell. Then, macrocell is divided into two regions according to the distance from MBS. 0~409m is inner region, and 409~577m is outer region. There are 30 users in each macrocell. We set the radius of femtocell to 10m; each femtocell has one user. We assume that FBS is loacaed at the center (0m) of each femtocells. Totally, 409 Femtocells are distributed uniformly in a macrocell and the distance between adjacent FBS is 40m. The table 1 shows parameters used in performance evaluation.

The path loss model used in performance evaluation is as follows. First, we use the modified COST 231-Walfisch-Ikegami urban macro model to know path loss of macrocell user [13].

$$P_{out}[dB] = 31.81 + 40.5 \log_{10}(d) + \psi_{out} \quad (3)$$

where the $d[m]$ is distance between MBS and the its user. $\psi_{out}[dB]$ is outdoor shadowing modeled as log-normal distribution with zero mean, 8dB standard deviation.

Second, we use Modified COST 231-multiwall model to know path loss of femtocell user [14].

$$P_{in}[dB] = 37 + 32 \log_{10}(d) + L + \psi_{in} \tag{4}$$

where the $d[m]$ is distance between the FBS and the its user, $L[dB]$ is path loss by wall located in the edge of femtocell. We set to $L = 8dB$ in this simulaiton. $\psi_{in}[dB]$ is indoor shadowing modeled as log-normal distribution with zero mean, 4dB standard deviation.

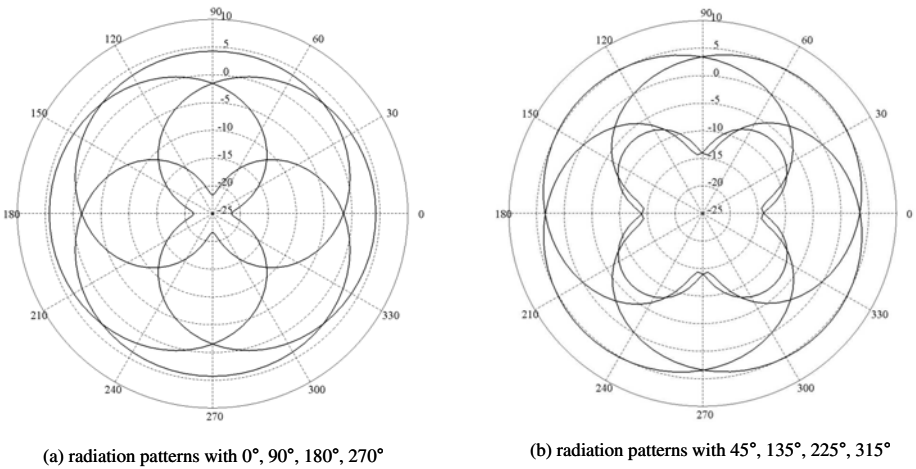


Fig. 4. The radiation patterns of SPA antenna

Third, we use indoor-to-outdoor model to know the effect of femtocell signal to macrocell user [15].

$$P_{in-to-out} = P_{out}(d_{out} + d_{in}) + 14 + 0.5d_{in} + \psi_{in} \tag{5}$$

where $P_{out}[dB]$ is outdoor path loss model in Eq. (3). d_{out} is the distance between macrocell user and the wall located at the edge of femtocell. d_{in} is distance between the FBS and the wall located in the edge of femtocell.

Figure 4 shows beam patterns of SPA antenna. The mesh size is $\lambda/10$ (λ is wavelength.), the simulation is executed by the method of moments (MoM) which is a full wave solution of Maxwell’s equation in the frequency domain. As we can see in Figure 4, SPA antenna can radiate in one direction which has about 10dB front-to-back ratio symmetrically [10-12]. SPA has symmetric radiation pattern for all directions and reasonable front-back ratio for each direction. Thus, by controlling beam pattern of FBS according to its user’s location, we can

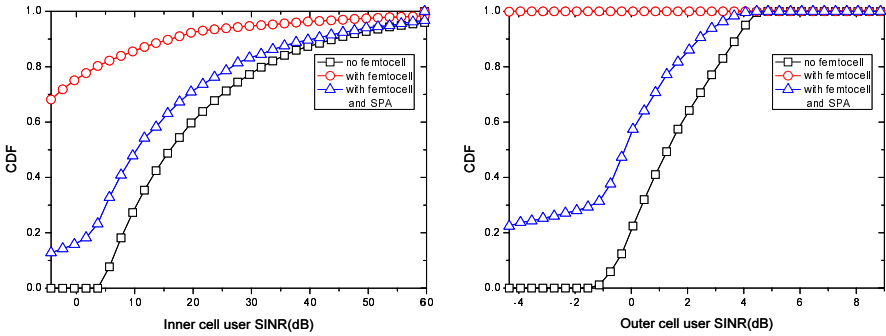


Fig. 5. CDF of average macrocell user SINR

decrease interference to macrocell user with same performance that of omnidirectional antenna.

Figure 5 shows that CDF of the macrocell user’s average SINR in the following three cases: conventional OFDMA macrocell with no femtocells; macrocell with femtocells using omni-directional antenna; and macrocell with femtocells using SPA antenna. First case has the highest user’s SINR with no doubt. Because, there are only inter-cell interference among the adjacent macrocells. In conventional OFDMA macrocell, outer users have low SINR because, outer users get severe interference from neighboring cells.

Second case shows that macrocell users have the lowest SINR among three cases. Because femtocell causes co-channel interference with macrocell, macrocell users get interfered from FBS using omni-directional antenna. It degrades the macrocell user’s SINR severely. In this simulation, we set the minimum required SINR to -4.34dB for macrocell user to communicate with MBS. In outer region, most users have SINR lower than -4dB . It denotes that most users cannot communicate with serving MBS because of the co-channel interference from FBS. The third case shows the macrocell user’s SINR when our proposed scheme is applied. Macrocell user can reduce the co-channel interference from FBS. Simulation results show that our proposed scheme improves the macrocell user’s SINR about 60% in inner region and 80% in outer region compared to the FBS using omni-directional antenna. It means FBS avoids the co-channel interference to adjacent macrocell user effectively because, FBS steers beam pattern to its user’s location.

Figure 6 shows the cell throughput of the above conditions. First case has the highest throughput as shown in Figure 5. Second case, conventional OFDMA macrocell with femtocell using omni-directional antenna, has much smaller throughput than the first case, because macrocell users are interfered from FBS severely. Especially macrocell users in outer region get intense interference not only inter-cell interference from adjacent macrocells, but also co-channel interference from nearby femtocells. Since SINR is lower than the minimum requirement, outer user’s throughput is close to zero. Third case shows the macrocell throughput when we adopt our proposed scheme. It has higher throughput than second

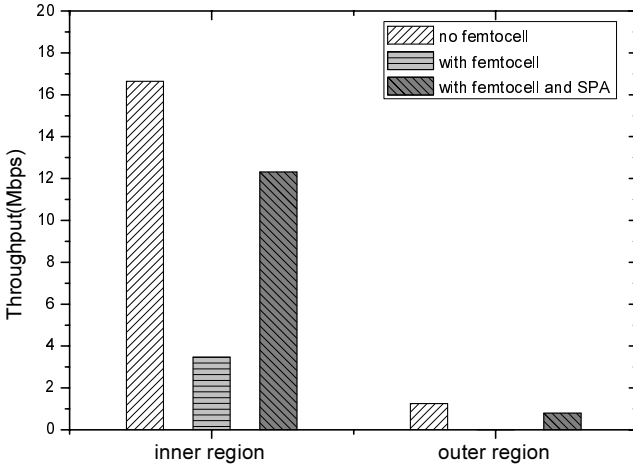


Fig. 6. Macrocell throughput

case, because there is less probability that macrocell users get interfered from FBS than the second case. Third case shows that our proposed scheme improves the throughput of macrocell users more than 4 times compared to the second case.

5 Conclusion

Femtocells are small cellular base stations which extend indoor coverage and provide high data rate. However, there exists co-channel interference between macrocell and femtocells when they are operating in the same frequency band. This co-channel interference causes severe performance degradation both macrocell and femtocells. Especially as FBS are using omni-directional antenna, they cause co-channel interference to nearby macrocell users regardless of their locations. And macrocell user's SINR gets lower as it gets closer to FBS. Therefore, in this paper we have proposed procedure of data transmission and power control for FBS using SPA antenna to reduce the co-channel interference with macrocell users. In performance evaluation, we confirmed that macrocell user's SINR is improved by 60 ~ 80% and macrocell throughput is increased more than four times compared to FBS using omni-directional antenna.

Acknowledgment

This research was supported by MKE, Korea under ITRC NIPA-2009-(C1090-0902-0046) and by MEST(Korea), under WCU Program supervised by the KOSEF (No.R31-2008-000-10062-0). Profs. B.Kim and H.Cho are corresponding authors.

References

1. Mansfield, G.: Femtocells in the US Market -Business Drivers and Consumer Propositions. FemtoCells Europe, ATT, London (2008)
2. Femto Forum, <http://www.femtoforum.org/femto/>
3. Chandrasekhar, V., Andrews, J., Gatherer, A.: Femtocell Networks: A Survey. *IEEE Comm. Mag.* 46, 59–67 (2008)
4. Lopez-Perez, D., Valcarce, A., Roche, G., Zhang, J.: OFDMA Femtocells: A Roadmap on Interference Avoidance. *IEEE Comm. Mag.* 47, 41–48 (2009)
5. Wang, Y., Kumar, S., Garcia, L., Pedersen, K., Lovacs, I., Frattasi, S., Marchetti, N., Mogensen, P.: Fixed Frequency Reuse for LTE Advanced Systems in Local Area Scenarios. In: *VTC 2009-Spring*, pp. 1–5 (2009)
6. Yavuz, M., Meshkati, F., Nanda, S., Pokhariyal, A., Johnson, N., Raghothaman, B., Richardson, A.: Interference management and performance analysis of UMTS/HSPA+ femtocells. *IEEE Comm. Mag.* 47, 102–109 (2009)
7. Claussen, H.: Performance of Macro- and Co-Channel Femtocells in a Hierarchical Cell Structure. In: *IEEE PIMRC 2007*, pp. 1–5 (2007)
8. Guvenc, I., Jeong, M., Watanabe, F., Inamura, H.: A hybrid frequency assignment for femtocells and coverage area analysis for co-channel operation. *IEEE Commun. Lett.* 12, 880–882 (2008)
9. Claussen, H., Pivitt, F.: Femtocell Coverage Optimization Using Switched Multi-Element Antennas. In: *IEEE ICC 2009*, pp. 1–6 (2009)
10. Scott, N., Miles, O.: Diversity Gain from a single-Port Adaptive Antenna Using Switched Parasitic Elements Illustrated with a Wire and Monopole Prototype. *IEEE Trans. Antennas Propagat.* 47, 1–10 (1999)
11. Wennstom, M., Svantesson, T.: An antenna solution MIMO channels: The switched parasitic antenna. *Signal, Systems and Computer* 2000 2, 1617–1621 (2000)
12. Mitilineos, S., Thomopoulos, S.: Development of a Compact SPA for 2.4GHz Applications using Commercially Available Elements: Design and Experimental Validation. In: *Loughborough Antennas and Propagation Conference 2009*, pp. 16–17 (2009)
13. Baum, D., Hansen, J., Salo, J.: An interim channel model for beyond-3G systems: extending the 3GPP spatial channel model (SCM). In: *VTC 2005-Spring*, pp. 3132–3136 (2005)
14. European Cooperation in the Field of Scientific and Technical Reserch, EURO-COST231, Digital Mobile Radio Towards Future Generation Systems, COST 231 final report, <http://www.lx.it.pt/cost231/>
15. IEEE 802.16m-08/004r5, IEEE 802.16m Evaluation Methodology Document, EMD (2009)
16. Alipour, A., Hassani, H.: A Novel Omni-Directional UWB Monopole Antenna. *IEEE Trans. Antennas Propagat.* 56, 3554–3857 (2008)

User Policy Based Transmission Control Method in Cognitive Wireless Network

Noriki Uchida¹, Yoshitaka Shibata², and Kazuo Takahata³

¹ Faculty of Software and Information science, Iwate Prefectural University
152-52 Sugo, Takizawa, Iwate 020-0193 Japan
2362006003@sb.soft.iwate-pu.ac.jp

² Faculty of Software and Information science, Iwate Prefectural University
152-52 Sugo, Takizawa, Iwate 020-0193 Japan
shibata@iwate-pu.ac.jp

³ Dept. of Informational Social Studies
Saitama Institute of Technology
1690 Fuzenji, Fukaya, Saitama 369-2093 Japan
takahata@sit.ac.jp

Abstract. Remarkable wireless networks technology developments have made us to expect the realization of new applications like the advanced traffic system, the disaster prevention system, and the adhoc network system. However the resources of wireless bandwidths are not enough to use for such new applications because it is not efficient usage. Therefore, it is necessary to develop with new efficient wireless transmission methods like cognitive wireless network. In this paper, the transmission control methods in cognitive wireless network considering with cross layers including user policies are discussed. First, at the observation stage, the physical data such as user policy, electric field strength, bit error rate, jitter, latency, packet error rate, and throughput are observed. Then, at the decision stage, AHP (Analytic hierarchy process) is applied for decision making process with those parameters. Finally, the action stage, one of the suitable link is chosen and changed links and networks.

In the simulation, ns2 are used for the computational results to the effectiveness of the suggested transmission methods in cognitive wireless networks.

Keywords: Cognitive Wireless Network; QoS; AHP; AODV.

1 Introduction

Recently the development of wireless networks has been growing in various fields, such as cellular phone, digital TV, and wireless Internet services. Those technologies make us to expect the new application like the advanced traffic system, the disaster prevention system, and so on. However, such a remarkable spread of wireless services cause lack and inefficient usage of radio frequency wireless spectrum, we're in front of a spectrum explosion, that is, one may not have enough wireless resources for new wireless applications. According to FCC reports [6], wireless resources are limited

and may not be efficient, because wireless device depends on access to the radio frequency (RF) wireless spectrum, and spectrum has been chronically limited ever since transmissions were first regulated in the early 20th century. Therefore, new technologies that use spectrum more efficiently and more cooperatively, unleashed by regulatory reforms, will soon overcome the spectrum shortage.

Cognitive wireless network (CWN) by J. Mitola III [8] from software defined wireless was originally considered to improve spectrum utilization, and generally considered as a technology to identify the opportunities using the “spectrum holes” for telecommunications [1] [2]. In other words, Mitola defined that CWN is an intelligent wireless communication system that is aware of its surrounding environment, and uses the methodology of understanding by-building to learn from the environment and adapt its internal states to statistical variations in the incoming RF stimuli by making corresponding changes in certain operating parameters in real-time, with two primary objectives in mind: highly reliable communication whenever and wherever needed; efficient utilization of the wireless spectrum.

Also, CWN is consisted of the cogitation cycle; observing the environment, orienting itself, creating plans, deciding, and then acting reconfigurations. The cognition cycle is continually responding stimuli from environment, and it makes successful link or path of network.

However, CWN still has some problems to realize like some algorithms, control methods, and technical problems to attain efficient transmission. First of all, it is recently known that CWN needs not only ordinal spectrum management or Physical/Mac layer observation control but also new control method including upper layers. Also, CWN still has technical problems like how to set radio frequency, how to observe its environment, and how to select proper transmission protocol or radio frequency, and so on [3][4][5]. Moreover, CWN needs the QoS algorithm to select one of wireless links and network route for each application or user environment.

In this paper, we consider the selecting method of the proper wireless link which each node has several wireless links such as IEEE802.11a/b, and IEEE802.16 (WiMax), and the proper route in which network environment is changing through the time. First, user policy and network parameters are observed at the observation stage. At this stage, we set various policies for video, VoIP, text, disaster applications, and so on. By each policy, the weight parameter is decided for the decision algorithm which based on Analytic Hierarchy Process (AHP). Network parameters like throughput, latency, jitter, packet error rate, bit error rate, and electric field strength are continuously measured and also used for the calculation of AHP. Secondly, at the decision stage, the results of AHP for each communication link are compared, and the proper link is selected when the results are changed. Also, if there is no proper wireless links at the observing nodes, proper route is analyzed by extend Ad hoc On-Demand Distance Vector (AODV) algorithm. Finally, in the acting stage, the selected link or route is applied for supposed network, and then simulation is carried out for our supposed methods.

In the followings, in section 2, network model and system architecture of our proposal communication method are defined. Section3 deals with the observation stage which is the method of observing network parameters and user policy. Section 4 describes wireless link selection, route selection and decision algorithm which based on AHP. Section 5 explains how to change link and route at supposed network. At section 6, the

simulation is held for the model by the calculation for the effectiveness of the suggested control method. Finally section 6 derived our conclusion and future works.

2 Network Model

A. Network Configuration

The network configuration of our suggested cognitive wireless network consists of several wireless nodes as shown in Fig. 1.

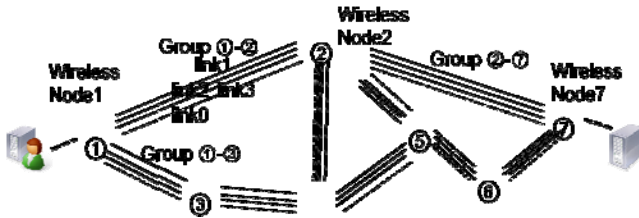


Fig. 1. Wireless Network Architecture

Wireless nodes are supposed as wireless terminal, and they have routing, adhoc, and multi-hop functions. Every node has multi wireless links like IEEE802.11a, IEEE802.11b, IEEE802.16 (WiMax), and IMT2000. The lowest RF link is used as the control links for the purpose of sending a network characteristic data or message of switching links. Also, Antenna of each node is supposed as non-directional.

Moreover, the network condition is changed over time by node’s movement or radio interference like trees or buildings. A data transmission is carried out by ordering user request from a server to a user terminal.

B. System Architecture

The system architecture is organized three layers, including the physical layer, the network layer, and the application layer in Fig 2.

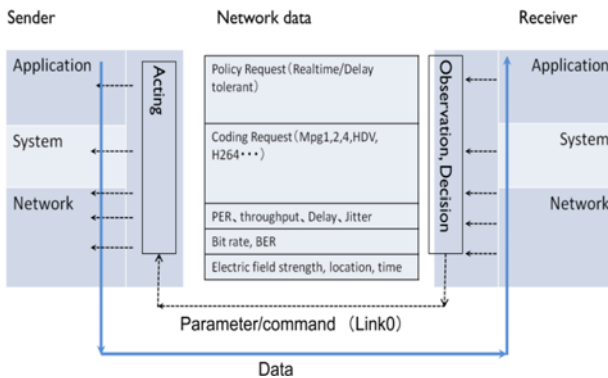


Fig. 2. System Architecture

When transmission data sent from sender to receiver, network data from each layer is observed through link0 that is the control link as explained in the previous section.

In the application layer, a user application is detected for the decision of user policy. A user policy is assumed to various types like video, VoIP, text, connectivity (emergence purpose like disaster) and so on. These policies are used for the decision formula which will be mentioned in the later section.

The system layer collects type of coding and performs coding functions including transcoding of a video coding to another one, such as from/to Motion JPEG to/from MPEG1.

The network layer observes network conditions such as the value of PER, throughput, delay, jitter, BER, electric field strength, and so on.

Network data from each layer is observed, and the decision making is held with crossing through layers. Then, the message of prefer link or route is sent through link0, and a reconfiguration procedure is acted at both sender and receiver nodes.

3 Methodology

As the previous research [8], the cognition cycle is introduced to our proposal method. The cognition cycle is consisted of three stages; the observation stage, the decision stage, and the acting stage. Each stage is continuously cycled in order to perform link or route configuration.

A. Observation Stage

Network data is continuously observed through each layer at this observation stage. Wireless network condition varies depending on the movement of nodes or radio interference. Therefore, CWN needs to parse these stimuli to select the available solution for providing the performance from user requests.

In this paper, our supposed system observes application types in order to decide user policy which is depending on the specific services or media. Also, physical characteristics like coding, PER, throughput, delay, jitter, BER, and electric field strength are observed. Those parameters are used for understanding stimuli from user environment.

B. Decision Stage

Decision making is held to maintain QoS in this decision stage. When the network condition is changed, the proposal system will seek the suitable link and route by the calculating from the values of network characteristics like user policy, throughput, BER, and so on. We introduce AHP for the calculation of link, and extend AODV for the decision of suitable route.

1) Link Selection

AHP is one of multi-attribute decision making and structured techniques for dealing with complex decisions. It was developed by Thomas L. Saaty in the 1970s [9]. By structuring a decision problem hierarchy, AHP provides us quantifications of its elements and evaluations of alternative solutions.

For example, when the suitable link between neighbor nodes is solved by AHP, the hierarchy of the problem is first structured. That is, goal (To decide suitable link of wireless node), criteria (network characteristics such as delay, PER, throughput, and

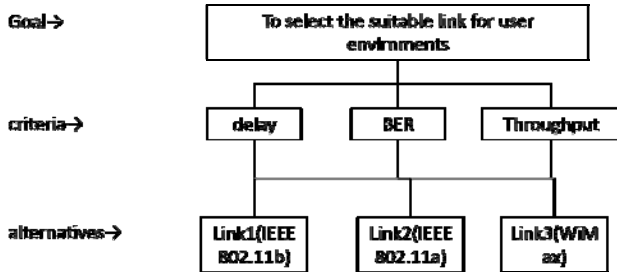


Fig. 3. Example of Hierarchy on AHP

so on), and alternatives (wireless links such IEEE802.11b, IMT2000, and satellite network, and so on).

Then, paired comparisons of criteria and alternatives are calculated for the priority value on AHP. Assume that we are given n as the number of criteria/alternatives, and w_1, w_2, \dots, w_n as the weight of each criteria/alternative, and paired comparison of each element is express as $a_{ij} = w_i/w_j$. These paired comparisons are expressed as the following paired comparison matrix A .

$$A = \begin{bmatrix} \frac{w_1}{w_1} & \dots & \frac{w_1}{w_j} & \dots & \frac{w_1}{w_n} \\ w_1 & & w_j & & w_n \\ \dots & & \dots & & \dots \\ \frac{w_i}{w_1} & \dots & \frac{w_i}{w_j} & \dots & \frac{w_i}{w_n} \\ w_1 & & w_j & & w_n \\ \dots & & \dots & & \dots \\ \frac{w_n}{w_1} & \dots & \frac{w_n}{w_j} & \dots & \frac{w_n}{w_n} \\ w_1 & & w_j & & w_n \end{bmatrix} \quad (1)$$

We have multiplied A on the right by the vector of weights $w=(w_1, w_2, \dots, w_n)^T$, and then the result of this multiplication is λw . That is λ is the eigenvalue of A .

$$Aw = \begin{bmatrix} \frac{w_1}{w_1} & \dots & \frac{w_1}{w_j} & \dots & \frac{w_1}{w_n} \\ w_1 & & w_j & & w_n \\ \dots & & \dots & & \dots \\ \frac{w_i}{w_1} & \dots & \frac{w_i}{w_j} & \dots & \frac{w_i}{w_n} \\ w_1 & & w_j & & w_n \\ \dots & & \dots & & \dots \\ \frac{w_n}{w_1} & \dots & \frac{w_n}{w_j} & \dots & \frac{w_n}{w_n} \\ w_1 & & w_j & & w_n \end{bmatrix} \begin{bmatrix} w_1 \\ \dots \\ w_n \end{bmatrix} = \lambda \begin{bmatrix} w_1 \\ \dots \\ w_n \end{bmatrix} = \lambda w \quad (2)$$

Then, priority is defined as;

$$priority = \frac{w_i}{\sum w_i} \quad (3)$$

At this time, we introduce network data like the following in order to acquire network environment for the alternatives.

$$S_i = \left(\frac{n_i - l_i}{u_i - l_i} \right) \times 10 \text{ (In case of throughput)} \tag{4}$$

$$S_i = \left(1 - \frac{n_i - l_i}{u_i - l_i} \right) \times 10 \text{ (In case of BER, PER, latency, etc)} \tag{5}$$

In the above formula, S_i is the weight of each alternative, u_i and l_i are the upper and lower limits of the alternative, and n_i is the observed value from network.

We also adopt the value of consistency index (C.I.). C.I. expresses the characteristics of polynomial of A, and we accept the estimation of w if C.I. is less than certain threshold.

$$C.I. = (\lambda_{\max} - n) / (n - 1) \tag{6}$$

In this paper, we propose the calculation of AHP based on each user policy like video, VoIP, or connectivity for decision making of the suitable link. For example, in case of emergency like disaster, connectivity is more important as user policy. Thus, the criteria like electric field strength and BER are weighed and calculated by the formula (3). Then, priorities of alternatives is calculated by (4) and (5), and the results are inserted for (3). Finally, each value of alternative is calculated by alternative priority multiplied with criteria priory, and the link with largest value will be decided as the best suitable link.

2) Route Selection

Our proposed system would change network route if the suitable link is not found or minimum requirement for user is not satisfied. When a network route needs to change, we introduced extend AODV for the decision of suitable route.

AODV is a routing protocol for mobile ad hoc networks (MANETs) and other wireless ad-hoc networks.[10][11] It is a reactive routing protocol, which means that it establishes a route to a destination only on demand base. Therefore, the connection is slower than proactive routing protocols like OLSR or TBRPF. But AODV is superior when the network condition changes so often or changes so slowly. [7].

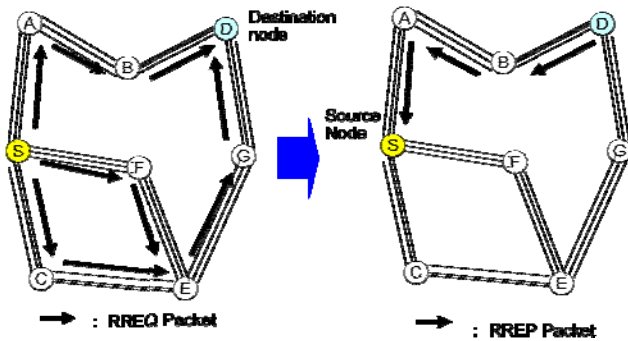


Fig. 4. RREQ and RREP

AODV builds routes using a route request (RREQ) and route reply (RREP). When a source node requires a route to a destination for which it does not have a route, the source node broadcasts RREQ packets across the network. When the other nodes receive those packets, and then update their routing information for the source node and set up backwards pointers to the source node in the route tables.

The RREQ contains the source node's IP address, current sequence number, broadcast ID, and the most recent sequence number for the destination node. A node receiving the RREQ may send RREP packets if it is either the destination or it has a route to the destination with corresponding sequence number greater than or equal to that contained in the RREQ. If this is the case, RREP packets are sent back to the source node with setting the routing tables relatively.

3) *Extended AOCV*

We propose a extended AODV protocol by adding the link values and network conditions on each node for RREQ and RREP packets on each node from the source to the destination. The link values of each node are calculated by AHP as the previous sections, and those are set to the information of next node ID. Also, network characteristics values of each node such as delay, PER, throughput, and so on are added to RREQ and RREP packets. Then, a destination node receives all of the RREQ packets from the possible routes during certain unit time. Then, possible routes are compared to select the best route,

- 1) that can provide the maximum End-to-End throughput among all of the routes for video service, or
- 2) that can provide the minimum End-to-End delay time among all of the routes for VoIP, and
- 3) that can optimize the policy based AHP for Web service.

Through those comparisons, the best suitable route is decided by the policy of transmission data.

We propose to add the link values by AHP for RREQ and RREP packets, and those are set in the routing tables like Fig.5.

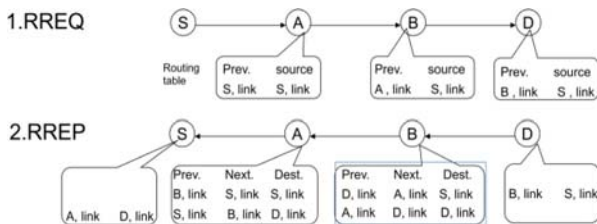


Fig. 5. Routing Tables by Extend AODV

C. Acting Stage

After the decision making of the suitable link or the route, a link or route will be changed in the acting stage.

In the proposed system, a control wireless link is assumed to use for a transmission of the control information. A control link is set to the lowest frequency bandwidth, and a node sends the information with a suitable link to the next node through the

control link. Then, both nodes act the link change at the same time. A route change also acts as the same way of a link change.

The cognition cycle goes back to the observation stage after the acting stage, and it works to a link or route change relatively.

4 Simulation

In the simulation, ns2 (Network Simulator 2) was examined to evaluate the effectiveness of our suggested method. IEEE802.11a, IEEE802.11b, IEEE802.11g is used as the wireless links, and IMT2000 is used as the control link. The number of wireless nodes is set to two, and every node has the same wireless links. Then, only a source node moved with the same speed, and video data were transmitted from the source node to the destination node. The simulation conditions are the following table.

Table 1. Simulated Conditions

Item	Simulation Content
Nodes	Node0 moves from (10,5,5) to (300,5,5) with 2m/s Node1 is set to (10,5,5)
Antenna	Non-directional
Transmission Data	320x240 MJPEG (15fps, 1/15 Compressed data) is send from node0 to node1.
Link Interface	IEEE802.11a/b/g. Each link has one channel.
Space	Free Space
Observed Parameters	Signal Strength, Throughput, Jitter, PER

In this paper, a link change is especially focus on the simulation. First, the simulation of IEEE802.11a, b and g is simulated under table1 conditions.

The Fig.6 shows the results of each node simulation. At Fig.6, throughput of IEEE802.11a is quickly down about 18sec, and the link is disconnected about 23sec. Also, IEEE802.11g is gradually down and disconnected about 145sec. IEEE802.11b keeps link connection although throughput is not so high.

On the contrary, the results of our proposal method are showed in Fig.7. Fig.7 shows that our proposed method keeps high throughput because the link selection is effectively worked under video transmission scenario. The wireless link is switched from IEEE802.11a to IEEE802.11g at 14 sec, and also switched IEEE802.11g to IEEE802.11b at 132 sec. This results showed the proposed methods change the best suitable link to maximize the transmission rate with user policy.

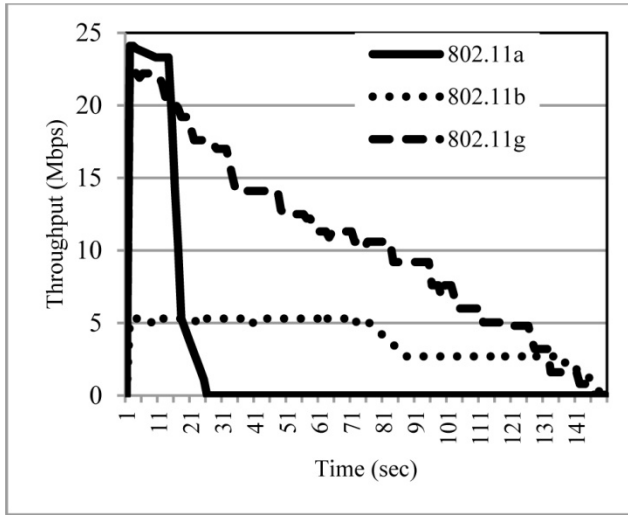


Fig. 6. IEEE 802.11a/b/g Throughput

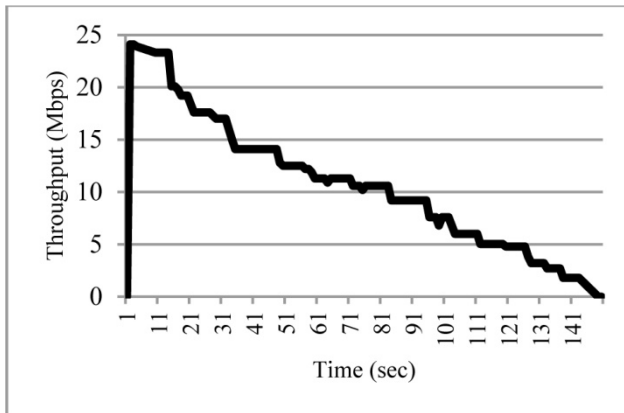


Fig. 7. Proposed Cognitive Wireless Network Throughput

5 Conclusions

In this paper, the transmission control methods in cognitive radio network considering with cross layers including user policies are introduced. Compared with giving fixed transmission rate, the our suggested methods can determines the transmission rate dynamically to achieve the maximum throughput by changing wireless links. The simulation result based on the video based data show that the proposed methods change the best suitable link to maximize the transmission rate with user policy. Therefore, the proposed method is shown in practical, efficient and reliable to support data transmission.

References

- [1] Chen, K.C., Peng, Y.J., Prasad, N., Liang, Y.C., Sun, S.: Cognitive radio network architecture: part I – general structure. In: ICUIMC 2008: Proceedings of the 2nd international conference on Ubiquitous information management and communication (January 2008)
- [2] Chen, K.C., Peng, Y.J., Prasad, N., Liang, Y.C., Sun, S.: Cognitive radio network architecture: part II – trusted network layer structure. In: ICUIMC 2008: Proceedings of the 2nd international conference on Ubiquitous information management and communication (January 2008)
- [3] Cordeiro, C., Challapali, K., Ghosh, M.: Cognitive PHY and MAC layers for dynamic spectrum access and sharing of TV bands. In: TAPAS 2006: Proceedings of the first international workshop on Technology and policy for accessing spectrum (August 2006)
- [4] Weingart, T., Sicker, D.C., Grunwald, D.: Evaluation of cross-layer interactions for reconfigurable radio platforms. In: TAPAS 2006: Proceedings of the first international workshop on Technology and policy for accessing spectrum (August 2006)
- [5] Kliazovich, D., Granelli, F.: Packet concatenation at the IP level for performance enhancement in wireless local area networks. *Wireless Networks* 14(4) (2008)
- [6] Renk, T., Kloock, C., Burgkhardt, D., Jondral, F.K., Grandblaise, D., Gault, S., Dunat, J.-C.: Bio-inspired algorithms for dynamic resource allocation in cognitive wireless networks. In: *Mobile Networks and Applications*, October 2008, vol. 13(5) (2008)
- [7] Staple, G., Werbach, K.: IEEE Spectrum: The End of Spectrum Scarcity, <http://www.spectrum.ieee.org/telecom/wireless/the-end-of-spectrum-scarcity>
- [8] Sugimoto, T., Yamaguhi, S., Asatani, K.: A Proposal for System Selection Scheme Using Multiple Handover Triggers in Heterogeneous Wireless Networks, Technical Report of IEICE
- [9] He, W., Nahrstedt, K., Liu, X.: End to End Delay Control of Multimedia Applications over Multihop Wireless Links. *ACM Transaction on Multimedia Computing, Communications and Applications* 5(2), Article 16 (November 2008)
- [10] Perkins, C.E., Royer, E.M.: Ad hoc On-Demand Distance Vector Routing. In: *Proceedings of the 2nd IEEE Workshop on Mobile Computing Systems and Applications*, February 1999, pp. 90–100 (1999)
- [11] Chakeres, I.D., Royer, E.M.: AODV Routing Protocol Implementation Design. In: *Proceedings of the International Workshop on Wireless Ad Hoc Networking (WWAN)* (March 2004)

Development of a Hybrid Case-Based Reasoning for Bankruptcy Prediction

Rong-Ho Lin^{1,*} and Chun-Ling Chuang²

¹ Department of Industrial Engineering and Management,
National Taipei University of Technology, Taipei 10608, Taiwan, R.O.C.
rhlin@ntut.edu.tw

² Department of Information Management, Kainan University,
Taoyuan, 33857, Taiwan, R.O.C.
clchuang@mail.knu.edu.tw

Abstract. This paper aims to develop an integrated model of predicting business failure, using business financial and non-financial factors to diagnose the status of business, thereby providing useful references for business operation. This study applied Rough Set Theory to extract key financial and non-financial factors and Grey Relational Analysis (GRA) as the approach of assigning weights. In addition, Case-Based Reasoning (CBR) are adopted to propose a new hybrid models entitled RG-CBR (combining RST and CBR with GRA) to compare the accuracy rates in predicting failure. After exploring the TEJ (*Taiwan Economic Journal*) database and conducting various experiments with CBR, RST-CBR and RG-CBR the study finds CBR, RST-CBR and RG-CBR reporting an accuracy rate in predicting business failure of 49.2%, 59.8% and 83.3% respectively. The RG-CBR boasts the highest accuracy rate while also effectively reducing Type I and Type II error rates.

Keywords: decision analysis, data mining, business failure prediction, rough set theory, case-based reasoning.

1 Introduction

Bankruptcy prediction is one of the most important issues influencing financial and investment decision-making [1]. While affecting the entire life span of a corporation or an organization, a business failure can even undermine social functions and agitate economical structures [2]. Therefore, the evaluation of business failure has emerged as a professional area prompting academics and professionals to develop optimal prediction models based on their specific interest for practical applications in the business community.

Methods developed, such as univariate statistical model, Multiple Discriminant analysis, Linear Probability models, Logistic Regression, and Probit analysis [3][4], are now generally regarded as obsolete statistical methods [5]. Statistical methods

* Corresponding author.

require a few assumptions and normally being constrained by the linearity, thus they are quietly difficult to deal with the mass and complicated data [6]. To meet the varied and restrictive assumptions, such as a large number of samples, normal distributed independent variables, and linear relationship between all variables, scholars in related fields have come up with alternatives methods to predict the risk of business failure, notably Artificial Neural Network (ANN), Case Based-Reasoning (CBR), and Data Mining techniques [7][8][9][2].

Although reported to perform fairly well in terms of predictive accuracy, ANN classifies every object by its computational characteristic and is therefore often criticized as a “black box” approach for its lack of transparency [10]. In contrast to the “black box” scenario, CBR, developed by Schank and Abelson [11] in the 1970s, uses “similar” cases stored in a case base to build a solution to a new case based on the adaptation of solutions to past similar cases [12]. Literature indicates that CBR has been widely used in diagnosis domain and enhanced some of the deficiencies in statistical models and ANN [13][14]. However, accuracy and effectiveness may be reduced when CBR deals with too many attributes.

To solve the problem of excessive variables (attributes), Rough Set Theory (RST), a data mining technique that has been successfully applied to solve classification problems, can be adopted to reduce unnecessary attributes. Attributes, once identified, need to be further prioritized by taking their weights into due consideration and this task can be accomplished by performing Grey relational analysis (GRA) that compares quantitative analysis to the development between every variable in the grey system dynamically. The method describes the relation between the main factor (the business status) and other variables in the grey system.

This study accordingly proposes a hybrid model for predicting business failure, called RG-CBR, which incorporates three stages starting with applying RST to extract key attributes, adopting GRA to obtain the weights of the key attributes, and culminating in feeding the retrieved key attributes and weights into CBR to improve the model’s prediction accuracy.

Following the present section as the introduction, four major sections are incorporated in the paper. Section 2 outlines the backgrounds to RST, GRA, and CBR. Section 3 explains the research methodology. Section 4 centers on experiment results and model evaluations while Section 5 presents our conclusions.

2 Research Background

2.1 Rough Set Theory (RST)

RST, a powerful mathematical tool introduced by Pawlak [15], can discover facts from imperfect data [16]. Studies showed that RST has played an important role in coping with classification problems [9][2][7]. It also has been successfully applied to real-world classification problems with the following advantages [7]: (a). Identifying potentially important information contained in a large number of cases and reducing the information into a model including a generalized description of the knowledge; (b). Requiring no interpretation as the model is a set of easily understandable

decision rules; and (c). Demanding no additional information, such as probability in statistics.

Many researchers [7] [17][2] have used RST to construct a prediction model or combined it with other methods as a data preprocessor. In particular, when RST further integrated with other methods, such as the proposed hybrid model, the results tend to be more accurate than those obtained by using RST alone [17][2].

2.2 Grey Relational Analysis (GRA)

Deng [18] pioneered a mathematical method called Grey Relational Analysis (GRA) which is characterized by good performance in analyzing few data and many variables to examine the relationship among factors in observable systems and construct the prediction model. GRA has been successfully used in a wide range of fields [19], including agriculture, traffic, industrial engineering, education. Only a few studies, however, have applied it to tackle business prediction problems.

2.3 Case-Based Reasoning (CBR)

An analogical reasoning method developed by Schank & Abelson [11] in the 1970s, CBR can serve as a useful research paradigm for exploring complex problems [12] by adapting past experiences to acquire potential solutions. As a flourishing problem-solving method, CBR has been successfully applied to a wide spectrum of different fields, such as medical diagnosis [20], bankruptcy prediction [8], fault diagnosis [21], risk analysis [22], and marketing [2]. The study by Bryant [8] that adopts CBR to bankruptcy prediction modeling in particular underlines the importance of selecting proper variables in the successful application of the method. Complete reviews of CBR with comprehensive details and applications for further works can be found in the study of Watson & Marir [23].

3 Research Model Development

The study first applies CBR to analyze previous experiences of failed companies and knowledge about business failure prediction. However, it is not easy to obtain successful results with high prediction accuracy by applying CBR alone. The measures of success of a CBR system depend on its ability to index cases and retrieve the most relevant ones in support of the solution to a new case. RG-CBR is accordingly proposed to strengthen the capacity of CBR. The model integrates RST as a data preprocessor to identify key attributes first and proceeds to compute and assign different weights to the identified key attributes by using GRA as an alternative method, and then culminate in feeding the retrieved key attributes and weights into CBR. To evaluate the prediction accuracy of RG-CBR model, the Figure 1 shows the models of CBR, RST-CBR and RG-CBR in (a), (b), and (c), respectively, and the study further compares its performance with those of CBR, and CBR combined with RST. RG-CBR model comprises three stages each of which is introduced as follows:

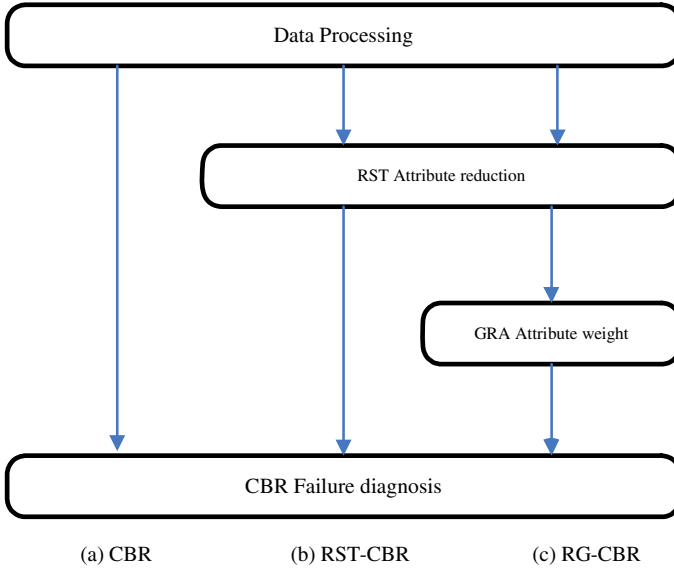


Fig. 1. Configuration of CBR, RST-CBR, and RG-CBR

3.1 RST

At the beginning, data processing is initiated to collect samples and establish criteria for sampling and coding. RST is adopted to reduce attributes; some redundant attributes are eliminated without any information loss to induct the minimal variables subsets named “reducts.” The following sections present some basic concepts of RST.

A finite set of objects U and a finite set of attributes A can be considered as a relational system $S = (U, A)$. Each attribute a belongs to the considered set of attributes A ($a \in A$); $a(x)$ represents the value of attribute a of object x . For every set of attributes $B \in A$, an indiscernibility relation $Ind(B)$ is defined in the following way:

$$Ind(B) = \left\{ (x_i, x_j) \in UXU; b(x_i) = b(x_j); b \in B \right\} \quad (1)$$

The equivalence class of x_i in relation to $Ind(B)$ is represented as $[x_i]_{Ind(B)}$. Let X denote the subset of elements of U , the upper and lower approximations can be represented as

$$\overline{B}(X) = \{x_i \in U \mid [x_i]_{Ind(B)} \cap X \neq \emptyset\} \quad (2)$$

$$\underline{B}(X) = \{x_i \in U \mid [x_i]_{Ind(B)} \subseteq X\} \quad (3)$$

The boundary of X in U is defined as $Bnd(X) = \overline{B}(X) - \underline{B}(X)$. If $Bnd(X)$ is an empty set, then the set X is definable with respect to B ; otherwise, X will be referred to as a rough set. “Reduct” and “core” are two fundamental RST concepts. Assume that the

set of attributes $R \subseteq A$ depends on the set of attributes $B \subseteq A$ in S (denotation $B \rightarrow R$) iff $Ind(B) \subseteq Ind(R)$. The minimal subset is $R \subseteq B \subseteq A$ such that $\mu_B F = \mu_R F$ ¹ is called F -reduct of B , which means a reduct is the minimal subset of attributes to represent the whole set of attributes. In addition, S might have more than one F -reduct. Intersection of all F -reducts is called F -core of B , i.e. $CORE_B(F) = \bigcap RED_B(F)$. The core is a collection of the most significant attributes in the system. In other words, the attributes in B are absolutely necessary attributes.

3.2 GRA Attribute Weighting

Grey relational analysis (GRA) can measure the degree of interrelationships between one major sequence and another sequence in a given system [24]. The relational grade is regarded as high if the two sequences tend towards concordance; otherwise, the relational grade is regarded as low. Assume that the objective sequence is defined as

$$X_0(k) = [x_0(1), x_0(2), \dots, x_0(k)], \quad k = 1, 2, \dots, n \tag{4}$$

The reference sequence can be represented as follows:

$$X_i(k) = [x_i(1), x_i(2), \dots, x_i(k)], \quad i = 1, 2, \dots, l \quad k = 1, 2, \dots, n \tag{5}$$

Given the objective sequence $X_0^*(k)$ and the reference sequence $X_i^*(k)$ with the normalized form, the grey relational coefficient $r_{0i}(k)$ between them at any data point k is defined as

$$r_{0i}(k) = \frac{\Delta_{\min} + \rho \Delta_{\max}}{\Delta_{0i}(k) + \rho \Delta_{\max}} \tag{6}$$

Where $\Delta_{0i}(k) = |x_0^*(k) - x_i^*(k)|$ is termed as the relative grey relational space;

$$\Delta_{\min} = \min_i \min_k |x_0^*(k) - x_i^*(k)| \tag{7}$$

$$\Delta_{\max} = \max_i \max_k |x_0^*(k) - x_i^*(k)| \tag{8}$$

Δ_{\min} and Δ_{\max} are respectively the minimum and maximum distances for all factors in all sequences. ρ is the distinguishing coefficient ($0 \leq \rho \leq 1$), and usually $\rho = 0.5$. After the grey relational coefficient is obtained, the grey relational grade, i.e. the mean of the coefficient, can be measured by

¹ $\mu_B F$ and $\mu_R F$ is called the quality of approximation of partition F by set of attributes B and R .

$$\gamma(X_0, X_i) = \frac{1}{n} \sum_{k=1}^n r_{0i}(k) \quad (9)$$

3.3 CBR Failure Diagnosis

A case is a contextualized piece of knowledge representing an experience. It contains past lessons that are the content of the case and the context in which the lessons can be used [25]. It is difficult to share experience and knowledge for future reuse, and CBR is developed as a problem solving technique by making use of past cases and experiences to solve a new problem. It can retrieve a set of similar cases and identify the most similar case for reuse to help solve the problem at hand. The process of CBR includes case representation, case retrieval, and case adaptation. In particular, the measures of success of a CBR system depend on its ability to index cases and retrieve the most relevant ones in support of the solution to a new case.

A case should contain both content and context, typically composed of the problem, solution, and outcome [26]. Cases can be represented in different forms. In this study, a structured case is defined as a vector of features: $x = \{x_1, x_2, \dots, x_n\}$, where n is the number of features. Given a description of a problem, a retrieval algorithm, using the indices in the case-memory, needs to be adopted to retrieve the cases most similar to the current problem or situation. The retrieval algorithm relies on the indices and the organization of the memory to direct the search to potentially useful cases. Shin and Han [26] categorize case indexing into three types: nearest neighbor, inductive, and knowledge-guided. Nearest neighbor is the most commonly used approach [27].

The method is used to measure the similarity between two cases and known for its easy application to numerical data such as financial ratios. In general, use of the nearest neighbor method leads to the retrieval time increasing linearly with the number of cases. Therefore, this approach is more effective when the case base is relatively small. A typical algorithm for calculating nearest neighbor matching is defined by [12] as:

$$\text{Similarity Score} = \frac{\sum_{i=1}^n w_i \times \text{sim}(f_i^I, f_i^R)}{\sum_{i=1}^n w_i} \quad (10)$$

where $\text{sim}(f_i^I, f_i^R) \in [0,1]$ is the similarity function, and f_i^I and f_i^R are the values for attribute i in the input and retrieved cases, respectively.

4 Experimental Design and Results

Based on information and data acquired from Taiwan Stock Exchange Corporation (TSEC) and Taiwan Economic Journal (TEJ) database, the firms in Taiwan's IT industry judged as failed ones during the period from 1999-2006 are taken as the

samples. A qualified sample, i.e. a failed firm, refers to a company with the transaction of its listed securities placed under the “altered-trading method” category, suspended or terminated. More specifically, a firm is cited as failed for (a). Suffering from credit crisis, (b). Having net operating loss, (c). Failing to pay debts, or (d). Violating related regulations.

In order to improve the comprehensiveness and accuracy in predicting, this study will concern simultaneously non-financial (qualitative) and financial (quantitative) attributes and then apply them to diagnose the business status with different models. 21 attributes (13 financial ones and 8 non- financial) used in previous researches are selected for our study and coded and described in Table 1.

The sample data sets consist of an equal number of healthy and failed companies (321 companies in each category). Sample companies are grouped by binary assignment into a decision class (healthy or failed, coded by 1 and 0, respectively). The data sampling follows the “pairwisd sampling method” first proposed by Beaver [3]. Each data set is split into two subsets: a training set of 70% and a validation set of 30% of the data, respectively. The training data are used as a case base, and the validation data not used in constructing the case base are adopted to test the model. Five different validation sets, each with six entries of data, are sampled and then used to validate the model performance; that is, a total of 97 test samples are used for validation.

As RST is concerned with discrete values, the experimental study uses the software ROSETTA that includes a method for discretizing continuous attributes called entropy heuristic approach. The approach will subject continuous data to a process of discretization.

Table 1. List of selected attributes

Attributes	Description
A1	Return on assets (%)
A2	Return on equity (%)
A3	Net income (%)—except disposed
A4	Gross margin (%)
A5	Net income (%)
A6	Current ratio (%)
A7	Acid test ratio (%)
A8	Liabilities ratio (%)
A9	TCRI credit ranking
A10	Cash flow operation to current liabilities (%)
A11	Total equity growth ratio (%)
A12	Return on total assets growth ratio (%)
A13	Accounts receivable turnover (day)
A14	Inventory turnover (%)
A15	Earning per share
A16	Added value per person
A17	Manager-director
A18	Director and supervisor shareholding
A19	Inflation rate (%)
A20	Business cycle
A21	Rediscount Rate (%)

Table 2. Measuring Data of each key attribute

D	A6	A7	A8	A9	A16	A18
1	57.91	27.46	26.5	8	1289.5	19.41
1	32.57	21.25	80.2	9	590.4	12.37
1	63.69	54.94	100.55	10	1188.81	7.45
1	144.5	88.71	43.55	8	2494.46	8.59
1	80.98	35.69	78.36	9	181.06	11.4
1	143.61	85.9	71.51	10	47.41	24.49
1	249.79	86.01	85.27	10	619.21	35.29
1	71.01	51.43	82.89	10	1105.08	25.36
1	2.35	1.24	94.2	10	0.78	8.89
1	512.16	438.14	11.32	8	489.64	14.74
1	199.93	151.11	44.92	5	8.68	5.7
1	117.09	63.13	52.87	9	978.6	13.46
1	225.23	143.6	44.52	6	669.41	10.47
1	30.87	64.94	97.29	10	635.85	11.61
1	209.79	158.66	50.24	10	2017.14	14.75
1	298.29	242.31	18.77	5	424.59	35.36
1	135.29	128.63	63.15	7	243.78	8.88
1	133.45	121.01	55.77	8	1377.4	9.4
1	147.06	123.46	69.68	6	747.69	16.41
1	167.64	137.2	46.55	2	1852.75	13.61

The discernibility matrix can be used to find the minimal subset(s) of attributes called reducts who have the same quality as the complete set of attributes. As a result, we further apply RST to reduce the redundant attributes and to combine with CBR for lowering misclassification rate. At the beginning, there were 78 reducts computed. These reducts were employed from the total 21 attributes. According to the results obtained from RST, the measuring data of each key attribute is shown in Table 2. The weights of those key attributes are developed by GRA in Excel, as described in Table 3 which shows the grey relational coefficient of each attribute adjusted its weight. The six attributes are A6 (Current ratio), A7 (Acid test ratio), A8 (Liabilities ratio), A9 (TCRI credit ranking), A16 (Added value per person) and A18 (Director and supervisor shareholding).

Table 3. Key attributes with GRA-weights

Attributes	Description	Weight
A6	Current ratio	0.032
A7	Acid ratio	0.048
A8	Liability ratio	0.333
A9	TCRI credit ranking	0.529
A16	Added value per person	0.005
A18	Director and supervisor shareholding	0.053
total		1.00

The six attributes are ranked based on their weights into the following order: TCRI credit ranking, Liabilities ratio, Director and supervisor share-holding, Acid test ratio, Current ratio, and Added value per person. An attribute is ignored during similarity computation if it is found to have zero weight, whereas an attribute with a higher weight exerts a greater impact on determining similarity. As indicated above, TCRI credit ranking emerges as the most significant attribute influencing business failure prediction with a highest weight of 0.529.

Identifying key attributes is a sub-process of investigating the important attributes of factors critical for identifying analogous cases and of predicting the value of the target variable. RST for identifying key features and GRA for ranking key features by their weights are applied to compensate the weaknesses of CBR. And the enhanced CBR forms the core of the proposed RG-CBR model. The model then proceeds to search in the case base the case that best matches with the target case and to adapt the retrieved case to fit the problem and produce the solution for the problem. In the experiment, the overall predictive accuracy of the validation data emerges to be 83.3% as shown in Table 4.

To evaluate the prediction accuracy of RG-CBR model, Table 4 compares the results of the three different models of RG-CBR, CBR, and RST-CBR (combining RST and CBR with equal weights). Among the three models, RG-CBR shows the highest level of accuracy (83.3%) in the given validation data. Therefore, the result concludes that RG-CBR outperforms CBR and RST-CBR and testifies to its qualification to serve as a promising alternative for business failure prediction.

Table 4. Performance Comparison of four different Models

Model	Accuracy (%)	Type I error (%)	Type II error (%)
CBR	49.2	20	30.8
RST-CBR	59.8	20	20.2
RG-CBR	83.3**	20.2	13.3

** signifies the best performance.

5 Conclusions

CBR, as indicated in the literature, has been widely used in diagnosis domain by retrieving most similar cases more accurately and effectively. The accuracy and effectiveness, however, may be reduced when there are too many attributes. It is therefore not easy to obtain successful results by adopting CBR alone to predict business failure. In order to enhance the performance of the CBR system, this study proposes a hybrid model, RG-CBR, for predicting business failure. RG-CBR can refer the user to the similar cases as references for making decisions, and the attributes retrieved and assigned different weights by mathematical methods (RST with GRA) are more objective than those decided by experts' subjective experiences. As suggested by the results of the experiment, RG-CBR outperforms other comparative algorithms such as CBR, and RST-CBR in terms of accuracy and serves a promising alternative for business failure prediction.

References

- [1] O'Leary, D.E.: Using neural networks to predict corporate failure. *International Journal of Intelligent Systems in Accounting, Finance and Management* 7, 187–197 (1998)
- [2] Ahn, B.S., Cho, S.S., Kim, C.Y.: The integrated methodology of rough set theory and artificial neural network for business failure prediction. *Expert Systems with Applications* 18, 65–74 (2000)
- [3] Beaver, W.H.: Financial ratios as predictors of failure- empirical research in accounting: selected studies. *Journal of Accounting Research* 4, 71–111 (1966)
- [4] Ohlson, J.A.: Financial ratios and the probabilistic prediction of bankruptcy. *Journal of Accounting Research* 18, 109–131 (1980)
- [5] Balcaen, S., Ooghe, H.: 35 years of studies on business failure an overview. *The British Accounting Review* 38, 63–93 (2006)
- [6] Akahoshi, M., Amasaki, Y., Soda, M.: Correlation between fatty liver and coronary risk factors: a population study of elderly men and women in Nagasaki, Japan. *Hypertens Research* 24, 337–343 (2001)
- [7] Slowinski, R., Zopounidis, C.: Application of the rough set approach to evaluation of bankruptcy risk. *International Journal of Intelligent Systems in Accounting, Finance and Management* 4, 27–41 (1995)
- [8] Bryant, S.M.: A case-based reasoning approach to bankruptcy prediction Modeling. *International Journal of Intelligent Systems in Accounting, Financial and Management* 6, 195–214 (1997)
- [9] McKee, T.E.: Developing a bankruptcy prediction model via rough set theory. *International Journal of Intelligent Systems in Accounting, Finance and Management* 9, 159–173 (2000)
- [10] Kumar, P.R., Ravi, V.: Bankruptcy prediction in banks and firms via statistical and intelligent techniques – A review. *European Journal of Operational Research* 180, 1–28 (2007)
- [11] Schank, R. (ed.): *Dynamic Memory: A Theory of Learning in Computers and People*. Cambridge University Press, New York (1982)
- [12] Kolodner, J.: Improving human decision making through case-based decision aiding. *AI Magazine* 12, 52–68 (1991)
- [13] Barletta, R.: An introduction to case-based reasoning. *AI Expert* 6, 42–49 (1991)
- [14] Park, T.S., Han, J.Y.: Derivation and characterization of pluripotent embryonic germ cells in chicken. *Molecular Reproduction and Development* 56, 475–482 (2000)
- [15] Pawlak, Z.: Rough sets. *International Journal of Information and Computer Sciences* 11, 341–356 (1982)
- [16] Walczak, B., Massart, D.L.: Rough sets theory. *Chemometrics and Intelligent Laboratory Systems* 47, 1–16 (1999)
- [17] McKee, T.E., Lensberg, T.: Genetic programming and rough sets: A hybrid approach to bankruptcy classification. *European Journal of Operational Research* 138, 436–451 (2002)
- [18] Deng, J.: Control problems of grey systems. *System and Control Letters* 1, 288–294 (1982)
- [19] Lai, H.H., Lin, Y.C., Yeh, C.H.: Form design of product image using grey relational analysis and neural network models. *Computers & Operations Research* 32, 2689–2711 (2005)
- [20] Varma, A., Roddy, N.: ICARUS: Design and deployment of a case-based reasoning system for locomotive diagnostics. *Engineering Applications of Artificial Intelligence* 12, 681–690 (1999)

- [21] Yang, B.S., Han, T.H., Kim, Y.S.: Integration of art-kohonen neural network and case-based reasoning for intelligent fault diagnosis. *Expert Systems with Applications* 26, 387–395 (2004)
- [22] Jung, C., Han, I., Suh, B.: Risk analysis for electronic commerce using case-based reasoning. *International Journal of Intelligent Systems in Accounting, Finance and Management* 8, 61–73 (1999)
- [23] Watson, I., Marir, F.: Case-based reasoning: A review. *The Knowledge Engineering Review* 9, 327–354 (1994)
- [24] Hsu, Y.T., Chen, C.M.: A novel fuzzy logic system based on N-version programming. *IEEE Transactions on Fuzzy Systems* 8, 155–170 (2000)
- [25] Kolodner, J.: *Case-based reasoning*. Morgan Kaufmann, San Mateo (1993)
- [26] Shin, K.S., Han, I.: A case-based approach using inductive indexing for corporate bond rating. *Decision Support Systems* 32, 41–52 (2001)
- [27] Burke, E.K., MacCarthy, B., Petrovic, S., Qu, R.: Structured cases in case-based reasoning-re-using and adapting cases for timetabling problems. *Knowledge-Based Systems* 13, 159–165 (2000)

The Time Series Image Analysis of the HeLa Cell Using Viscous Fluid Registration

Soichiro Tokuhisa and Kunihiko Kaneko

Department of Intelligent Systems, Graduate School of Information Science and Electrical Engineering, Kyushu University
744 Motoooka, Nishi, Fukuoka, Japan
{tokuhisa, kaneko}@db.is.kyushu-u.ac.jp

Abstract. Optical microscopy image analysis is important in the life science research. To obtain the motion of the cell, we use the viscous fluid registration method based on fluid dynamics. Viscous fluid registration deforms an image at time t to the next image at time $t+1$. In this algorithm, there is a problem that an object cannot be divided into two. In other words, the divided objects from one are connected by thin line because the velocity field on the connected thin line is zero. To solve this problem, we suggest a new viscous fluid registration algorithm for the object division. This algorithm is only added similarity maximization step to correct the displacement in the near pixels in the original viscous fluid registration. Using this method, one object is divided into two, and divided objects are not connected. We experiment the anaphase detection based on a nucleus identification using laser scanning microscope HeLa cell images. Experimental result shows that 74 in 76 cells are tracking well and 6 cells in the anaphase are detected. In three scenes in the cell division which can not be divided into two using original viscous fluid registration, suggested algorithm can be divided into two cells completely.

Keywords: HeLa cell, cell division, cell tracking, viscous fluid registration.

1 Introduction

Microscopy image analysis is important in the life science research. Cell phase identification and analysis of the cell features are attracted. The abnormal cell detection method helps early detection of a cancer, diabetes, and the gene deficit disease. Therefore early treatment for disease is enabled and helps medical care.

In the field of the cell analysis, various image processing methods are performed in cell segmentation and the phase identification. For example, it is suggested the accurate automated method that cell segmentation using watershed algorithm and the cell phase identification using Markov model by Zhou et al. [1] and Chen et al. [2]. Dzyubachyk et al. suggested the multiple level-set framework that each cell have each level-set function [3].

We think that the cell movement can be expressed by the viscous fluid, so the cell movement can be described as the time series image matching technique. We use the

image registration method applied in the medical image processing such as the CT image and MRI image. Though this method has a large calculation cost, it is obtained the smooth displacement because it is based on the fluid dynamics. The displacement contains vectors $\mathbf{x}_t - \mathbf{x}_{t-1}$. It means that a pixel at position \mathbf{x}_t at time t is moved from a pixel at position \mathbf{x}_{t-1} at time $t-1$.

There are some kinds in a microscopy images such as an optical microscopy image and confocal laser scanning microscope. The three-dimensional cell image can be constructed using confocal laser scanning microscope. In this experiment, we use these cell images in Fig.1 that are obtained by live cell imaging [4] of the HeLa cell with mCherry-NLS[5] via the confocal laser scanning microscope created by Yuki Tsujimura in RIKEN.

When an object is divided into two objects, viscous fluid registration cannot to be separate correctly. So, we suggest the new algorithm that is added one calculation step in the viscous fluid registration which step is the similarity maximization step to correct the displacement in the near pixels.

When this method is used, we can divide an object to two objects completely. Using viscous fluid registration with similarity maximization step, it will be performed a smooth cell movement analysis. The displacement can be use for the time series image analysis like Fig.2. The displacement is useful for anaphase detection because the corresponding time series cell at time t and cell at time $t+1$ can be detected. The cell in anaphase is defined the first time of the cells called daughter cell which was divided one cell into two in this article. From the cell identification result, it can be confirmed the movement of the daughter cells after the cell division. Cell identification result shows that 74 in 76 cell in 17 images is success in tracking.

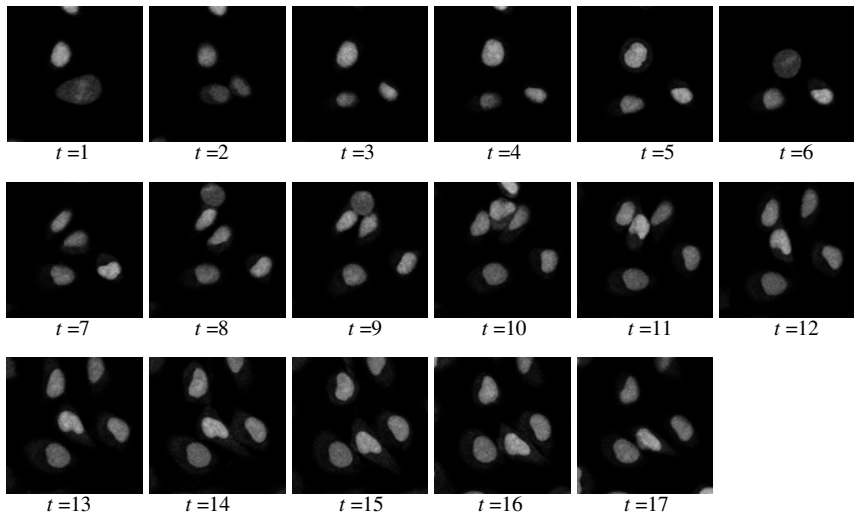


Fig. 1. Time series image data created by Yuki Tsujimura in RIKEN

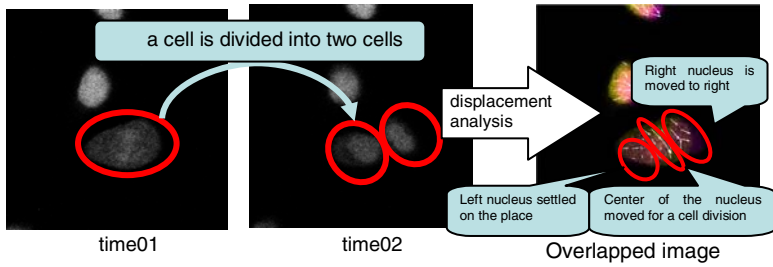


Fig. 2. Example of the time series image analysis. The displacement calculated by viscous fluid registration shows the white arrow in right image named “Overlapped image”. These arrows show a pixel in time01 is moved to a pixel in time02.

2 Viscous Fluid Registration

Viscous fluid registration is the non-rigid image registration method [6] to deformed one image to the other image based on viscous fluid dynamics. It is used for adjust the shape and appearance of two medical images. When this application implemented by C++ language is given a template image, a study image and parameters as input data, it put the output files containing a transformed image and a total displacement file. The deformed image is deformed a template image into a study image. Viscosity fluid registration is calculated three fields the number of repetition times. A force field, a velocity field and a displacement are included in three fields (Fig.3). The number of repetition is defined in an input parameter. In addition, a multi-resolution method is implemented for speedup of calculation and robustness of a registration. Using the example image set that are sample01 (contains in a black square) and sample02 (contains in a white “C”), sample01 is deformed the shape, and the application put the deformed image that is resembles the shape of “C” (Fig.4). The parameters are viscosity=20, rMin=1, rMax=8, calculation number=100. The parameters of the image of the cell are viscosity=60, rMin=1, rMax=8, calculation number=10. There is the explanation of the parameters to 2.1 and 2.2, and we explain the viscous fluid registration algorithm in the section 2.1. In the section 2.2, we explain the two-dimensional viscous image registration workflow using multi-resolution method.

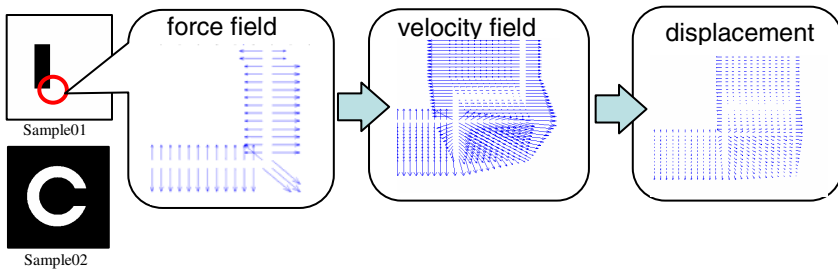


Fig. 3. Example of a force field, a velocity field and a displacement calculated first time

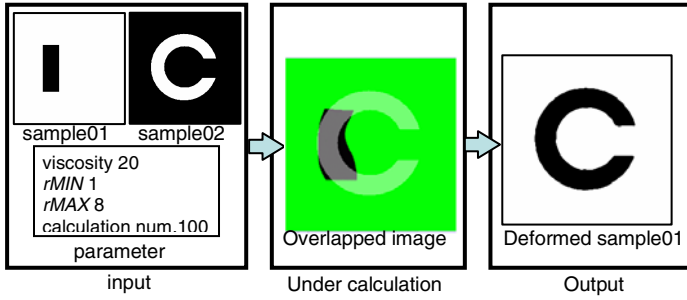


Fig. 4. Example of the viscous fluid registration

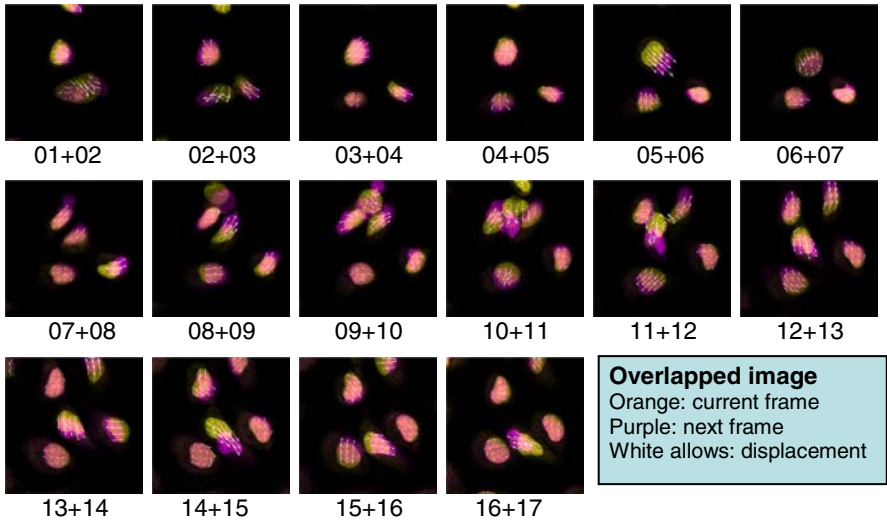


Fig. 5. Overlapped current frame image (orange) on next frame image (purple)

2.1 Viscous Fluid Registration Algorithm

Viscous fluid registration is the image registration method to align two images by deforming one image to the other [7] [8]. A viscous fluid registration is based on the theory of a viscous fluid dynamics. This algorithm contains the calculation of three vector fields: a force field, a velocity field, and a displacement. These vector fields are a set of two-dimensional vectors that are x-axis and y-axis values because we perform two-dimensional viscous fluid registration for all images. We use sum of squared difference as an image similarity measure. Sum of squared difference is used to calculate the force field. The force field is equal to the derivative of image similarity measure. The force field equation can be written as follow:

$$f(x, u(x, t)) = -[T(x - u(x, t)) - S(x)]\nabla T(x - u(x, t)) \tag{1}$$

where $f(\mathbf{x}, \mathbf{u}(\mathbf{x}, t))$ is the force field at point \mathbf{x} and at time t that depends on the displacement $\mathbf{u}(\mathbf{x}, t)$, $T(\mathbf{x})$ is the pixel value at point \mathbf{x} on template image T , and $S(\mathbf{x})$ is the pixel value at point \mathbf{x} on study image S . ∇T is calculated by sobel filter. The velocity field is equal to the convolution of the force field with the Gaussian kernel. The velocity field $\mathbf{v}(\mathbf{x}, \mathbf{u}(\mathbf{x}, t))$ can be written as follow:

$$\mathbf{v}(\mathbf{x}, \mathbf{u}(\mathbf{x}, t)) = G_\sigma (f(\mathbf{x}, \mathbf{u}(\mathbf{x}, t))) \tag{2}$$

The Gaussian kernel G_σ has two parameters: window size w and standard deviation σ . We define this standard deviation σ as the viscosity. The value of the viscosity determines the degree of freedom of the deformation to limit the movement of the pixel at point \mathbf{x} . The window size w is defined as follow:

$$w = \frac{6\sigma}{r} \tag{3}$$

where r is the resolution level that allows for reducing the image size. If the multi-resolution method is not used, the resolution level is set to 1. Using this window size w , the velocity at \mathbf{x} is convolved in 99.74% force field affected to the viscosity at \mathbf{x} according to Gaussian distribution. Displacement $\mathbf{u}(\mathbf{x}, t)$ is calculated as follow:

$$\mathbf{u}(\mathbf{x}, t_{i+1}) = \mathbf{u}(\mathbf{x}, t_i) - (t_{i+1} - t_i) (\mathbf{I} - \nabla \mathbf{u}(\mathbf{x}, t_i)) \mathbf{v}(\mathbf{x}, t_i) \tag{4}$$

where the velocity field $\mathbf{v}(\mathbf{x}, \mathbf{u}(\mathbf{x}, t))$ is equal to the Lagrange derivative of displacement $\mathbf{u}(\mathbf{x}, t)$. Matrix \mathbf{I} is the identity matrix. The time $(t_{i+1} - t_i)$ is calculated by the following equation:

$$MAX(|\mathbf{v}(\mathbf{x}, t_i)|)(t_{i+1} - t_i) = du_{max} \cdot r \tag{5}$$

where du_{max} is the maximal movement allowed in one iteration. We set du_{max} as a fixed value of 0.7 (pixel). If $\mathbf{I} - \nabla \mathbf{u}(\mathbf{x}, t)$ is less than 0.5, the displacement $\mathbf{u}(\mathbf{x}, t)$ is applied to the image, and the displacement $\mathbf{u}(\mathbf{x}, t)$ and time t are initially set to 0, because the displacement is singular for large curved deformation. The total displacement is calculated as follow:

$$\mathbf{u}_{total}(\mathbf{x}, t_i) = \mathbf{u}(\mathbf{x}, t_i) + \mathbf{u}_{total}(\mathbf{x} - \mathbf{u}(\mathbf{x}, t_i), t_{i-1}) \tag{6}$$

2.2 Two-Dimensional Viscous Fluid Registration

Viscous fluid registration algorithm requires four input parameters: the maximum resolution level $rMAX$, the minimum resolution level $rMin$, the calculation number k , and the viscosity σ . We use a multi-resolution method that enables fast and robust image registration. In the multi-resolution method, image I_r in the current resolution level r is used. Image I_r is reduced in size by resolution level r . If resolution level r is 2, the image size is reduced by half. Image I_r is defined as follows:

$$I_r(x_r, y_r) = I_1(r \cdot x_r, r \cdot y_r) \tag{7}$$

$$1 \leq x_r \leq \frac{width}{r} \quad 1 \leq y_r \leq \frac{height}{r} \tag{8}$$

where x_r and y_r are array numbers and they are constrained by (9). The current resolution level r is taken as 1, 2, 4, 8, ..., 2^n . The values of the *width* and *height* are obtained by the image width and image height at resolution level 1.

The viscous fluid registration method is calculated three vector field. In the first step, it calculates the force field value using data from two images. In the second step, it calculates the velocity field using the force field. In the third step, it calculates the displacement using the velocity field. When these three steps are finished, it increments the iteration number and repeats the three steps. If the iteration number is equal to the calculation number, the resolution level is divided by 2, and the iteration number is initialized to 0. After initialization, the three calculation steps are repeated. The resolution level indicates the scale of the image. For example, if the resolution level is 1, it indicates the original image, and if the resolution parameter is 4, it indicates that an image size is quarter. At the beginning of a calculation, the resolution level is set to the maximum resolution level. If the resolution level is equal to the minimum resolution level, it generates output data, including the total displacement and the image deformed using displacement.

2.3 Similarity Maximization Step

To divide one object into two, we suggest adding the similarity maximization step in viscous fluid registration. To confirm this technique, we used two image set in Fig.6, left image set is shown that one circle is divided into two circles, and right image set is shown that one square is divided into two squares. Using original viscous fluid registration, two objects are connected by the thin line (see 40calculation in Fig.7). The force field over and under the line have in the opposite direction (see Fig.8). Because the velocity is calculated by the convolution of the force field, the velocity field on the line is calculated to zero. And the displacement is in proportion to the velocity, displacement on the connected line is zero. After all calculation, there are gaps between the displacement on the connected line and over and under the line. Therefore I introduce a similarity maximization step to approximate the displacement at the point x by near the point x . The similarity maximization step is described as follow:

$$\mathbf{u}_{total}(\mathbf{x}, t_i) = MIN\{(T(\mathbf{x} - \mathbf{u}(\mathbf{x}, t)) - S(\mathbf{x}))^2 \mid \mathbf{x} - 1 \leq \mathbf{x} \leq \mathbf{x} + 1\} \tag{9}$$

This equation maximizes an image similarity measure by minimizing the square error of the difference. This step is added after the calculation of the total displacement \mathbf{u}_{total} in (6). Using this step, deformed result is shown in Fig.9, and the deformed grid image is shown in Fig.10. Fig.10 shows that image edge is enhanced in circle grid image using the new step. This means that fluency of the motion vector is decreased, and the issue of connected line was solved.



Fig. 6. two image set. Left is circle division example. Right is square division example

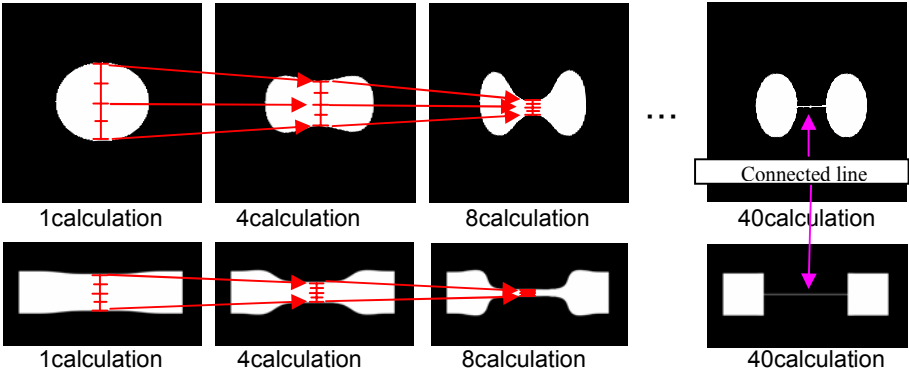


Fig. 7. Connected line is not go disappear using original viscous fluid registration. Top images are circle image result. Bottom images are a square image result.

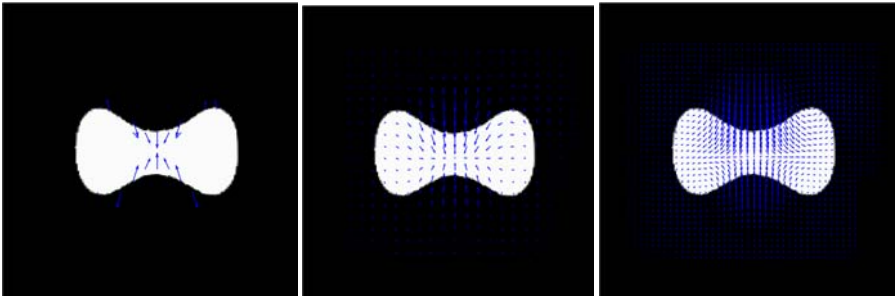


Fig. 8. Force field, velocity field, and displacement at 7calculation in circle image

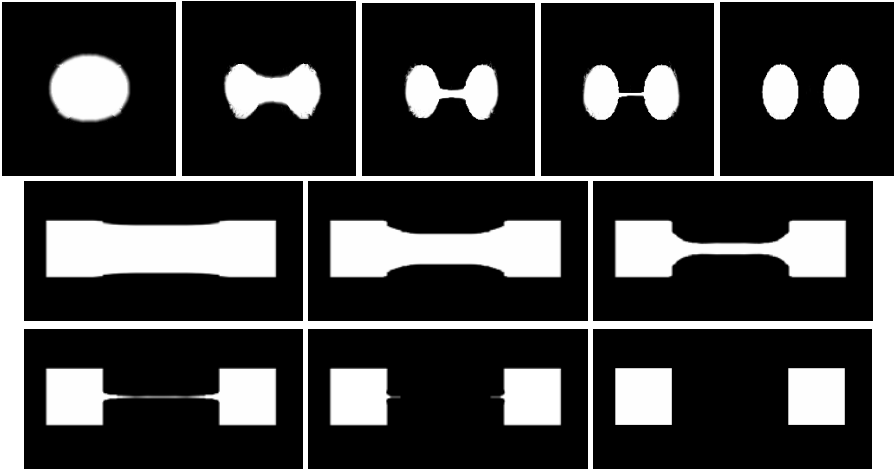


Fig. 9. Disappear connected line adding Similarity maximization step in registration. Top line images are circle image result. center and bottom line images are sqare image result.

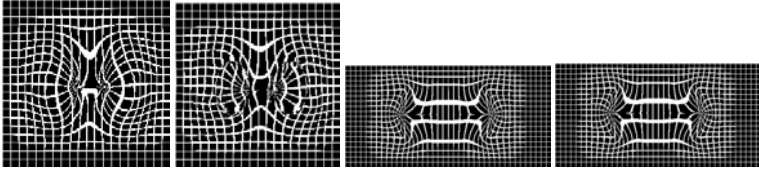


Fig. 10. Deformed grid image. Left two images are deformed circle grid result. Right two images are deformed square grid result. In these set, left is not using step result, right is using step result.

3 Experiment

We experiment the detection of the anaphase in the time series HeLa cell images. Anaphase detection consists of the three steps: image processing, viscous fluid registration and cell identification. We use 17 bitmap images in Fig.1. These images are 350*350 size and 8bit grayscale data created by Yuki Tsujimura in RIKEN. These images contain 76 cells. The resolution of the image is 0.286um in x and y direction and time span between images is 1 hour.

At first we generate a label image by image processing. This step is described in 3.1. In the next step, we perform viscous fluid registration time series images. This second step is described in 3.2. Finally, we perform time to time cell identification by using of the displacement calculated by viscous fluid registration. If corresponding cell is not found after the registration, these cells make pair of a nearest cell in the previous time image. After three steps, we make cell identification in all time, so we can detect a cell in anaphase because these cells are divided into two. This last step and anaphase detection are described in 3.3.

3.1 Image Processing

For the cell identification, we generate a label image from grayscale image using hybrid image processing. This image processing is implemented by the octave language. This method contains five processes: thresholding, median filter, erosion, labeling and dilation(Fig.11). Thresholding process make binary image from grayscale image. Median filter process removes white noise in the cell from binary image. Erosion process makes the cell area shrink. Labeling process add label data on the cell. Dilation process makes the cell area inflated. Dilation and erosion is used for division cell because some cells are touching together. After labeling, label image is displayed by a color using color table. So, label image have pixel values such as 0 on the background and 1, 2, 3... on each independent cell. In this process, 17 label images are created. Regulating parameters at each time are described in Table 1.

Table 1. Parameters list in Image processing at each time

time	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
threshold	30	30	30	30	30	30	30	50	70	70	70	30	30	30	30	60	60
erosion and diration repeat numver	0	0	0	0	0	0	0	5	6	5	5	0	0	0	0	0	0

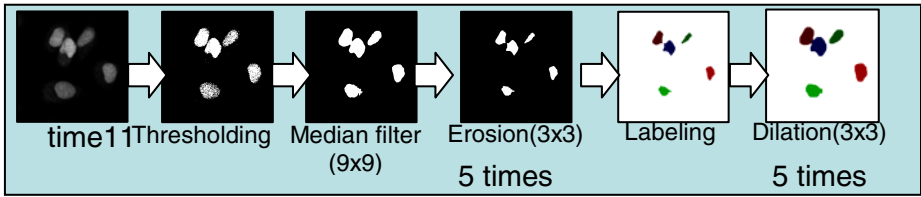


Fig. 11. Image processing flow example at time 10

3.2 Viscous Fluid Registration

In this step, we input two time series images into viscous fluid registration to get the displacement. The label images are not used in this step. The result of the original viscous fluid registration is shown in Fig.12, and the result of the suggested viscous fluid registration added the similarity maximization step is shown in Fig.13. We can find that a cell is divided obviously in suggested technique. As for this technique, limitation of the transformation becomes lax, so the deformed cell shape almost same as input study image.

By the next step, a label image in Fig.11 is deformed using the displacement, and determined the correspondence of the pair of the cells in time t and time $t+1$ image. If the cells in time t and time $t+1$ image exist in the same position, these cells are treated as the correspondence of the pair. And when cell division occurs, cell at time t is divided into two cells at time $t+1$, so we can distinguish the cells in Anaphase by checking label numbers are same value or not.

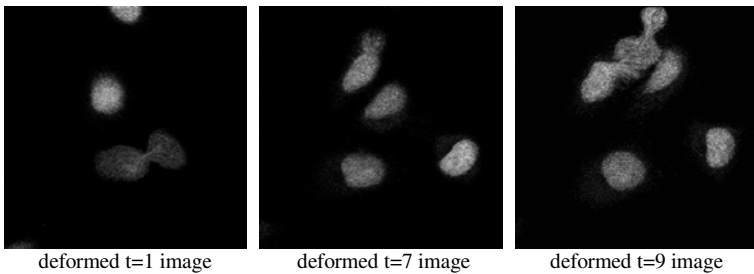


Fig. 12. Original viscous fluid registration result when cell division is occurred

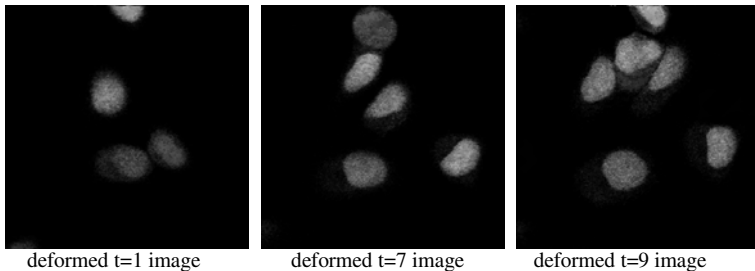


Fig. 13. Suggested viscous fluid registration result when cell division is occurred

3.3 Cell Identification and Anaphase Detection

In this step, cell identification and anaphase detection is performed. The label images and the displacement is used for cell identification. Current time label of cell is taken the correspondence of the next time label. The correspondence of the label is found from the relations of a cell position between deformed current time label and next time label. We calculate histogram of deformed current time label image and next time label image. And next time label value is changed the value of the most overlapped label from deformed current time label image (Fig.14 left). Most overlapped label is found from the histogram. If there are same labels in the next time label image, it found that anaphase is detected. Once anaphase is detected in one label value,

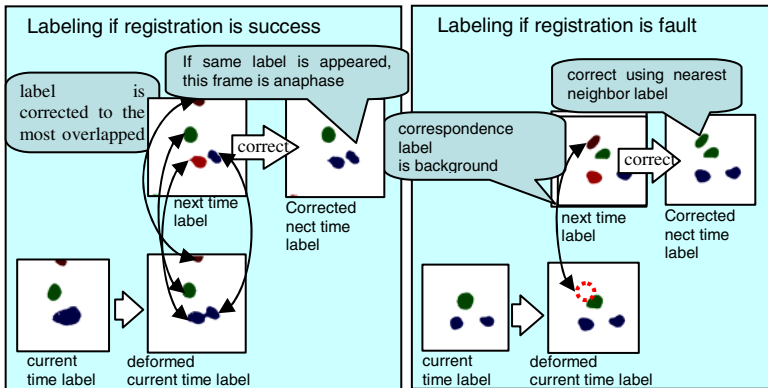


Fig. 14. Procedure of a correcting label

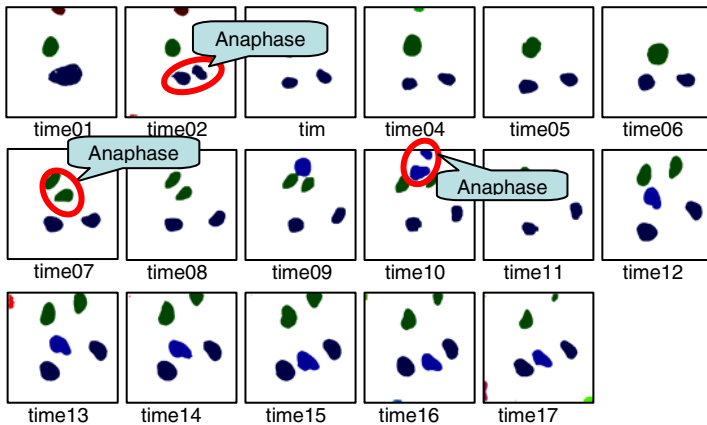


Fig. 15. Result of a detection of anaphase. Dividing cell in anaphase is shown in the red circle

create flag on this label value. This mean anaphase is already found on this label. So, if flag is created, anaphase is not detected on this label. There is the case that registration is fault because the next time cell and deformed cell is not overlapped well. In other words, the corrected label value candidate has the background value (Fig.14 right). When a label candidate became the background, nearest neighbor label is given. In 17 images, two cells are detected by nearest neighbors. This nearest neighbor process is adapted only center of the image. If the center of gravity of a label exists on the area of 40 pixels from boundary of the image, this cell is not regard as in anaphase because this cell moves image boundary to inside and outside. Fig.15 shows the result of corrected all label images.

4 Time Series Image Analysis

To observe the movement of the cell division, all label images are overlapped adapted each image in transparent percentage of 1/17, except for the cell which appeared first is transparency 50% and the cell boundary displays in white and which appeared last time is transparency 0% and the cell boundary displays in black. From this image, it can be observed that three cells are divided into two, and daughter cells are moved to different direction.

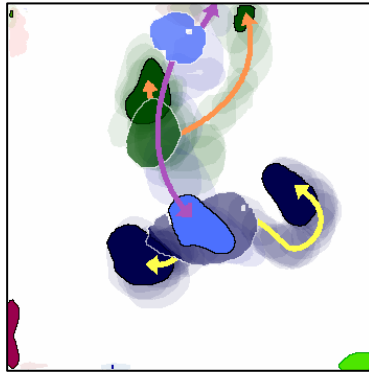


Fig. 16. Overlapped time series label images. This image shows the dividing nucleus movement. The boundary line of nucleus that appears the first is white and that appear the last time is black. A cell color from the second and before last nucleus is semitransparent. Dark blue nucleus is divided and daughter nuclei are moved to left and right (yellow allow), green daughter ones are moved to top and right-top (red allow), and blue ones are moved to top and bottom (purple allow).

5 Conclusion

We can find that divided objects are connected by a thin line using original viscous fluid registration. So we suggest a new viscous fluid registration algorithm for an object division. This algorithm is to add new step in original viscous fluid registration.

We confirmed that connected line which appeared in original viscous fluid registration result is disappeared by adding new step in two sample image set and time series HeLa cell images. We examine the anaphase detection to confirm that new algorithm is worked well. By creating label image and deforming label image using displacement, the cell identification is success 74 in 76 cells. 2 of the remainder cell can be tracked by nearest neighbor processing. A cell in anaphase is detected by cell identification, and daughter cell movement is observed by the overlapped time series label image. Displacement calculated by viscous fluid registration is useful in various ways. Using displacement, we can observe cell movement like Fig.2.

References

1. Zhou, X., Li, F., Yan, J., Wong, S.T.: A novel cell segmentation method and cell phase identification using Markov model. *IEEE transactions on information technology in biomedicine* 13, 152–157 (2009)
2. Chen, X., Zhou, X., Wong, S.T.: Automated segmentation, classification, and tracking of cancer cell nuclei in time-lapse microscopy. *IEEE Transactions on Biomedical Engineering* 53, 762–766 (2006)
3. Dzyubachyk, O., Niessen, W., Meijering, E.: Advanced level-set based multiple-cell segmentation and tracking in time-lapse fluorescence microscopy images. *Biomedical Imaging* 21, 185–188 (2008)
4. Goldman, R.D., Spector, D.L.: *Live Cell Imaging: A Laboratory Manual*. Cold Spring Harbor Laboratory Press (2005)
5. Shu, X., Shaner, N.C., Tarbrough, C.A., Tsien, R.Y., Remington, S.J.: Novel Chromophores and Buried Charges Control Color in mFruits. *Biochemistry* 45, 9639–9647 (2006)
6. Zitova, B., Flusser, J.: Image registration methods: a survey. *Image and Vision Computing* 21, 977–1000 (2003)
7. D'Agostino, E., Maes, F., Vandermeulen, D., Suetens, P.: A Viscous Fluid Model for Multimodal Non-rigid Image Registration Using Mutual Information. In: Dohi, T., Kikinis, R. (eds.) *MICCAI 2002*. LNCS, vol. 2489, pp. 541–548. Springer, Heidelberg (2002)
8. Bro-Nielsen, M., Gramkow, C.: Fast fluid registration of medical images. *Visualization in Biomedical Computing*, 267–276 (1996)

Matrices Representation of Multi Soft-Sets and Its Application

Tutut Herawan^{1,2}, Mustafa Mat Deris¹, and Jemal H. Abawajy³

¹ Faculty of Information Technology and Multimedia
Universiti Tun Hussein Onn Malaysia, Johor, Malaysia

² Department of Mathematics Education
Universitas Ahmad Dahlan, Yogyakarta, Indonesia

³ School of Engineering and Information Technology
Deakin University, Geelong, VIC, Australia

tutut81@uad.ac.id, mmustafa@uthm.edu.my, jemal@deakin.edu.au

Abstract. In previous paper, we introduced a concept of multi-soft sets and used it for finding reducts. However, the comparison of the proposed reduct has not been presented yet, especially with rough-set based reduct. In this paper, we present matrices representation of multi-soft sets. We define AND and OR operations on a collection of such matrices and apply it for finding reducts and core of attributes in a multi-valued information system. Finally, we prove that our proposed technique for reduct is equivalent to Pawlak's rough reduct.

Keywords: Multi-valued information system; Multi-soft sets, Matrices representation, Reducts and core of attributes.

1 Introduction

Soft set theory [1], proposed by Molodtsov in 1999, is a new general method for dealing with uncertain data. In recent years, research on soft set theory has been active, and great progress has been achieved, including the works of theoretical soft set, soft set theory in abstract algebra, parameterization reduction, decision making and forecasting. Let $S = (U, A, V, f)$ be an information system as in [2]. The "standard" soft set deals with a binary-valued information system $S = (U, A, V_{\{0,1\}}, f)$. For a multi valued information system, we introduced a concept of multi soft-sets [3]. The idea is based on a decomposition of a multi-valued information system $S = (U, A, V, f)$, into $|A|$ number of binary-valued information systems $S = (U, A, V_{\{0,1\}}, f)$, where $|A|$ denotes the cardinality of A . Consequently, the $|A|$ binary-valued information systems define *multi-soft sets*, denoted by $(F, A) = \{(F, a_i) : 1 \leq i \leq |A|\}$. In [4], we used the concept of multi-soft sets and AND operation for finding reducts in a multi-valued information system. However, the comparison of the proposed reduct has not been presented yet, especially with rough-set based reduct [2,5]. In this paper, we present the notion of matrices representation of multi-soft sets. We present a definition of AND and OR operations on a collection

of such matrices. Further, we apply the notion of AND operation for finding reducts and core attributes in a multi-valued information system. We prove that the proposed technique of reduct is equivalent with Pawlak’s rough reduct.

The rest of this paper is organized as follows. Section 2 describes the notion of information system. Section 3 describes fundamental concept of soft set theory. Section 4 describes multi soft sets construction in a multi-valued information system. Section 5 describes matrices representation of multi-soft sets, AND and OR operations, as well a simple example. Section 6 describes applications of matrices representation of multi-soft sets for finding reducts and core of attributes. Further, we prove that our proposed technique for reduction is equivalent to Pawlak’s rough reduction. Finally, we conclude and describe future activities of our works in section 7.

2 Information System

Data are often presented as a table, columns of which are labeled by *attributes*, rows by *objects* of interest and entries of the table are *attribute values*. By an *information system*, we mean a 4-tuple (quadruple) $S = (U, A, V, f)$, where $U = \{u_1, u_2, u_3, \dots, u_{|U|}\}$ is a non-empty finite set of objects, $A = \{a_1, a_2, a_3, \dots, a_{|A|}\}$ is a non-empty finite set of attributes, $V = \bigcup_{a \in A} V_a$, V_a is the domain (value set) of attribute a , $f : U \times A \rightarrow V$ is an information function such that $f(u, a) \in V_a$, for every $(u, a) \in U \times A$, called information (knowledge) function. An information system can be intuitively expressed in terms of an information table (refers to Table 1).

Table 1. An information system

U	a_1	\dots	a_k	\dots	$a_{ A }$
u_1	$f(u_1, a_1)$	\dots	$f(u_1, a_k)$	\dots	$f(u_1, a_{ A })$
u_2	$f(u_2, a_1)$	\dots	$f(u_2, a_k)$	\dots	$f(u_2, a_{ A })$
\vdots	\vdots	\ddots	\vdots	\ddots	\vdots
$u_{ U }$	$f(u_{ U }, a_1)$	\dots	$f(u_{ U }, a_k)$	\dots	$f(u_{ U }, a_{ A })$

A relational database may be considered as an information system in which rows are labeled by the objects (entities), columns are labeled by attributes and the entry in row u and column a has the value $f(u, a)$. We note that a each map $f(u, a) : U \times A \rightarrow V$ is a tuple $t_i = (f(u_i, a_1), f(u_i, a_2), f(u_i, a_3), \dots, f(u_i, a_{|A|}))$, for $1 \leq i \leq |U|$. In many applications, there is an outcome of classification that is known. Note that the tuple t is not necessarily associated with entity uniquely (refer to objects 3 and 4 in Table 3). In an information table, two distinct entities could have the same tuple representation (duplicated/redundant tuple), which is *not permissible* in relational databases. Thus, the concepts in information systems are a generalization of the same concepts in relational databases.

This a posteriori knowledge is expressed by one (or more) distinguished attribute called decision attribute; the process is known as supervised learning. An information system of this kind is called a decision system. A *decision system* is an information system of the form $D = (U, A \cup \{d\}, V, f)$, where $d \notin A$ is the decision attribute. The elements of A are called condition attributes. An example of a decision system is given in Table 2.

3 Soft Set Theory

The theory of soft set [1] proposed, by Molodtsov 1999 is a new method for handling uncertain data. The soft set is a mapping from parameter to the crisp subset of universe. From such case, we may see the structure of a soft set can classify the objects into two classes (yes/1 or no/0).

Definition 1. (See [1].) A pair (F, A) is called a soft set over U , where F is a mapping given by

$$F : A \rightarrow P(U).$$

In other words, a soft set over U is a parameterized family of subsets 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 of ε -approximate elements of the soft set. Clearly, a soft set is not a (crisp) set.

Based on the definition of an information system and a soft set, in this section we show that a soft set is a special type of information systems, i.e., a binary-valued information system.

Proposition 2. If (F, A) is a soft set over the universe U , then (F, A) is a binary-valued information system $S = (U, A, V_{\{0,1\}}, f)$.

Proof. Let (F, A) be a soft set over the universe U , we define a mapping

$$F = \{f_1, f_2, \dots, f_n\},$$

where

$$f_i : U \rightarrow V_i \text{ and } f_i(x) = \begin{cases} 1, & x \in F(a_i) \\ 0, & x \notin F(a_i) \end{cases}, \text{ for } 1 \leq i \leq |A|.$$

Hence, if $V = \bigcup_{a_i \in A} V_{a_i}$, where $V_{a_i} = \{0,1\}$, then a soft set (F, A) can be considered as a binary-valued information system $S = (U, A, V_{\{0,1\}}, f)$. □

From Proposition 2, it is easily to understand that a binary-valued information system can be represented as a soft set. Thus, we can make a one-to-one correspondence between (F, E) over U and $S = (U, A, V_{\{0,1\}}, f)$.

Definition 3. (See [6].) The class of all value sets of a soft set (F, E) is called value-class of the soft set and is denoted by $C_{(F,E)}$.

4 Multi-Soft Sets Construction in Information Systems

The idea of multi-soft sets is based on a decomposition of a multi-valued information system $S = (U, A, V, f)$, into $|A|$ number of binary-valued information systems $S = (U, A, V_{\{0,1\}}, f)$, where $|A|$ denotes the cardinality of A . Consequently, the $|A|$ binary-valued information systems define *multi-soft sets* $(F, A) = \{(F, a_i) : 1 \leq i \leq |A|\}$.

4.1 Decomposition of Multi-valued Information Systems

The decomposition of $S = (U, A, V, f)$ is based on decomposition of $A = \{a_1, a_2, \dots, a_{|A|}\}$ into the disjoint-singleton attribute $\{a_1\}, \{a_2\}, \dots, \{a_{|A|}\}$. Here, we only consider for complete information systems. Let $S = (U, A, V, f)$ be an information system such that for every $a \in A$, $V_a = f(U, A)$ is a finite non-empty set and for every $u \in U$, $|f(u, a)| = 1$. For every a_i under i^{th} -attribute consideration, $a_i \in A$ and $v \in V_a$, we define the map $a_v^i : U \rightarrow \{0,1\}$ such that $a_v^i(u) = 1$ if $f(u, a) = v$, otherwise $a_v^i(u) = 0$. The next result, we define a binary-valued information system as a quadruple $S^i = (U, a_i, V_{\{0,1\}}, f)$. The information systems $S^i = (U, a_i, V_{\{0,1\}}, f)$, $1 \leq i \leq |A|$ is referred to as a decomposition of a multi-valued information system $S = (U, A, V, f)$ into $|A|$ binary-valued information systems, as depicted in Figure 1. Every information system $S^i = (U, a_i, V_{a_i}, f)$, $1 \leq i \leq |A|$ is a deterministic information system since for every $a \in A$ and for every $u \in U$, $|f(u, a)| = 1$ such that the structure of a multi-valued information system and $|A|$ number of binary-valued information systems give the same value of attribute related to objects.

4.2 Multi-soft Sets in Information Systems

Based on the notion of a decomposition of a multi-valued information system in the previous sub-section, in this sub-section we present the notion of multi-soft set representing multi-valued information systems. Let $S = (U, A, V, f)$ be a multi-valued information system and $S^i = (U, a_i, V_{a_i}, f)$, $1 \leq i \leq |A|$ be the $|A|$ binary-valued information systems. From Proposition 2, we have

$$\begin{aligned}
 S = (U, A, V, f) &= \begin{cases} S^1 = (U, a_1, V_{\{0,1\}}, f) & \Leftrightarrow (F, a_1) \\ S^2 = (U, a_2, V_{\{0,1\}}, f) & \Leftrightarrow (F, a_2) \\ \vdots & \vdots \\ S^{|A|} = (U, a_{|A|}, V_{\{0,1\}}, f) & \Leftrightarrow (F, a_{|A|}) \end{cases} \\
 &= ((F, a_1), (F, a_2), \dots, (F, a_{|A|}))
 \end{aligned}$$

We define $(F, A) = ((F, a_1), (F, a_2), \dots, (F, a_{|A|}))$ as a *multi-soft sets* over universe U representing a multi-valued information system $S = (U, A, V, f)$.

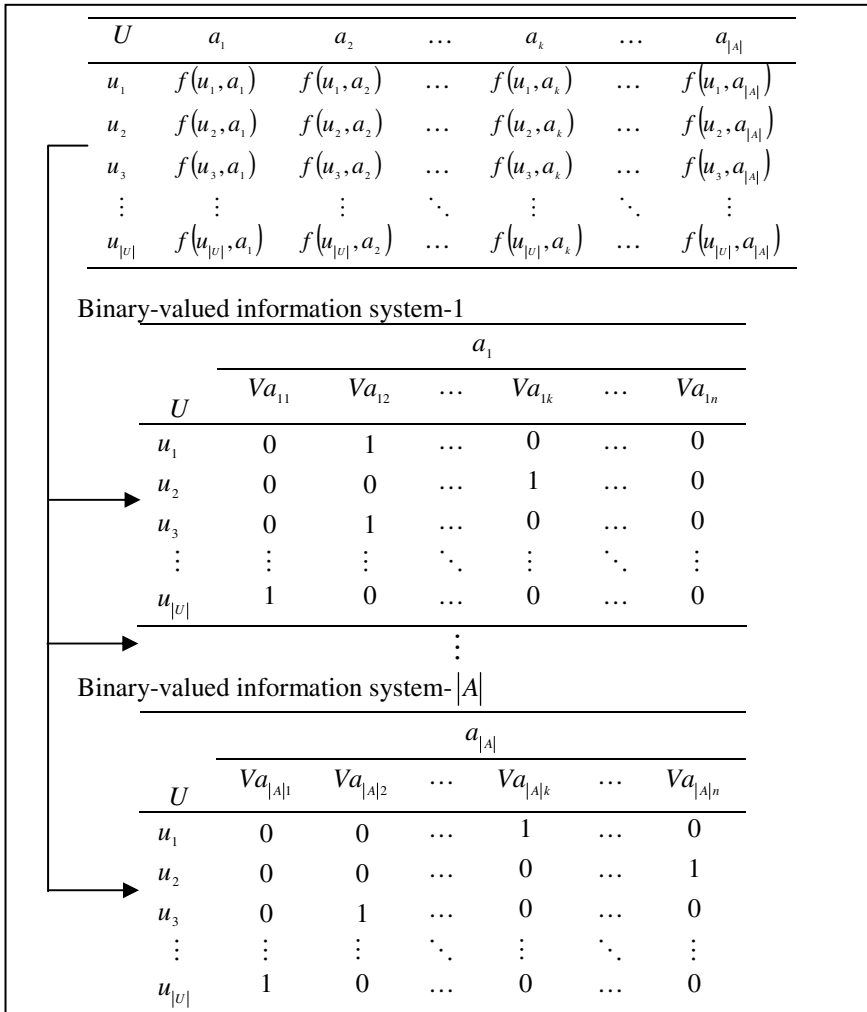


Fig. 1. A decomposition of information systems

Example 4. *Hiring* [7]: an example of a multi-valued information system is presented in Table 2. In Table 2, six students are characterized by four conditional attributes; Diploma, Experience, French, Reference and one decision attribute; Decision.

Let $A = \{\text{Diploma, Experience, French, Reference, Decision}\}$. Therefore, the multi-soft set representing Table 2 is given in Figure 2. Note that the class value of every soft set is a partition of U .

Table 2. *Hiring*: an information system from [7]

Student	Diploma	Experience	French	Reference	Decision
1	MBA	Medium	Yes	Excellent	Accept
2	MBA	Low	Yes	Neutral	Reject
3	MCE	Low	Yes	Good	Reject
4	MSc	High	Yes	Neutral	Accept
5	MSc	Medium	Yes	Neutral	Reject
6	MSc	High	Yes	Excellent	Accept
7	MBA	High	No	Good	Accept
8	MCE	Low	No	Excellent	Reject

$$(F, A) = \left(\begin{array}{l} \{MBA = \{1,2,7\}, MCE = \{3,8\}, MSc = \{4,5,6\}\}, \\ \{Medium = \{1,5\}, Low = \{2,3,8\}, High = \{4,6,7\}\}, \\ \{Yes = \{1,2,3,4,5,6\}, No = \{7,8\}\}, \\ \{Excellent = \{1,6,8\}, Neutral = \{2,4,5\}, Good = \{3,7\}\}, \\ \{Accept = \{1,4,6,7\}, Reject = \{2,3,5,8\}\} \end{array} \right)$$

Fig. 2. Multi soft-sets representing Table 2

5 Matrices Representation of Multi-soft Sets

5.1 Matrix Representation

The concept of matrix representing multi-soft sets is described here. In this subsection, the notation (F, A) represents a multi-soft set over universe U of a multi-valued information system $S = (U, A, V, f)$.

$$(F, A) \Leftrightarrow \left\{ \begin{array}{l} (F, a_1) \Leftrightarrow S^1 = (U, a_1, V_{\{0,1\}}, f) \Leftrightarrow M_{a_1} \\ (F, a_2) \Leftrightarrow S^2 = (U, a_2, V_{\{0,1\}}, f) \Leftrightarrow M_{a_2} \\ \vdots \\ (F, a_{|A|}) \Leftrightarrow S^{|A|} = (U, a_{|A|}, V_{\{0,1\}}, f) \Leftrightarrow M_{|A|} \end{array} \right.$$

Definition 5. Matrix $M_{a_i}, 1 \leq i \leq |A|$ is called matrix representation of soft set (F, a_i) over universe U . The dimension of matrices is defined by $\dim(M_{a_i}) = |U| \times |Va_i|$. All entries of $M_{a_i} = [a_{ij}]$ is belong to a set $\{0,1\}$, where

$$a_{ij} = \begin{cases} 0, & \text{if } |f(u, \alpha)| = 0 \\ 1, & \text{if } |f(u, \alpha)| = 1 \end{cases}, \quad 1 \leq i \leq |U|, 1 \leq j \leq |Va_i|, u \in U \text{ and } \alpha \in Va_i.$$

The collection of all matrices representing (F, A) is denoted by \mathcal{M}_A , i.e.,

$$\mathcal{M}_A = \{M_{a_i} : 1 \leq i \leq |A|\}.$$

Definition 6. Let $M_{a_i} \in \mathcal{M}_A$ be a matrix representation of a soft set (F, a_i) over U . The value-class of M_{a_i} , i.e., class of all value sets of M_{a_i} , denoted $C_{M_{a_i}}$ is defined by

$$C_{M_{a_i}} = \left\{ \left\{ u : |f(u, \alpha_1)| = 1 \right\}, \dots, \left\{ u : |f(u, \alpha_{|V_{a_i}|})| = 1 \right\} \right\},$$

where $1 \leq i \leq |V_{a_i}|$, $u \in U$ and $\alpha \in V_{a_i}$.

Clearly $C_{M_{a_i}} \subseteq P(U)$.

Example 7. The collection of matrices representing (F, A) , is given as

$$\mathcal{M}_A = \{M_{a_1}, M_{a_2}, M_{a_3}, M_{a_4}, M_{a_5}\},$$

where

$$M_{\text{Diploma}} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, M_{\text{Experience}} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}, M_{\text{French}} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix},$$

$$M_{\text{Reference}} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \text{ and } M_{\text{Decision}} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

where

$$C_{M_{\text{Diploma}}} = \{\{1,2,7\}, \{3,8\}, \{4,5,6\}\}, C_{M_{\text{Experience}}} = \{\{1,5\}, \{2,3,8\}, \{4,6,7\}\},$$

$$C_{M_{\text{French}}} = \{\{1,2,3,4,5,6\}, \{7,8\}\}, C_{M_{\text{Reference}}} = \{\{1,6,8\}, \{2,4,5\}, \{3,7\}\},$$

$$\text{and } C_{M_{\text{Decision}}} = \{\{1,4,6,7\}, \{2,3,5,8\}\},$$

respectively.

5.2 AND and OR Operations in a Collection of Matrices

The AND and OR operation of the matrices are described in this sub-section.

Definition 8. Let $M_{a_i} = [a_{kl}]$, $1 \leq k \leq |U|$, $1 \leq l \leq |V_{a_i}|$ and $M_{a_j} = [a_{mn}]$, $1 \leq m \leq |U|$, $1 \leq n \leq |V_{a_j}|$ be two matrices in \mathcal{M}_A . The "AND" operation $M_{a_i} \text{AND} M_{a_j}$ of matrices M_{a_i} and M_{a_j} is defined as follows

$$M_{a_i} \text{AND} M_{a_j} = M_{a_{ij}} = [a_{pq}] \text{ with } \dim(M_{a_{ij}}) = |U| \times (|V_{a_i}| \times |V_{a_j}|),$$

where

$$a_{p1} = \min\{a_{k1}, a_{m1}\}, a_{p2} = \min\{a_{k1}, a_{m2}\}, \dots, a_{p(|V_{a_i}| \times |V_{a_j}|)} = \min\{a_{k|V_{a_i}|}, a_{m|V_{a_j}|}\}.$$

Example 9. Let $M_{\text{Diploma}}, M_{\text{Experience}} \in \mathcal{M}_E$, from Definition 8, we have

$$M_{\text{Diploma}} \text{AND} M_{\text{Experience}} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \text{AND} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix},$$

where

$$C_{M_{\text{Diploma}} \text{AND} M_{\text{Experience}}} = \{\{1\}, \{2\}, \{3,8\}, \{4,6\}, \{5\}, \{7\}\}.$$

Definition 10. Let $M_{a_i} = [a_{kl}]$ and $M_{a_j} = [b_{mn}]$ be two matrices in \mathcal{M}_A . The "OR" operation $M_{a_i} \text{OR} M_{a_j}$ of matrices M_{a_i} and M_{a_j} is defined as follows

$$M_{a_i} \text{OR} M_{a_j} = M_{a_{ij}} = [c_{pq}] \text{ with } \dim(M_{a_{ij}}) = |U| \times (|V_{a_i}| \times |V_{a_j}|),$$

where

$$a_{p1} = \max\{a_{k1}, b_{m1}\}, a_{p2} = \max\{a_{k1}, b_{m2}\}, \dots, a_{p(|V_{a_i}| \times |V_{a_j}|)} = \max\{a_{k|V_{a_i}|}, b_{m|V_{a_j}|}\}.$$

Example 11. Let $M_{\text{Diploma}}, M_{\text{Experience}} \in \mathcal{M}_E$, from Definition 10, we have

$$M_{\text{Diploma}} \text{OR} M_{\text{Experience}} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \text{OR} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 1 & 0 \end{bmatrix},$$

where

$$C_{M_{\text{Diploma}} \text{ORM} M_{\text{Experience}}} = \{\{1\}, \{2\}, \{3,8\}, \{4,6\}, \{5\}, \{7\}\}.$$

We can see that, since $\dim(M_{a_j}) = |U| \times |Va_i \cup Va_j|$, thus \mathcal{M}_A is not closed under AND and OR operations, thus $(\mathcal{M}_A, \text{AND})$ and $(\mathcal{M}_A, \text{OR})$ are not groupoids.

Proposition 12. *Let M_{a_i} , M_{a_j} and M_{a_k} be matrices in \mathcal{M}_A . The following properties are holds.*

- a. $M_{a_i} \text{AND} M_{a_i} = M_{a_i}$, *idempotent*
- b. $M_{a_i} \text{ORM} M_{a_i} = M_{a_i}$, *idempotent*
- c. $M_{a_i} \text{AND} (M_{a_j} \text{AND} M_{a_k}) = (M_{a_i} \text{AND} M_{a_j}) \text{AND} M_{a_k}$ *associative*
- d. $M_{a_i} \text{OR} (M_{a_j} \text{OR} M_{a_k}) = (M_{a_i} \text{OR} M_{a_j}) \text{OR} M_{a_k}$ *associative*

AND and OR operations are not commutative.

6 Application

In this section, we present an application of matrices representation of multi-soft sets for dimensionality reduction and finding core attributes.

6.1 Reduct and Core

In [4], we presented the application of multi-soft sets for finding reducts. The proposed approach is based on AND operation in multi-soft sets [3]. In this subsection, we explore the applicability of matrices representation for finding reducts and core attributes. We show that the reducts obtained are equivalent with that in [4]. The notion of reduct in [4] is redefined and given as follow.

Definition 13. *Let $(F, A) = ((F, a_i); 1 \leq i \leq |A|)$ be multi-soft set over U representing a multi-valued information system $S = (U, A, V, f)$. A set of attributes $B \subseteq A$ is called a reduct for A if*

$$C_{F(b_1 \times \dots \times b_{|B|})} = C_{F(a_1 \times \dots \times a_{|A|})} \text{ and } C_{F(b_1 \times \dots \times b_{|B|})} \neq C_{F(a_1 \times \dots \times a_{|A|})}, \quad B^* \subset B,$$

where

$$(F, a_i \times a_j) = (F, a_i) \text{AND} (F, a_j).$$

A core of A is defined as

$$\text{CORE}(A) = \bigcap \text{RED}(A),$$

where $\text{RED}(A)$ is the set of all reducts of A .

Let, given $B = \{a_1, a_2, a_3\}$ and $C = \{a_3, a_4\}$, then we have

$$M_{a_1} \text{ANDM}_{a_2} \text{ANDM}_{a_3} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

and

$$M_{a_3} \text{ANDM}_{a_4} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

where

$$C_{M_{a_1} \text{ANDM}_{a_2} \text{ANDM}_{a_3}} = \{\{1\}, \{2\}, \{3,4\}, \{5\}\} \text{ and } C_{M_{a_3} \text{ANDM}_{a_4}} = \{\{1\}, \{2\}, \{3,4\}, \{5\}\},$$

respectively. Thus, we have $\{a_1, a_2, a_3\}$ and $\{a_3, a_4\}$ are reducts of A . Furthermore, core is

$$\{a_1, a_2, a_3\} \cap \{a_3, a_4\} = \{a_3\}.$$

6.2 Equivalence with Pawlak’s Rough Reduct

In rough set theory [2,5], the aim of dimensionality reduction is to keep only those attributes that preserve the indiscernibility relation and, consequently, set approximation. The remaining attributes are redundant since their removal does not worsen the classification. There are usually several such subsets of attributes and those which are minimal are called reducts. While computing equivalence classes is straightforward, finding all reducts is NP-hard. This reduction is different in soft set theory. The existing techniques of soft parameterization are still based on a binary information system [8,9]. Soft parameterization reduction is obtained based on the optimal and sub-optimal choice related to each object. Thus, the idea of reduction under rough set theory generally cannot be applied directly in reduction under soft set theory. In this sub-section, we show that our proposed technique on dimensionality reduction in a multi-valued information system is equivalent with that in [2,5]. Let, firstly, we recall the notion of rough reduction as follow.

Let $S = (U, A, V, f)$ be an information system and let B be any subset of A . Two elements $x, y \in U$ are said to be B -indiscernible (indiscernible by the set of attribute B) if and only if $f(x, a) = f(y, a)$, for every $a \in B$. Obviously, every subset of A induces unique indiscernibility relation. Notice that, an indiscernibility relation is an *equivalence relation*. The *partition* of U induced by $IND(B)$ in $S = (U, A, V, f)$ is denoted by U/B .

Definition 16. Let $S = (U, A, V, f)$ be an information system and let B be any subset of A . A subset $B \subseteq A$ is called a *reduct* of A if B satisfies the following conditions

- a. $U/B = U/A$
- b. $U/(B - \{b\}) \neq U/A, \forall b \in B$

A *core* of A is defined as

$$\text{CORE}(A) = \bigcap \text{RED}(A),$$

where $\text{RED}(A)$ is the set of all reducts of A .

Proposition 17. Soft reduction in Definition 13 is equivalent with rough reduction in Definition 16.

Proof. (\Rightarrow) It is clear that $(F, A) = (U, A, V, f)$, thus $C_{F(a_i)} = U/a_i, 1 \leq i \leq |A|$. From the hypothesis, since B is a reduct of A , then

$$\begin{aligned} C_{F(b_1 \times \dots \times b_{|B|})} &= C_{F(a_1 \times \dots \times a_{|A|})} \\ (F, b_1) \text{AND} \dots \text{AND} (F, b_{|B|}) &= (F, a_1) \text{AND} \dots \text{AND} (F, a_{|A|}) \\ C_{F(b_1)} \cap \dots \cap C_{F(b_{|B|})} &= C_{F(a_1)} \cap \dots \cap C_{F(a_{|A|})} \\ \bigcap_{1 \leq i \leq |B|} C_{F(b_i)} &= \bigcap_{1 \leq i \leq |A|} C_{F(a_i)} \\ \bigcap_{1 \leq i \leq |B|} U/b_i &= \bigcap_{1 \leq i \leq |A|} U/a_i \\ U/B &= U/A \end{aligned}$$

Say, $B^* = B - \{b\}$, where $b \in B$, then

$$\begin{aligned} C_{F(b_1 \times \dots \times b_{|B|})} &\neq C_{F(a_1 \times \dots \times a_{|A|})} \\ U/(B - \{b\}) &\neq U/A \end{aligned}$$

(\Leftarrow) Obvious. □

Example 18. From Table 2, let

$$B = \{\text{Diploma, Experience, French, Reference}\}.$$

Using rough set approach, we have

$$U/B = \{\{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \{6\}, \{7\}, \{8\}\}.$$

Therefore, reduct of B is $B^* = \{\text{Experience, Reference}\}.$

Therefore, reduct of B is $B^* = \{\text{Experience, Reference}\}$. It is shown that, using rough and matrices approaches, reduct obtained is equivalent.

7 Conclusion

The “standard” soft set deals with a binary-valued information system. For a multi-valued information system, the notion of multi-soft sets has been proposed. The idea of multi-soft sets is based on a decomposition of a multi-valued information system into binary-valued information systems. In this paper, we have presented the notion of matrices representation of such multi-soft sets. The AND and OR operations in a collection of matrices sets also presented. The AND and OR operations satisfy idempotent and associative properties, neither closed nor commutative. Further, we have proven that AND operation can be applied for finding reducts. Finally, we have proven that our proposed technique provide the same results with that rough reduct. For future research, we further elaborate our approach for finding functional, identity and attribute dependencies in information systems. Since, the calculation of matrices is complex, we will use MATLAB to solve this problem. We believe that our proposed approach can be used for feature selection and data cleansing in multi-valued information system under soft set theory.

Acknowledgement

This work was supported by the FRGS under the Grant No. Vote 0402, Ministry of Higher Education, Malaysia.

References

1. Molodtsov, D.: Soft Set Theory-First Results. *Computers and Mathematics with Applications* 37, 19–31 (1999)
2. Pawlak, Z., Skowron, A.: Rudiments of Rough sets. *Information Sciences* 177(1), 3–27 (2007)
3. Herawan, T., Mustafa, M.D.: On Multi-soft Sets Construction in Information Systems. In: Huang, D.-S., Jo, K.-H., Lee, H.-H., Kang, H.-J., Bevilacqua, V. (eds.) ICIC 2009. LNCS (LNAI), vol. 5755, pp. 101–110. Springer, Heidelberg (2009)
4. Herawan, T., Rose, A.N.M., Mustafa, M.D.: Soft Set Theoretic Approach for Dimensionality Reduction. In: *Communication of Computer and Information Sciences*, vol. 64, pp. 180–187. Springer, Heidelberg (2009)
5. Pawlak, Z.: *Rough Sets: A Theoretical Aspect of Reasoning about Data*. Kluwer Academic Publisher, Dordrecht (1991)
6. Maji, P.K., Biswas, R., Roy, A.R.: Soft Set Theory. *Computers and Mathematics with Applications* 45, 555–562 (2003)
7. Komorowski, J., Ohn, A.: Modeling prognostic power of cardiac tests using rough sets. *Artificial Intelligence in Medicine* 15, 167–191 (1999)
8. Chen, D., Tsang, E.C.C., Yeung, D.S., Wang, X.: The Parameterization Reduction of Soft Sets and its Applications. *Computers and Mathematics with Applications* 49, 757–763 (2005)
9. Zhi, K., Gao, L., Wang, L., Li, S.: The Normal Parameter Reduction of Soft Sets and Its Algorithm. *Computers and Mathematics with Applications* 56, 3029–3037 (2008)

Clustering Analysis of Water Quality for Canals in Bangkok, Thailand

Sirilak Areerachakul and Siripun Sanguansintukul

Department of Mathematics, Faculty of Science, Chulalongkorn University,
Bangkok, Thailand
Sirilak.Ar@Student.chula.ac.th, siripun.s@chula.ac.th

Abstract. Two clustering techniques of water quality for canals in Bangkok were compared: K-means and Fuzzy c-means. The result illustrated that K-means has a better performance. As a result, K-means cluster was used to classify 24 canals of 344 records of surface water quality within Bangkok; the capital city of Thailand. The data was obtained from the Department of Drainage and Sewerage Bangkok Metropolitan Administration during 2005-2008. Water samples were collected and analyzed on 13 different parameters: temperature, pH value (pH), hydrogen sulfide (H_2S), dissolved oxygen (DO), biochemical oxygen demand (BOD), chemical oxygen demand (COD), substance solid (SS), total kjeldahl nitrogen (TKN), ammonia nitrogen (NH_3N), nitrite nitrogen (NO_2N), nitrate nitrogen (NO_3N), total phosphorous (T-P) and total coliform. The data were analyzed and clustered. The results of cluster analysis divided the canals into five clusters. The information from clustering could enhance the understanding of surface water usage in the area. Additionally, it can provide the useful information for better planning and watershed management of canals in Bangkok.

Keywords: surface water quality, watershed management, K-means clustering.

1 Introduction

Water will become the major constraining resource for sustainable development of large areas in the world. Water quality is one of the main characteristics of a river, even when its purpose is other than human water supply [4]. Water quality in the superficial waters has started to degenerate as a result of waste-water being let go to the receiving ground and surface water without any controls [5].

Bangkok is the capital city, as well as, the economic center of Thailand. It's activities, which include commercial, industrial and service have caused the expansion of the city and its population to accumulate environmental pollution to the point that nature cannot cope with the pollution loading, especially for water quality [21].

At present, the canal water pollution is very severe as the canals still act as sewers for direct discharge of waste-water in most areas of Bangkok even though there is already a requirement for large buildings to treat waste-water and houses must have at least septic tanks. However, other waste-waters are discharged without treatment to the public sewers which drain into the canals [22].

Multivariate statistical methods have been widely applied to investigate environmental phenomena, such as Laaksoharju *et al.* [11] who did hydrochemical evaluation using a multivariate mathematical tool. Anazawa *et al.* [6] employed factor analysis (FA); two major factors, to describe the chemical behaviors parameters. Güler and Thyne [2] applied HCA analysis to classify recharged area waters. Cluster analysis, principle component analysis/factor analysis, time series analysis, self-organizing maps and classification and regression trees (CART) strategy, etc., are powerful tools for deriving useful information from complicated data about water quality studies [1],[7],[9],[12],[14],[15]. The multivariate statistical methods have extensive applications in characterization and evaluation of surface water quality and are useful for evaluation of temporal and spatial variations caused by natural and anthropogenic factors[16],[17].

The purpose of this paper is to compare two clustering techniques: K-means and Fuzzy c-means. These two methodologies have been selected because they are the most classical, popular and efficient methodologies for clustering surface water quality [8],[19]. The better clustering approach will be used to summarize and analyze surface water quality of canals in Bangkok.

2 Data and Site Description

In this study, water quality data are provided by the Department of Drainage and Sewerage, Bangkok Metropolitan Administration, during the years 2005-2008. There are 24 canals of 344 records of sites. Each record consists of 13 parameters, namely: temperature, pH value (pH), hydrogen sulfide (H₂S), dissolved oxygen (DO), biochemical oxygen demand (BOD), chemical oxygen demand (COD), substance solid (SS), total kjeldahl nitrogen (TKN), ammonia nitrogen (NH₃N), nitrite nitrogen (NO₂N), nitrate nitrogen (NO₃N), total phosphorous (T-P) and total coliform. Table 1 shows unit of surface water quality parameters[20].

Table 1. List of surface water quality parameters

Name of Parameters	Unit of measurement
Temperature	Celsius
pH value	Standard Units
Hydrogen Sulfide	Milligrams per liter
Dissolved Oxygen	Milligrams per liter
Biochemical Oxygen Demand	Milligrams per liter
Chemical Oxygen Demand	Milligrams per liter
Substance Solid	Milligrams per liter
Total Kjeldahl Nitrogen	Milligrams per liter
Ammonia Nitrogen	Milligrams per liter
Nitrite Nitrogen	Milligrams per liter
Nitrate Nitrogen	Milligrams per liter
Total Phosphorous	Milligrams per liter
Total Coliform	Most Probable Number per 100 Milliliter

The network of canals (shown in Fig.3) is important for the daily life of the people in Bangkok. These canals are used for consumption, transportation and recreation. Therefore, the rapid growth of industry, condominium, high-rise and low-rise building, and other infrastructure, have a significant effect on the canal water. The major impact can be classified into 3 types: [21]

- Impact to tourism activities (shown in Fig.1): Bangkok is the capital city and the center of tourism. Therefore, deterioration of canal water has direct impact to tourism activity. As pollution in canals is in the inner area of Bangkok, it certainly gives negative impression to the tourists who travel and stay in Bangkok.
- Impact to aquatic life (shown in Fig.2): Basically, water pollution causes deaths of aquatic creatures either by its toxicity or reducing dissolved oxygen concentration. The toxicity may come from high concentration of sulfide or ammonia. It is observed that when water is polluted there are usually fall in the total number of species of organisms generally more sensitive than fishes, a change in the type of species present and a change in the numbers of individuals of each species in the water. In severe cases, where dissolved oxygen concentration becomes zero, all aquatic creatures die and this occurs in the highly polluted canals in Bangkok.
- Impact to public health: poor sanitary conditions prevail in many parts of Bangkok. Records of some water related diseases generally associated with sanitary conditions are discussed in the Annual Epidemiological Surveillance Report of the Epidemiological Subdivision, Disease Control Division Department of Health [23]. In case of acute diarrhea, about 45,000 and 50,000 cases were found in 1996 and 1997 respectively.



Fig. 1. The impact to tourism activities



Fig. 2. The impact to aquatic life

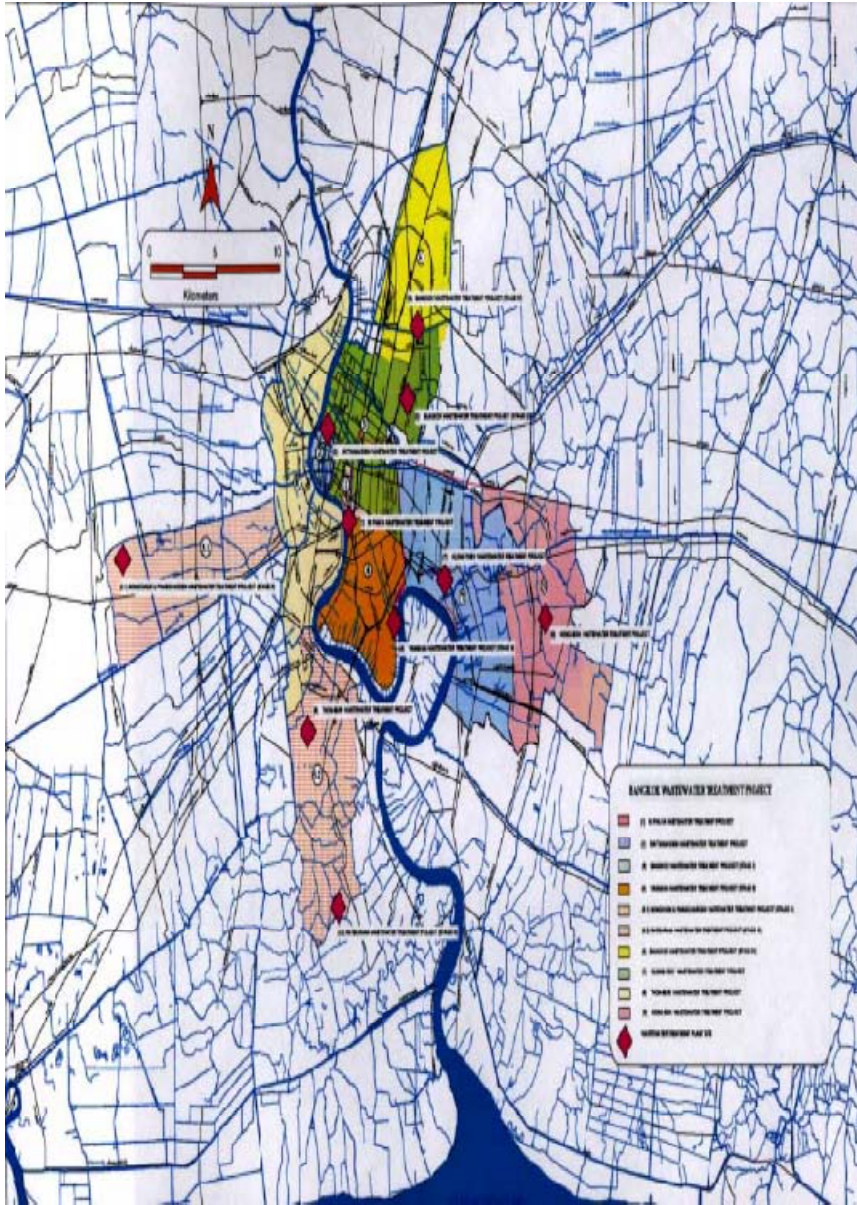


Fig. 3. The network of canals in Bangkok, Thailand

3 Methodology

Clustering is an important data mining technique that has a wide range of applications in many areas like biology, medicine, market research, image processing, and geographical information systems, among others [24].

In this study, the water qualities of canals in Bangkok were grouped according to their characteristic forming clusters. The clustering process was carried out using a K-means algorithm and Fuzzy c-means algorithm.

3.1 K-Means Algorithm

The K-means algorithm is a cluster analysis algorithm used as a partitioning method, and was developed by MacQueen in 1967 [3]. The K-means algorithm defines a random cluster centroid according to the initial parameter [10]. Each consecutive case is added to the cluster according to the proximity between the mean value of the case and the cluster centroid. The clusters are then re-analyzed to determine the new centroid point. This procedure is repeated for each data object.

The algorithm is composed of the following steps: [24]

1. Place K point into the space represented by the objects that are being clustered. These points represent initial group centroids.
2. Assign each object to the group that has the closest centroids.
3. When all objects have been assigned, recalculate the positions of the K centroids.
4. Repeat step 2 and 3 until the centroids no longer move. This produces a separation of the objects into groups from which the metric to be minimized can be calculated.

The main idea is to define k centroids, one for each cluster. The aim of that algorithm is to find cluster centroids for each group. The algorithm minimizes a dissimilarity function which is given in equation (1) [3].

$$J = \sum_{j=1}^k \sum_{i=1}^n \|x_i - c_j\|^2 \quad (1)$$

where x_i is the i^{th} data point, c_j is the j^{th} cluster centroid, n is the number of data

$\|x_i - c_j\|$ is a chosen distance measure between a data point x_i and the cluster centroids c_j .

3.2 Fuzzy c-Means Algorithm

Fuzzy c-means (FCM) is a method of clustering which allows one piece of data to belong to two or more clusters. This method is frequently used in pattern recognition [18]. It is based on minimization of the following objective function:

$$J_m = \sum_{i=1}^N \sum_{j=1}^C u_{ij}^m \|x_i - c_j\|^2 \quad (2)$$

where m is any real number greater than 1, u_{ij} is the degree of membership of x_i in the cluster j , x_i is the i^{th} of d -dimensional measured data, c_j is the d -dimensional center of the cluster and $\|*\|$ is any norm expressing the similarity between any measured data and the center.

Fuzzy partitioning is carried out through an iterative optimization of the objective function shown above with the update of membership u_{ij} and the c_j cluster center by:

$$u_{ij} = \frac{1}{\sum_{k=1}^c \left(\frac{\|x_i - c_j\|}{\|x_i - c_k\|} \right)^{\frac{2}{m-1}}} \tag{3}$$

$$c = \frac{\sum_{i=1}^N u_{ij}^m x_i}{\sum_{i=1}^N u_{ij}^m} \tag{4}$$

This iteration will stop when

$$\max_{ij} \{ |u_{ij}^{(k+1)} - u_{ij}^{(k)}| \} < \varepsilon \tag{5}$$

where ε is a termination criterion between 0 and 1 and k are the iteration steps. This procedure converges to a local minimum or a saddle point of J_m .

The algorithm is composed of the following steps.

1. Initialize $U = [u_{ij}]$ matrix, $U(0)$
2. At k -step: calculate the centers vectors $C(k)=[c_j]$ with $U(k)$
3. Update $U(k)$, $U(k+1)$
4. If $\|U(k + 1) - U(k)\| < \varepsilon$ then STOP; otherwise return to step 2.

3.3 Determination of Optimize Cluster

The criteria widely accepted for partitioning a data set into a number of clusters are the separation of the cluster and their compactness. The optimum case implies parameters that lead to partitions that are as close as possible (in term of similarity) to the real partitions of the data set [13]. A reliable quality assessment index should consider both the compactness and the separation. One of the quality measures that can be used in clustering is described as follows [8], [13].

The compactness of spatial data set x , called $\sigma(x)$, the value of the p -th dimension is defined as follows: [8], [13]

$$\sigma_x^p = \frac{1}{n} \sum_{k=1}^n (x_k^p - \bar{x}^p)^2 \tag{6}$$

where \bar{x}^p is the p -th dimension of

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k, \forall x_k \in x \tag{7}$$

The compactness computed from variance of cluster, Compactness of cluster i is called $\sigma(v_i)$ and its p -th dimension defined as:

$$\sigma_{v_i}^p = \frac{\sum_{k=1}^{n_i} (x_k^p - v_i^p)^2}{n_i} \tag{8}$$

The total compactness of spatial data set with respect to c cluster is:

$$\sigma = \sum_{i=1}^c \|\sigma(v_i)\| \quad (9)$$

Further the term $\|x\|$ is defined as: $\|x\| = (x^T x)^{1/2}$, where x is a vector

The average compactness of c cluster, $Comp$ is :

$$Comp = \sigma/c \quad (10)$$

The average scattering of data set compactness, $Scatt_Comp$ is:

$$Scatt_Comp = Comp/\|\sigma(x)\| \quad (11)$$

The more compact the cluster are the smaller the $Scatt_Comp$. Thus, for a given spatial data set, a smaller $Scatt_Comp$ indicates a better clustering scheme. The distance between cluster is defined by the average distance between the centroid of specified clusters, that is:

$$d = \frac{\sum_{i=1}^c \sum_{j=1}^c \|v_i - v_j\|}{c(c-1)} \quad (12)$$

The larger d is the more separated the clusters are according to definitions, a quality measure for clustering was defined as follows:

$$CD = Scat_Comp/d \quad (13)$$

The definition of CD indicated that both criteria of “good” clustering (i.e., compactness and separation) are properly combined, enabling reliable evaluation of clustering results [8]. Small values of CD indicate all the cluster in clustering scheme are overall compact and separated.

4 Experimental and Results

This section discusses data preprocessing, K-means clustering technique and result of experiment.

4.1 Data Preprocessing

At the initial stage of the experiment, data was scaled or normalized using

$$x_{new} = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \quad (14)$$

where x is the original data point, x_{\min} and x_{\max} are the minimum and maximum values in the data set, respectively. This is done in order to ensure that the minimum value in the data set is scaled to zero, and that the maximum value is scaled to one [5].

4.2 K-Means Clustering Compare with Fuzzy c-Means Clustering

The K-means clustering and Fuzzy c-means clustering methods were applied in this study intending to divide the groups of several sites of canals. Distances between site

and clusters are computed as Euclidean distances in the K-dimension space, a number of clusters N were designated arbitrary. As a result, very different means are obtained between clusters for most K parameters used in the analysis [10].

Because the distance in the K-dimensional space is greatly affected by differences in scale among the variables from which the distance are computed, each variable has been standardized before implementation of cluster analysis.

The cluster analysis of each initial matrix has been carried out in several steps, increasing the cluster number K from 2 to 6. From determination of optimized cluster small values of *CD* indicate all the cluster in clustering scheme are overall compact and separated that both criteria of “good” clustering.

Table 2. The result of optimize cluster

Number of K	Value of <i>CD</i>	
	K-means Clustering	Fuzzy c-means Clustering
K=2	1.599	2.256
K=3	1.445	2.433
K=4	1.193	2.463
K=5	0.682	1.963
K=6	0.965	1.971

Table 3. Descriptive mean of water quality parameters in five clusters

Cluster Parameters	Cluster 1	Cluster 2	Cluster 3	Cluster 4	Cluster 5
Temp	29.19	29.19	27.96	28.87	29.22
pH	6.91	6.97	7.14	6.98	7.10
DO	0.04	1.08	2.17	0.37	2.77
H ₂ S	1.76	0.11	0.08	0.76	0.01
BOD	40.92	12.40	9.33	23.81	6.05
COD	100.89	46.92	42.99	67.21	33.94
SS	26.68	33.05	39.08	23.61	28.56
TKN	13.98	7.15	4.64	9.59	4.34
NH ₃ N	9.42	2.15	1.40	3.71	0.78
NO ₂ N	0.04	0.12	0.11	0.05	0.11
NO ₃ N	2.28	1.69	0.66	1.48	1.67
T-P	1.75	0.99	0.67	1.37	0.47
T.Coliform	6.42E+09	2.18E+08	1.53E+06	2.64E+08	3.03E+07

Table 2 demonstrates the comparisons between *CD* values of K-means and Fuzzy c-means clustering based on number of clusters (K) from 2 to 6 clusters. It can be seen that for each number of cluster (K), K-means has a lower *CD* value than the

fuzzy c-means. Therefore, K-means provides a better performance than Fuzzy c-means through the dataset in this experiment. In addition, the CD value for each K using the K-means algorithm is the lowest when K is equal to 5 (in this case, the CD value is 0.682). Thus, the optimized cluster is equal to 5 from K-means clustering.

Each cluster consists of 13 parameters. Additional information for mean comparisons of water quality parameters in five clusters are presented in Table 3.

DO is the dissolved oxygen in the water while BOD is the Biological oxygen demand. The low DO will relate to high BOD. In table 3, cluster 1 shows lowest DO and highest BOD. It means that this cluster shows the worst quality of surface water. High COD (chemical oxygen demand) also indicates bad water quality. H_2S actually come from the decay of organic matter, which also causes a stinky smell. It should be implied that a high value of hydrogen sulfide also indicates a poor water quality. TKN (total nitrogen) comes from fertilizer, organic rubbish or from domestic water discharge. It means that high TKN represents a bad water quality while nitrate and nitrite can oxidize to nitrogen gas, making TKN a more significant value of water quality index by comparison. Total coliform shows microorganisms in the water. The high total coliform means high microorganisms or bacteria and also bad water quality. It can be summarized, from the table, that the best water quality is in cluster 5 and the worst is in cluster 1.

The mean of water quality indicated a relation of parameters in each cluster. Especially, regarding the values of Biochemical Oxygen Demand (BOD) and Dissolved Oxygen (DO), generally high value of BOD will indicate a quantity of low DO in surface water. Fig. 4 illustrates the comparison in means of these three selected parameters among 5 clusters. Dissolved Oxygen (DO) refers to the volume of oxygen that is contained in water and amount of oxygen that can be held by the water depends on the water temperature, salinity, and pressure. The amount of dissolved oxygen often determines the number and types of organisms living in that body of water. For example, fish like trout are sensitive to low DO levels (less than eight parts per million) and cannot survive in warm, slow-moving streams or rivers. Decay of organic material in water caused by either chemical processes or microbial action on

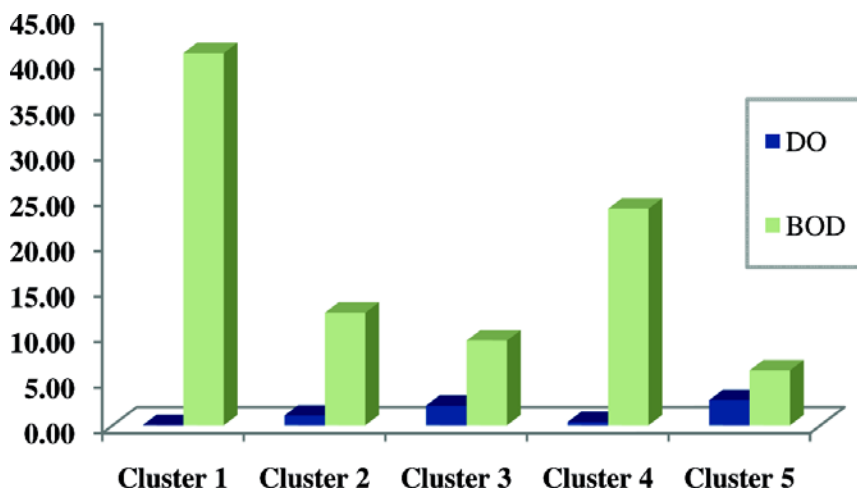


Fig. 4. A comparison of sample parameters in five clustering

Table 4. Water quality of records of site and number of sites in each cluster groups

Canal Number	Name of Canal	Total records of site	Number of sites in each Cluster:				
			1	2	3	4	5
1	Klong Kumuengderm	12	-	4	4	2	2
2	Klong Robkrung	24	1	4	7	8	4
3	Klong Mahanark	12	-	1	5	1	5
4	Klong PadungKrungkasem	20	-	2	9	4	5
5	Klong Samsen	16	2	3	8	2	1
6	Klong Sansab	40	-	9	21	4	6
7	Klong pramprachakorn	32	1	6	10	8	7
8	Klong Bangsue	12	1	1	3	5	2
9	Klong Satorn	8	1	2	2	1	2
10	Klong Somdetjaopraya	8	-	1	4	3	-
11	Klong Huai Khwang	16	-	4	7	4	1
12	Klong Lat phrao	16	-	4	11	-	1
13	Klong Mon	8	-	2	4	2	-
14	Klong Bangkoknoi	20	-	-	15	4	1
15	Klong Prapa	8	-	-	5	3	-
16	Klong Bangkokyai	20	-	2	16	2	-
17	Klong Bang khen	20	-	1	15	4	-
18	Klong Bangyikon	12	-	1	11	-	-
19	Klong Bang Po	8	-	-	8	-	-
20	Klong Tair	8	-	-	8	-	-
21	Klong Lat Yao	8	-	-	8	-	-
22	Klong Suan Lhong	8	-	-	8	-	-
23	Klong Bangkokabue	4	-	-	4	-	-
24	Klong Songtonnoon	4	-	-	4	-	-
	Total	344	6	47	197	57	37

untreated sewage or dead vegetation can severely reduce dissolved oxygen concentration. This is the most common cause of fish death, especially in summer months where the warmer water naturally holds less oxygen.

In this study, water quality data were collected during 2005-2008. There are 24 canals of 344 records of sites. Bangkok's canals are preferred for city transportation by many locals. The network of Bangkok Canals branches out of the Chao Phraya River, "the river of the kings", finally meeting with the ocean. The canals in Bangkok are hundreds of miles long. Not as many canals are used nowadays, and today both sides of the Chao Phraya River are totally different from what they used to be in the old days. Now, alongside the older homes and buildings, there are beautiful hotels. Table 4 shows cluster 1 consisting of 6 of 344 records of site, cluster 2 consists of 47 of 344 records of site, cluster 3 consists of 197 of 344 records of site, cluster 4 consists of 57 of 344 records of site and cluster 5 consists of 37 of 344 records of site, respectively.

Cluster 1 represents 6 records of site with a high value of BOD and low value of DO in surface water canals. The reasons for higher BOD and lower DO are site of canals connecting to industrial and populated accommodations. Most of the buildings in this area are the commercial and high rise buildings.

Cluster 2 represents 47 records of site and cluster 4 represents 57 records of site don't show the various site in the study area.

Cluster 3 represents 197 records of site and all sites of canal number 19 to 24 appeared in this cluster. The site of those canals are close to the sludge treatment centers that are knowledgeable about waste-water treatment including impact of pollution in surface water of canals. As a result, the canals of cluster 3 according to low value of BOD and high value of DO.

Cluster 5 represents 37 records of site and each site lowest value of BOD, COD and TKN but highest value of DO. Because, each site of canals in cluster 5 is near to pumping stations, non-point pollution from extractive and heavy industries, and some site collected to rainy season of high flow water, it can be summarized that the best water quality is in cluster 5.

6 Conclusions

In this paper, a comparison between clustering methods, K-means and Fuzzy c-means was performed. K-means algorithm is preferred since the CD values are lower than the Fuzzy c-means for all experimental clusters (2-6 clusters). Therefore, K-means clustering is applied to cluster surface water quality of canals in Bangkok, Thailand. The results indicate that five relatively cluster suitable for this sample site and conclude to parameters.

Some of the differences corresponded to different site of canals, or water surfaces separated by plugs or gates, that were held at different elevation and differences in source water. The BOD and DO value parameters can show water quality for each cluster. Better comprehension of the information may be helpful in further developing the water quality management. As a result, Cluster analysis is found to be a useful and efficient tool for stakeholders to manage natural resources, corresponding to the government's policy in water management.

Acknowledgments. The authors would like to thank Suan Sunandha Rajabhat University for scholarship support. Thanks to Department of Drainage and Sewerage Bangkok Metropolitan Administration for the provided data.

References

1. Bengraïne, K., Marhaba, T.F.: Using principle component analysis to monitor spatial and temporal changes in water quality. *Journal of Hazardous Materials* 100(1-3), 179–195 (2003)
2. Guler, C., Thyne, G.D.: Hydrologic and geologic factors controlling surface and groundwater chemistry in Indian wells-Owens Valley area, southeastern California. *Journal of Hydrology* 285(1-4), 177–198 (2004)
3. MacQueen, J.B., Foster, I., Kesselman, C.: Some Methods for classification and Analysis of Multivariate Observations: Procece C.: The Grid: Blueprint for a New Computing Infrastructure. In: Proc. Fifth Berkeley Symp. on Math. Statist. and Prob., vol. 1, pp. 281–297. Univ. of Calif. Press, Berkeley (1967)
4. Dogan, E., Sengorur, B., Koklu, R.: Modeling biological oxygen demand of the Melen River in Turkey using an artificial neural network technique. *Journal of Environmental Management* (2008)
5. Areerachakul, S., Sanguansintukul, S.: Water Classification Using Neural Network: A Case Study of Canals in Bangkok. In: The 4th International Conference for Internet Technology and Secured Transactions (ICITST 2009), Thailand. IEEE Press, United Kingdom (2009)
6. Anazawa, K., Ohmori, H.: The hydrochemistry of surface waters in Andesitic Volcanic area, Norikura Volcano, central Japan. *Chemosphere* 59(5), 605–615 (2005)
7. Astel, A., Tsakovski, S., Simeonov, V., Reisenhofer, E., Piselli, S., Barbieri, P.: Multivariate classification and modeling in surface water pollution estimation. *Analytical and Bioanalytical Chemistry* 390(5), 1283–1292 (2008)
8. Binbib, H., Fang, T., Guo, D.: Quality assessment and uncertainty handling in spatial data mining. In: Proc. 12th Conference on Geoinformations, Sweden (2004)
9. Helena, B., Pardo, R., Vega, M., Barrado, E., Fernández, J.M., Fernández, L.: Temporal evolution of groundwater composition in an alluvial aquifer (Pisuerga River, Spain) by principal component analysis. *Water Research* 34(3), 807–816 (2000)
10. Kanungo, T., Mount, D., Netanyahu, N., Piatko, D.C., Silverman, R.: An Efficient K-means Clustering Algorithm Analysis and Implementation. *IEEE Transaction on Pattern Analysis and Machine Intelligence* 24(7) (July 2002)
11. Laaksoharju, M., Gurban, I., Skarman, C., Skarman, E.: Multivariate mixing and mass balance (M3) calculations, a new tool for decoding hydrogeochemical information. *Applied Geochemistry* 14(7), 861–871 (1999)
12. Liu, C.W., Lin, K.H., Kuo, Y.M.: Application of factor analysis in the assessment of groundwater quality in a blackfoot disease area in Taiwan. *The Science of the Total Environment* 313(1-3), 77–89 (2003)
13. Halkidi, M., Vazirgiannis, M.: Cluster validity assessment: finding the optimal partitioning of a dataset. In: Proc. ICDM Conference, San Jose, California, USA, pp. 187–194 (1996)
14. Simeonov, V., Stratis, J.A., Samara, C., Zachariadis, G., Voutsas, D., Anthemidis, A., Sofoniou, M., Kouimtzi, T.: Assessment of the surface water quality in Northern Greece. *Water Research* 37(17), 4119–4124 (2003)
15. Simeonova, P., Simeonov, V., Andreev, G.: Water quality study of the Struma River Basin, Bulgaria. *Central European Journal of Chemistry* 2, 121–136 (2003)
16. Singh, K.P., Malik, A., Sinha, S.: Water quality assessment and apportionment of pollution sources of Gomti River (India) using multivariate statistical techniques: a case study. *Analytica Chimica Acta* 538(1-2), 355–374 (2005)

17. Vega, M., Pardo, R., Barrado, E., Deban, L.: Assessment of seasonal and polluting effects on the quality of river water by exploratory data analysis. *Water Research* 32(12), 3581–3592 (1998)
18. Bezdek, J.C.: Fuzzy partition and relations: an axiomatic basic for clustering. *Fuzzy Sets and Systems* 1, 111–127 (1978)
19. Liou, S., Lo, S., Hu, C.: Application of two-stage fuzzy set theory to river quality evaluation in Taiwan. *Water Research* 37, 1406–1416 (2003)
20. Ministry of Natural Resource and Environment, <http://www.mnre.go.th>
21. Department of Drainage and Sewerage Bangkok Metropolitan Administration, <http://dds.bangkok.go.th/wqm/thai/home.html>
22. Department of Pollution Control Bangkok, <http://www.pcd.go.th>
23. Ministry of Public Health, <http://eng.moph.go.th>
24. Luke, B.: K-Means Clustering, <http://fconyx.ncifcrf.gov/~lukeb/kmeans.html>

A Review of Routing Protocols for UWB MANETs

Yahia Hasan Jazyah and Martin Hope

Centre for Network and Telecommunications Research
The University of Salford
Salford, UK

y.h.jazyah @pgr.salford.ac.uk, m.d.hope@salford.ac.uk

Abstract. Mobile Ad hoc Networks (MANETs) have witnessed an increasing amount of interest during the last decade. Hosts (or nodes) within the network are mobile, and each node is equipped with a short range transmitter and receiver, antenna, and local power supply. Nodes then operate as a router to relay messages from the sender to the receiver, and they can be organised into different topologies; for example: they can be flat or hierarchical, they can move in any direction and speed, and they can communicate with each other through wireless routing protocols. More recently, systems based on Ultra-WideBand (UWB) technology have become a promising candidate in application to MANETs. This is mainly due to their powerful capabilities, such as their high data rates and low power consumption, and although many routing protocols have been designed for ad-hoc networks, few have considered their application to UWB based MANETs. This paper presents a review of different types of wireless routing protocols for MANETs with a concluding emphasis on their application to UWB based systems.

Keywords: MANET; UWB; routing protocol.

1 Introduction

The practical application of MANETs is sometimes preferable over infrastructure base networks due to the unavailability or high cost of wired infrastructure. They can be used effectively in military applications such as the battle field, where the existence of wired infrastructure is impossible. They are also very useful in industrial or commercial settings involving cooperative mobile data exchange, or when establishing communications in a disaster recovery situation such as a fire or rescue operation.

MANETs function in the absence of a base station where nodes communicate directly without the need of a centralised control system. Wireless routing protocols enable nodes within the MANET to relay data packets around the network. If nodes are within transmission range, they can communicate directly, or if they are away from each other, intermediate nodes are required to establish a multihop route between sender and receiver.

Now it is the time for UWB to be given more attention due to its great capability in high data rate, large BW, and low power consumption. Currently, hosts are equipped with UWB system in order to achieve short range communication and exploit the

previous mentioned advantages, besides UWB has the ability to measure distance between hosts accurately.

As a result, researchers have been encouraged to design new routing protocols for UWB networks to make it suitable for hosts to exploit features of UWB, but still needs more and more work to overcome drawbacks of network such as overhead, power consumption, unreliable route, etc...

UWB [1], [2] is a radio technology proposed to be used in personal area networks (PAN) and appeared in the IEEE 802.15.3a draft PAN standard. UWB systems consume very low energy levels, and can be used in short-range, high BW (BW>500 MHz, 20% of the center frequency). UWB operates in range 3.1–10.6 GHz (the minimum BW is 55MHz). The maximum allowable power spectral density is -41.3 dBm/MHz corresponding to average transmitted power of 0.5 mW when operating in band range 3.1-10.6 GHz.

UWB is best used for ad-hoc and sensor networks, it is used as part of location systems and real time location systems such as hospitals and healthcare, short broadcast time, “see-through-the-wall” precision radar imaging technology, and replacing cables between portable multimedia Consumer Electronics (CE) devices and in next-generation Bluetooth Technology devices [3].

It provides relatively long radio range which is around 150 meters indoor and 1 km outdoor and high data rate in excess of 100 Mbps with bit rates 55, 110 & 200 Mbps.

UWB uses pulse coded information with sharp carrier pulses (pulse-based systems), the signal consists of a train of pulses, one single pulse may be used to represent one bit.

There is duration between single pulses to avoid both pulse overlapping and ambiguity between pulse positions [3]. Each pulse occupies the entire UWB bandwidth, thus resist the multipath fading. UWB systems are carrierless systems based on pulses of very short duration transmission [4].

Generically, wireless routing protocols in MANETs can be classified into topology and position based protocols. These are presented in sections II and III, and of particular relevance to UWB systems, section IV then considers power aware routing protocols and then section V highlights the issues associated with routing protocols for UWB applications.

2 Topology Based Routing Protocols

Topology based routing protocols [5] use existing information about the network to flood (or forward) packets. There are two main routing strategies which are classified as topology based and these are: proactive and reactive routing protocols. In addition, hybrid protocols also exist which mix the previous two types.

A. Proactive Routing Protocols

Proactive (or table driven) protocols [5] maintain routing information for each node in the network and store information in routing tables. This information is then updated whenever the topology changes and so one or more routing tables are required by each node to store routing information. Most proactive strategies share the same features, but they differ in the number of routing tables and frequency of topological update. Examples of proactive routing protocols are Destination-Sequenced Distance Vector

(DSDV) [6], Cluster-head Gateway Switch Routing (CGSR) [7], Wireless Routing Protocol (WRP) [8], Optimized Link State Routing (OLSR) [9], and Hierarchical State Routing (HSR) [10]. The first three protocols are discussed below:

a) Destination-Sequenced Distance Vector (DSDV)

DSDV [6] is a flat proactive routing protocol which guarantees loop freedom, single shortest path to destination. Each node forwards its information among the network; neighbor nodes (intermediate nodes) receive this information and update their routing tables. Thus, all nodes have routes to all reachable nodes in the network in advance.

DSDV uses a sequence number to prevent loops; here a node will not forward information with the same sequence number twice. The higher the sequence number the newer the route; nodes update their routing tables whenever they receive route information that contains a higher sequence number than the stored one in routing table, or a route with same sequence number but shorter than the stored one.

When a source node wants to send data to destination, it checks its routing table and selects the shortest path to the destination. DSDV uses two types of packets; Full Dump packets which carry all available routing information and can require multiple network protocol data units (NPDUs), and a route Broadcast packet which contains the address of destination, number of hops, and the sequence number of the information received regarding the destination.

DSDV can not work in scalable network, especially when nodes move; this causes a high control packet update which increases network overhead.

b) Cluster-head Gateway Switch Routing (CGSR)

CGSR [7] is a hierarchical proactive routing protocol. Nodes are arranged in clusters, a cluster head (CH) is elected in each cluster by a distributed algorithm. Transmission is then achieved by CH inside clusters, which reduces overhead in network, but nodes within cluster are busy in selecting CH rather than packet relaying; which affects protocol performance. While transmission between clusters is achieved by gateway nodes that exist between the transmission ranges of two or more CHs as shown in Fig. 1.

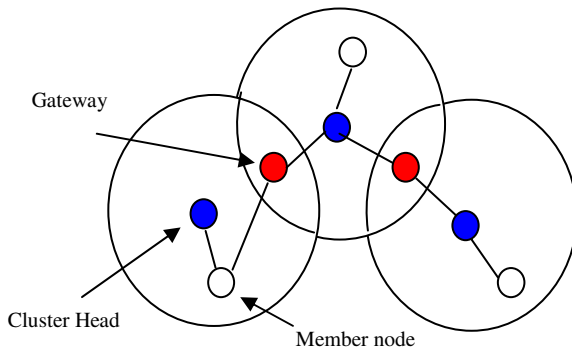


Fig. 1. Cluster Based in CGSR Routing Protocol

Each node maintains two tables; a cluster member table for storing destination CH of every node in the network, this table is broadcasted using the DSDV algorithm, and another routing table for storing neighbor nodes.

c) Wireless Routing Protocol (WRP)

WRP [8] is a flat proactive routing protocol which guarantees loop freedom by considering predecessor information. It uses four routing tables; distance table, routing table, link-cost table, and Message Retransmission List (MRL), table and thus it consumes memory.

Each record in MRL contains a sequence number, retransmission counter, acknowledgement required flag vector, and a list of update message that contains the destination, the distance to destination, and the predecessor of destination. An update message is sent in two cases; after processing updates from neighbours or when detecting a change in a link to neighbour.

WRP checks for nodes connectivity by sending hello messages. When a node receives hello message, it adds the node to its routing table, and then it sends a copy of its routing table to the node. This process causes overhead and consumes power.

d) Summary

Proactive routing protocols, in general, decrease end to end delay, but provide the shortest path. On the other hand, they increase network overhead, consume memory and power, and so, they are only suitable to small size networks.

B. Reactive Routing Protocols

Reactive routing protocols maintain route information on demand, i.e. when source node wants to send message to destination, it initiates route request (RREQ) to find route to destination. When a route fails, a route maintenance process is launched to repair the failed route.

Ad hoc On-Demand Distance Vector (AODV) [11], Dynamic Source Routing (DSR) [12], Temporally Ordered Routing Algorithm (TORA) [8], and Associativity-Based Routing (ABR) [8] are examples of reactive strategy.

More discussion about AODV and DSR and TORA is below:

a) Ad hoc On-Demand Distance Vector (AODV)

AODV [9] belongs to the class of Distance Vector (DV) routing protocols. In a DV, every node knows about its neighbours and the costs incurred in order to reach them using the Bellman-Ford algorithm [10].

AODV is a reactive shortest single path wireless routing protocol based on DSDV protocol. When a source wants to send a message to a destination, it checks its routing table, if there is a valid route to destination, it starts sending packets immediately. If not, it broadcasts RREQ to all neighbor nodes. It should be noted that RREQ contains the fields: hop count, source and destination sequence numbers, destination and source addresses, RREQ ID, and other pre-determined fields. When an intermediate node receives a RREQ, it checks its routing table for a path to the destination, if it exists, it unicasts a route reply (RREP) to the source, otherwise, it increases the hop count by one, and adds its ID to the RREQ and then re-broadcasts it to its neighbours; the process continues until the RREQ reaches its destination. The destination then selects the first coming RREQ and unicasts RREP using the reverse path to source node. When

the source receives several RREPs, it selects the route of the highest sequence number and minimum hop count, and then it establishes the route and starts sending packets.

To guarantee loop freedom, the source node uses a sequence number and includes it in RREQ. When a node receives a control message (RREQ, RREP, or RERR), it checks its routing table for an entry to the specified destination, if there is no entry in the routing table about the destination, it creates a new one. If there is an entry in the routing table, the route is only updated if the new sequence number is either higher than the destination sequence number in the routing table, the sequence numbers are equal, but the hop count plus one is smaller than the existing hop count in the routing table, or the sequence number is unknown. In addition, the source uses a time to live (TTL) count to limit the flooding of RREQ packets and control the overhead associated with the network. Finally, a HELLO message is broadcasted periodically to inform neighbor nodes about node existence. When an active node (node on the active route) detects a route failure (the neighbor node is unreachable; i.e. HELLO is not being received), it sends a route error (RERR) packet to the source node, which in turn, initiates a new RREQ.

Overhead is a major drawback of AODV because of the flooding of control messages (HELLO and RREQ packets). However, the overhead is low when compared to proactive routing protocols, and it also needs less memory, on the other hand it only stores the routes on demand so, it is preferable in application to large scale networks.

b) Dynamic Source Routing (DSR)

DSR [13] is a reactive routing protocol; when a source node wants to send data to a destination address; it initiates and broadcasts a RREQ that carries the full address (every hop). When intermediate receives the RREQ, it checks its routing table (nodes may store many routes to destination), if there is a valid path to destination, it sends RREP to source, or it forwards the packet after adding its address if its address is not existed in route record in order to avoid loop, and so on until RREQ reaches to destination which sends RREP to source using the reverse path.

A link failure is detected when no acknowledgement received by a node. RERR message is sent to source, which in turn, looks for a valid route in its cache, if no valid route is existed, source initiates RREQ again.

DSR does not use hello message, and so, it reduces the overhead and saves bandwidth, which makes it suitable for scalable network. On the other hand, packet size increases when route is long and network is large since every node adds its address to route record of RREQ.

c) Temporally Ordered Routing Algorithm (TORA)

TORA is a distributed and reactive routing protocol. It is suitable for highly dynamic topology network, it provides acyclic multipath but it suffers from long delay and does not support shortest path calculation. It is based on directed acyclic graph to maintain loop freedom.

When source want to forward packets to destination, it uses a height metric to establish direct acyclic graph and each node maintains information about one hop neighbors in order to localize control message to a very small set of nodes near the occurrence to topology change. When a link fails, a node generates a reference level which propagates by neighbor nodes.

d) Summary

Reactive routing protocols reduce network overhead, consume less bandwidth, and need less memory when compared to proactive protocols because they establish the route on demand. But they need a longer time than proactive protocols to establish the route between source and destination, and they also need to route provide maintenance because there are no ready available backup routes in general. In general, reactive routing protocols are suitable to large scale network compared to proactive ones.

Topology based routing protocols use flat addressing scheme as mentioned in the previous examples except CGSR which uses hierarchical addressing.

C. Hybrid Routing Protocols

Most networks are partitioned into zones (or zone-based) or nodes are grouped into trees or clusters. Hybrid routing protocols merge proactive and reactive strategies whereas proactive strategy is applied within zone (or intra zone), while a reactive strategy is applied outside zone (or inter zone). The following are examples of hybrid routing protocols:

a) Zone Routing Protocol (ZRP)

In ZRP [14], network is divided into routing zones, and each node belongs to a zone based on its geographical location. ZRP reduces the overhead of proactive protocols and the delay associated with reactive protocols, and so, it is suitable for large networks.

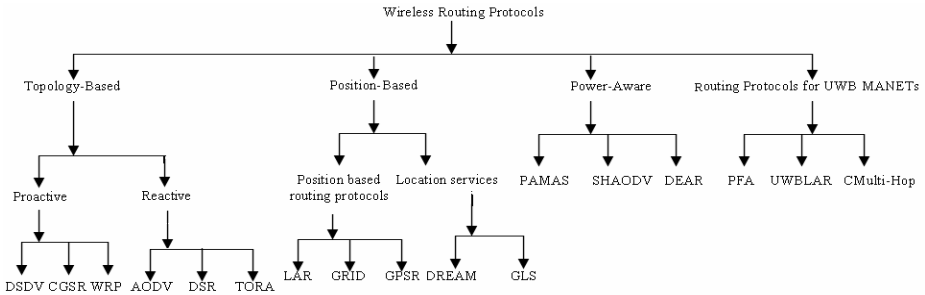


Fig. 2. Categorization of Ad hoc Routing Protocols

Zone radius is considered by number of hops. The selection of radius is a tradeoff between the routing efficiency of proactive routing and the increasing traffic in zone.

Source node looks for destination within zone proactively, if it is not existed, it sends RREQ reactively outside zone, if receiving node knows the path to destination, it sends RREP to source, otherwise, it rebroadcasts RREQ after adding its address until it reaches to destination.

In the case of route maintenance, packets can be directed around the broken link through an active multi-hop path.

b) Zone-based hierarchical link state (ZHLS)

ZHLS [7] is a hybrid hierarchical routing protocol acting similar to ZRP. It is proactive when Intra-zone Clustering, while it is reactive when Inter-zone Clustering. Each node has a node ID and zone ID based on Global Positioning System (GPS).

The network is divided in non overlapping zones without zone head which reduces overhead. Zone size depends on node mobility, network density, transmission power and propagation characteristics.

If destination is moving within the zone, no need to search in other locations since node ID and zone ID of destination are required for routing, and so, ZHLS may scale well to large networks; nevertheless, all nodes must have a pre-programmed static zone map, while geographical boundary of network is dynamic.

c) *Summary*

Hybrid routing protocols merge both proactive and reactive approaches together by utilizing advantages of both of them. They reduce the overhead of proactive protocols by applying reactive routing between zones and limiting flooding inside zones, and the delay associated with reactive protocols by applying proactive routing inside zones, thus it is suitable for large scale networks.

3 Position-Based Routing

Position based routing focuses on two issues: a location service to determine destination's position and forwarding strategy that can be greedy forwarding, restricted directional flooding, and hierarchical routing.

D. *Location Services*

Location service is used by sender node to learn the position information of a certain node in order to include it in the forwarded packet if it does not know the position of communication target. Each node registers its current position to a location service, when a node needs the position information of another one; it sends a request to location service. For example, Quorum-Based Location Service [5], Home-Zone [5], Distance Routing Effect Algorithm for Mobility (DREAM) [15], and Grid Location Service (GLS) [5], [16]. The last two location services are discussed below:

a) *Distance Routing Effect Algorithm for Mobility (DREAM)*

Each node maintains a position database that stores position information of each other nodes that are part of the network (all-for-all), where packets are regularly flooded to update position information [5].

DREAM [15] combines proactive and reactive approaches. Source node has to have positions of it, destination, one hop neighbours, and the maximum speed of destination. Source node uses this information to determine the direction towards destination in order to send packets by finding one hop neighbor within the range between source and destination according to a certain angle determined by speed and position of destination.

b) *Grid Location Service (GLS)*

GLS [5], [16] divides network into a hierarchy of squares. In each hierarchy, n -order squares contain exactly four $(n - 1)$ -order squares, forming a "quadtree". Each mobile node periodically updates a small set of other nodes with their current locations by maintaining a table of all other nodes within the local first-order square. (Classified as all-for-some)

GLS uses geographical forwarding. Source must know the geographical positions of any destination, and must label packets for that destination with its position, while intermediate node only needs to know its position and the positions of nearby nodes.

Geographical forwarding is similar to Cartesian routing. Each node broadcast a HELLO message that contains its position and velocity; where neighbours within transmission range store this information in a table. Node sends packet to the closest neighbor of destination after checking the neighbor table. When a node does not find closest node other than itself, it sends error indicating a hole in geographic distribution of nodes. Greedy Perimeter Stateless Routing (GPSR) is a solution for the problem of hole.

GLS is based on the idea that a node maintains its current location in a number of location servers (each node acts as a location server on behalf of some other nodes) distributed through over network, but they are no elected leaders. GLS makes balance when location servers are selected and dense of network is considered.

c) *Summary*

Location service is used by source node to determine position of destination. Location services are classified according to “how many servers maintain position information about mobile nodes “into all for all, all for some, some for all, and some for some. For example, when all nodes within the network have information about all nodes in the network, then the location service is classified as “all for all”.

E. *Position-Based Routing Protocols*

Position based routing protocols use position information gathered by GPS system. Each node in the wireless network knows its location, destination location, and other geographical information in order to forward packets in geographical direction instead of flooding as in topology based protocols.

Mainly, there are three forwarding strategies: First, greedy packet forwarding; whereas the sender of packet includes the approximate position of recipient in the packet, intermediate node forwards the packet to a neighbor which is lying in the general direction of the recipient. Second, restricted directional flooding; the packet is flooded in one-hop direction towards the destination where it is expected to be located; this process is repeated until the packet reaches the destination. The direction is determined by calculating the expected region that covers destination. Third, hierarchical Routing; where the complexity that each node has to handle can be reduced massively by establishing some form of hierarchy.

Many position based routing protocols have been designed such as, Location Aided Routing (LAR) [17], GRID [18], Greedy Perimeter Stateless Routing (GPSR) [19], Adaptive Location Aided Roaming from Mines Protocol (ALARM) [20], Location-Aware routing Protocol with Dynamic Adaptive of Request Zone (LARDAR) [21], Compass [22], Angle Routing PROTOCOL (ARP) [23], Position-based Selective Flooding (PSF) [24], Octopus [25], Improved Location-Aided Routing (ILAR) [26] , and Novel Routing Power Balance (NRMPB) [27] . The following is a discussion about position based routing protocols.

a) *Location Aided Routing Protocol (LAR)*

LAR [17] is an on-demand routing protocol that uses the modified Dijkstra's Algorithm to find the shortest path.

Destination lies in the centre of a circular region of certain radius at a certain time known as Expected Zone, as shown in Fig. 3, which indicates which zone of the network should be reached by RREQ. GPS enables terminals to know its own position and speed, while dissemination is performed by piggybacking location information in all routing packets.

“Piggybacking is a bi-directional data transmission technique in the network layer. It makes the most of the sent data frames from receiver to emitter, and adds the confirmation that the data frame sent by the sender was received successfully acknowledge.” [28].

There are two proposed schemes of LAR:

First, LAR1 defines request zone that includes sender and receiver on opposite corner of a rectangle. Rectangle dimensions are estimated according to the receiver average speed at a certain time. Only nodes within the zone respond to the RREQ of sender. LAR1 reduces network overhead, but causes delay when route is found.

Different shapes can be used other than rectangle such as circle, cone, and bar.

LAR2 is proposed based on the calculated distance between source position and the estimated position of destination. When node receives RREQ, it calculates the distance toward destination, if the distance is less than the distance from the previous sender node to destination; it forwards the packet as shown in fig. 4. In this scheme, intermediate receiving node may be the closest node to destination, and so the algorithm reaches to a dead-end.

b) *GRID*

GRID [18] is a reactive routing protocol which exploits location information. GRID has strong maintenance capability; where intermediate node that attends to move performs handoff operation to keep route alive.

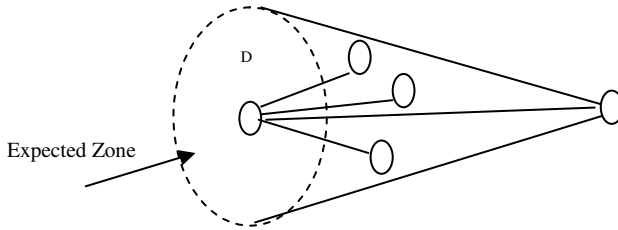


Fig. 3. Expected Zone in LAR

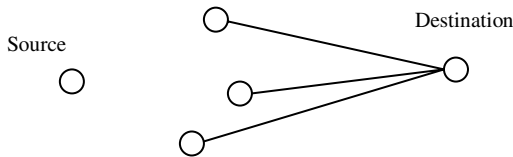


Fig. 4. LAR Scheme 2

Geographical area is divided into square grids; the nearest node to the centre of grid is elected as a leader (gateway) to maintain routing. If leader leaves a grid, it broadcasts routing table to the new leader. Position information of source and destination are used to confine.

The search area, leader is responsible for searching procedure and setting up routing table in grid by a grid manner; routing in GRID is performed grid by grid through grids' leaders.

Source node calculates the searching range depending on its location and the location, speed, and direction of destination in order to eliminate the effect of collection of redundant rebroadcast, heavy contention, and collisions (problem of broadcast storm).

c) *Greedy Perimeter Stateless Routing (GPSR)*

In GPSR [19], the closest intermediate node to destination is selected as next hop, but GPSR does not guarantee the existence of path, and so perimeter forwarding selects the neighbor which is farther than itself from the destination in order to solve the dead-end caused by greedy forwarding, GPSR switches between these two strategies. A modification to GPSR is made by using flooding to solve the problem of dead-end.

Each node needs information about next hop only, which is exchanged by means of periodic beacons broadcasted by each node.

d) *Octopus*

Octopus [25] is a fault-tolerant and efficient position-based routing protocol. Fault-tolerance is achieved by employing redundancy (storing the location of each node at many other nodes) and by keeping frequently refreshed soft state; nevertheless, it achieves a low location update overhead by employing a novel aggregation technique.

Octopus is a typical routing scheme, assumes that each node gets its own location information using like GPS. It divides the network into grid by horizontal and vertical strips, each node maintains a table for neighbors (updated by hello message) and a table for strips. Octopus forward RREQ in north and south directions to find destination, if this process fails, it sends the RREQ to the horizontal direction. Octopus uses Greedy Forwarding to forward data packets.

e) *Improved Location-Aided Routing (ILAR)*

ILAR [26] improves location-aided routing (LAR); when intermediate node receives RREQ, it checks whether it is within the requested zone, if it is, it will send RREQ revise (RREQ_R) to the transmitting node, RREQ_R includes the VDIST (vertical distance, it is the distance from node to the line connecting source and destination; base line) as shown in fig. 5. The node which receives RREQ_R will select the next hop which has the smallest VDIST and has a distance to source larger than to the previous node.

This technique has less packets and better performance than LAR.

f) *Novel Routing Power Balance (NRMPB)*

NRMPB [27] resolves the problems of ILAR which are: the movement of nodes near border of coverage broadcast which cause breakage of route, the other problem is selecting the next hop node; NPMPB comply position information and power saving to balance the energy of each nodes and prolong lifetime of network.

In ILAR protocol, the intermediate node which has the smallest vertical distance is selected as next hop, but the selected node could be located near the boundary of

transmission range or predecessor node, which cause unreliable (unstable) link. NRMPB solves this problem besides considering power saving.

When source wants to send packet to destination, it initiates RREQ packet including a field of wait time, each receiving node calculates the wait time and replaces it in the RREQ packet before forwarding. The wait time is reduced and when it counts to zero, node detects an error in path.

When source or intermediate node wants to find out the next hop, it sends route request transmit (RREQ_T) packet, which contains source and destination's locations, all nodes receive the packet will calculate their position with respect to source's position, if they are located within 80% of transmission range of source, they will reply by route request revise (RREQ_R) including the power of node which is calculated by a certain equation, the equation depends on distance between source and destination, distance between node and source, distance between node and the line connecting source and destination (baseline), see fig. 5, and node's power.

Source selects the node of maximum power within the 80% of its coverage range, and unicasts RREQ to it, the process continues until RREQ reaches to destination which replies by RREP using the reversed path to source node.

When a node which sent RREQ can not reach the next hop, it will try to find another next hop, if it fails, it will send RERR to the previous node, and so on. Another technique is used when node does not receive RREP from destination, it will forward RREQ_T.

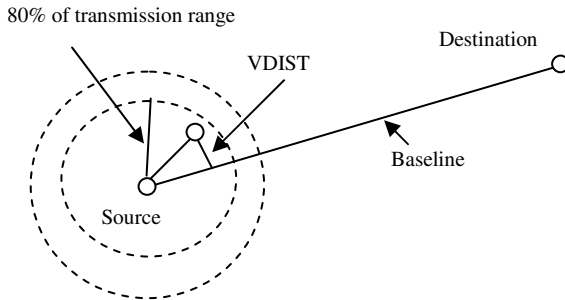


Fig. 5. VDIST in NRMPB and ILAR Routing Protocols

g) *Summary*

Position based routing protocols eliminates overhead of reactive protocols by directing flood of route request towards the expected location of destination, which leads to better performance in all aspects; power consumption, stability, reliability, maintenance, and bandwidth utilization. These specifications make position-based routing widely trusted to large scale networks.

4 Power-Aware Routing

Power metrics are important issues when routing protocols are designed since mobile nodes have limited resource of power which is represented by local batteries.

Generally, energy efficient issues are: 1) Sleeping mode; whereas nodes turn themselves off when not actively transmitting or receiving (sleep/awake). 2) Transmission Power Control; which aims to minimize the power per bit required to transmit packet from source to destination, this technique always selects the least-power cost routes; as a result, nodes along those routes tend to die soon because of battery energy exhaustion. 3) Power-aware route selection that find route at certain route discovery time such that a cost function is minimized. Cost function depends on energy consumed per packet, time to network partition, variance in node power levels, cost per packet, and node cost metrics. And 4) Broadcast control.

Examples of power-aware routing protocols are: Power-Aware Multiple Access Protocol (PAMAS) [29], Self-Healing-AODV (SHAODV) [30], Energy-Efficient Ad hoc Routing (DEAR) [31], Energy-Aware Ad hoc On-Demand Distance Vector (EAAODV) [32], Active Route Maintenance Protocol (ARMP) [33], and Distributed Power Control (DPC) [34]. The following are more details about routing protocols related to power.

a) Power-Aware Multiple Access Protocol (PAMAS)

PAMAS [29] is a shortest-cost routing algorithm, power-Aware Multiple Access protocol with Signalling, it uses a new power-aware metrics for determining routes in wireless ad hoc networks based on battery power consumption of node.

Nodes' radios are turned off (sleep mode) when they can not transmit or receive packets during transmission period of other nodes; which in turn saves their power.

b) Self-Healing-AODV (SHAODV)

SHAODV [30] overcomes the problem of edge effect and route break without causing heavy overhead. The path can be established for long lifetime by connection admission control of route request. The unstable route can be healed without route break, flood over the network, and overhead by link state inspection and active route local repair.

Signal to Noise Ratio (SNR) is measured and compared to a threshold in order to overcome the overhead; it is calculated according to the received power which is based on transmitted power, height and gain of antenna at both sides of transmitter and receiver, and distance between them.

Movable node requests the source to heal the unstable path before it moves and causes break in the route. Process of reduction the repaired path is initiated after that by broadcasting RREQ for reduction along the repaired path and accepting the first received RREQ packet with link state value (SNR) greater than a threshold.

c) Energy-Efficient Ad hoc Routing (DEAR)

DEAR [31] controls the rebroadcast time of RREQ packets to prolong network lifetime by balancing between minimum transmission energy consumption and fair node energy consumption in a distributed manner without additional control packets. It also improves data packet delivery ratio.

The physical layer can save energy by adapting transmission power according to the distance between nodes. While at the data link layer, energy conservation can be achieved by sleep mode operation.

Network lifetime is defined as the time when a node runs out of its own battery power for the first time.

DEAR utilizes only the RREQ message to distribute the average residual battery level of the entire network. Intermediate nodes control the rebroadcast time of the RREQ packet whereas retransmission time is proportional to the ratio of average residual battery power of the entire network to its own residual battery power, and so, nodes with larger battery energy will rebroadcast RREQ packets earlier.

Residual battery level of the entire network depends on the average residual power of route, which in turn depends on the residual power of node and hop count. DEAR control the rebroadcast time; if the residual battery power of node is smaller than the average network residual power, the retransmission time will be longer and vice versa. DEAR compromises between the min-hop path and the fair energy consumption path.

d) Energy-Aware Ad hoc On-Demand Distance Vector (EAAODV)

EAAODV [32] combines the on-demand routing capability of the AODV routing protocol with a distributed topology discovery mechanism – mobile agents, thus reaching the energy balance condition. EAAODV eliminates the problem of low efficiency and link breaking in routing due to the energy exhaustion of some nodes.

"A mobile agent is a software program that can suspend its execution on the host computer, transfer itself to another agent-enabled host through the network, and resume its execution on the new host." The key features of mobile agents which distinguish them from traditional distributed programming are: mobility; network awareness; communication; intelligence; reactivity; autonomous; goal-oriented; temporally continuous; learning; and flexible.

When intermediate node receives a RREQ, it decided to broadcast it according to its residual energy, if it is under a threshold level, the node terminates the RREQ, otherwise it broadcasts the RREQ. During the search process, the node checks its routing table, if the table is fresh; the residual energy information of the next-hop node indicated by the table is acquired by mobile agents, if the energy is above the threshold, the routing table is valid and it will be used, otherwise, searching for the routing will be finished, if the routing table is old, nodes will broadcast RREQ to their neighbors and update the information of the table.

EAAODV is able to achieve longer surviving network duration and a larger throughput compared to AODV.

e) Summary

Methods that concern power control issue can be classified into three categories: methods find the optimal transmitted power to control connectivity properties of network, power aware routing, and algorithms based on shortest path rather than hop count. In general, route maintenance is detected according to a threshold value which is compared to a certain power metric.

5 Routing Protocols for UWB Networks

UWB is a radio technology that can be used at a very low energy levels, short-range, high-bandwidth (>500 MHz, 20% of centre frequency) communications by using a large portion of the radio spectrum. UWB has been proposed as technology to use in Personal Area Networks (PANs). Many routing protocols have been designed to exploit features of UWB system such as Piconet Formation Algorithm (PFA) [35],

UWBLAR [36], Clustered-Multi-Hop (CMulti-Hop) [37], and as mentioned in papers [38] [39] [40] [41] [42]. The following are more details about some of routing protocols dedicated to UWB MANETs.

a) Piconet Formation Algorithm (PFA)

PFA [35] is a routing protocol for UWB networks that adopts a master-slave configuration and tries to minimize the total transmission power of master nodes and interference within piconets. Master nodes are selected according to the minimum average distance to neighbouring nodes and forms piconets that have the minimum total emission power for a given radio range. Master-slave configuration provides good performance in terms of control overhead, power savings, and overall efficiency. On the other hand, centralized schemes cannot be scaled to a large number of nodes or over a large distance. Hybrid approach can overcome this weak point.

PFA can reduce average interference generated per piconet, but it suffers from a load balancing problem; if there are too many nodes joining one cluster, the cluster head (master) may have to limit data rates to an undesirable level, thus performance may degrade. Moreover, a cluster with high total transmit power will generate high interference in neighbouring clusters. The solution of this problem is to impose an upper limit of nodes joining the cluster as proposed in [42].

b) UWBLAR

UWBLAR [36] is a position based routing strategy that sets up a dynamic routing list for an ad hoc network based on a UWB system.

The network is organised in a cluster manner where each node knows the exact position of every node in its cluster, when a connection request is generated; the cluster head sets up a table of distance information by broadcasting RREQ to all network nodes, the shortest path is then selected by the destination when it receives several RREQ.

c) Clustered-Multi-Hop (CMulti-Hop)

CMulti-Hop [37] is a new strategy for routing in UWB based ad-hoc networks that depend on minimizing power-dependent global cost function.

The model considers fixed nodes equipped by UWB systems, each node belongs to a cluster and knows the position of other nodes within the cluster, where cluster's size depends on radio channel condition and network load.

The first node of cluster is labelled by zero, the second level by one, on so on. A node set up a link with only a node of lower level, while no link between nodes within the same cluster is established in order to guarantee loop freedom.

Path of the lowest cost function is selected. If two paths have the same value, the path of minimum hop count is selected. If none of paths has a cost value greater than predefined cost value, the connection request is rejected.

d) Other Routing Strategies for UWB MANET

Paper [38] presents a position based routing strategy which exploits the high precision ranging capabilities offered by UWB In order to reduce the emitted power and multi-user interference.

Author considers the cost of route is the sum of the costs of all links that form the route, which comprises signaling cost for setting-up a new link, cost for transmitting

data which depends upon the requested data rate, both represent power consumption, and distance between two nodes. Besides Signal to Noise Ratio (SNR), Symbol Error Probability (SEP), and Packet Error Probability (PEP) are taken into account when multi-user interference (MUI) is considered.

The above technique is applied to DSR routing protocol. Whereas each intermediate node adds the cost of previous hop to the path cost. Destination selects the path according to the cost information included in RREQ. Also it is applied to LAR routing protocol.

Adoption of position-aware routing strategies significantly improves the efficiency in power management for UWB networks, for example, by optimizing the path selection procedure in DSR and trade-off between performance and power efficiency can be achieved by adopting different shapes for the Request Zone of LAR.

Paper [39] proposes power-efficient routing strategy based on the ranging capabilities offered by UWB.

The cost function of link is considered as in [38]. Minimum Cost strategy leads to more robust system in a high-interference environment; the Minimum Cost strategy is not effective in the case of low network connectivity since it increases the average number of hops with no advantage in system performance. When network connectivity grows, the Minimum Cost effectively increases the percentage of found connections by reducing the effect of multi-user interference.

Paper [41] proposes a novel packet routing scheme for UWB ad-hoc network. This approach utilizes the distance measurement between nodes that UWB provides, adopting the multihop transmission, determined spreading factor and the signal to interference and noise ratio (SINR).

Spreading Factor (SF) is the ratio of the chips to baseband information rate. Spreading factors vary from 4 to 512 in FDD UMTS. Spreading factor in dBs indicates the process gain, whereas the lower the spreading factor the higher the data rate.

Source node searches for node within its transmission range provided that the distance to destination is greater than the distance between the previous node and destination, the process is repeated by each node until destination is reached.

e) Summary

UWB is a new technology that is in MANET to transmit data packets in high speed. Routing protocols for UWB MANET have been designed to use the high performance and capabilities of UWB; such routing protocols are designed as master slave in general to overcome network overhead and reduce power consumption, other routing techniques consider a cost function that is accumulated until the RREQ reaches its destination, whereas the destination selects the route of cost function that satisfy the criteria.

Many routing protocols have been designed for UWB MANET, nevertheless the field needs more and more effort by researcher to design new protocols that can use the capabilities of UWB.

6 Conclusion

This work summarises wireless routing protocols for MANETs which can be classified into topology based and position based protocols. Topology based protocols are either

proactive or reactive, while hybrid strategies merge proactive and reactive approaches in order to enhance the drawback of proactive protocols which utilize the advantages of reactive ones.

Position based strategies are based on positional information gathered by GPS systems. They eliminate the overhead caused by topology based strategies by directing flooding towards the destination.

Power metrics are included in both topology and position based routing in order to enhance the performance of routing protocols in terms of power saving. They also take part in route maintenance as detector of route failure before it occurs, or as a metric to select the next hop.

UWB is a radio technology proposed to use in PANs and has appeared in the IEEE 802.15.3a draft PAN standard. UWB systems consume very low energy levels, and they are applicable to short-range high data rate systems which make it particularly useful to MANETs and sensor networks. Many conventional routing protocols have been designed for MANETs. However routing strategies for UWB networks need more work but will most likely employ aspect of power aware routing and geographic position.

References

- [1] Xia, L., Lo, A., Niemegeers, I., Baugé, T.: An Ultra-Wide Band Based Ad Hoc Networking Scheme for Personnel Tracking in Emergencies. In: This work was performed in the IST Framework 6 EUROPCOM project
- [2] Kuladinithi, K., Timm-Giel, A., Görg, C.: Mobile Ad Hoc Communications in AEC Industry (August 2004), <http://www.itcon.org/2004/22/>
- [3] Benedetto, M., Giancola, G., Domenicali, D.: Ultra Wide Band radio in Distributed Wireless Networks. IEEE, Los Alamitos (2006)
- [4] Project SPEARHUW, Enhanced Sensor and WLAN Networks using UWB Communications (April 2005)
- [5] Mauve, M., Widmer, J., Hartenstein, H.: A Survey on Position-Based Routing in Mobile Ad Hoc Networks. *IEEE Networks* 15(6), 30–39 (2001)
- [6] Perkins, C.E., Bhagwat, P.: Highly Dynamic Destination-Sequenced Distance-Vector Routing (DSDV) for Mobile Computers. *ACM SIGCOMM Computer Communication Review* 24(4), 234–244 (1994)
- [7] Abolhasan, M., Wysocki, T., Dutkiewicz, E.: A review of routing protocols for mobile ad hoc networks. *Ad Hoc Networks* 2, 1–22 (2003)
- [8] Royer, E.M., Santa, B., Toh, C.: A Review of Current Routing Protocols for Ad hoc Mobile Wireless Networks. *IEEE Wireless Communications* 6(2), 46–55 (1999)
- [9] Clausen, T., Jacquet, P., Laouiti, A., Muhlethaler, P., Qayyum, A., Viennot, L.: Optimized Link State Routing Protocol. In: *Proceedings of IEEE INMIC* (2001)
- [10] Kilinkaridis, T.: Routing Protocols for Wireless Ad Hoc Networks Hierarchical routing protocols
- [11] Perkins, C., Royer, E.B., Das, S.: Ad hoc On-Demand Distance Vector (AODV) Routing (2003), <http://www.ietf.org/rfc/rfc3561.txt>
- [12] Baumann, R.: AODV, Ad hoc On Demand Distance Vector Routing Protocol. AODV Presentation at ETH, Zürich (2002)

- [13] Johnson, D.B., Maltz, D.A., Broch, J.: DSR: The Dynamic Source Routing Protocol for Multi-Hop Wireless Ad Hoc Networks. In: *Ad Hoc Networking*, ch. 5, pp. 139–172 (2001)
- [14] Beijar, N.: Zone Routing Protocol (ZRP), <http://www.netlab.tkk.fi/opetus/s38030/k02/Papers/08-Nicklas.pdf>
- [15] Nardis, L.D., Giancola, G., Benedetto, M.: Power-aware design of MAC and routing for UWB networks. In: *Globecom Workshops*. IEEE Communications Society, Los Alamitos (2004)
- [16] Jinyang, L., Jannotti, J., Douglas, S.J., Karger, D.R., Morris, R.: A scalable location service for geographic ad hoc routing. In: *Proc. MOBICOM*, pp. 120–130 (2000)
- [17] Ko, Y.B., Vaidya, N.H.: Location-Aided Routing (LAR) in mobile ad hoc networks. *Wireless Networks* 6, 307–321 (2000)
- [18] Liao, W.H., Tseng, Y.C., Sheu, J.P.: GRID: A fully location-aware routing protocols for mobile ad hoc networks. *Telecommunication Systems* 18(1-3), 37–60 (2001)
- [19] Karp, B., Kung, H.T.: GPSR: Greedy Perimeter Stateless Routing for Wireless Networks. In: *MobiCom* (2000)
- [20] Boleng, J., Camp, T.: Adaptive Location Aided Mobile Ad Hoc Network Routing. In: *IEEE International Conference on Performance, Computing, and Communications*, pp. 423–432 (2004)
- [21] Shih, T.F., Yen, H.C.: Location-aware routing protocol with dynamic adaptation of request zone for mobile ad hoc networks, October 9. Springer, Heidelberg (2006)
- [22] Kranakis, E., Singh, H., Urrutia, J.: Compass Routing on Geometric Networks. In: *Proc. 11th Canadian Conf. Computational Geometry* (August 1999)
- [23] Banka, R.K., Xue, G.: Angle Routing Protocol: Location Aided routing for Mobile Ad-hoc Networks Using Dynamic Angle Selection. IEEE, Los Alamitos (2002)
- [24] Abolhasan, M., Wysocki, T.: GPS-based Route Discovery Algorithms for On-demand Routing Protocols. In: *MANETs* (January 2004)
- [25] Nakagawa, H., Ohta, T., Ishida, K., Kakuda, Y.: A Hybrid Routing with Location Information for Mobile Ad Hoc Networks. In: *Eighth International Symposium on Autonomous Decentralized Systems (ISADS 2007)*. IEEE, Los Alamitos (2007)
- [26] Wang, E.N.C., Wang, S.M.: An Efficient Location-Aided Routing Protocol for Mobile Ad Hoc Networks. In: *Proceedings of the 2005 11th International Conference on Parallel and Distributed Systems (ICPADS 2005)*. IEEE, Los Alamitos (2005)
- [27] Li, J.P., Wu, L.Y., Lin, S.Y., Wu, G.M.: A Novel Routing Protocol Integrate Power Balance for Ad Hoc Networks. In: *Third International Conference on Networking and Services (ICNS 2007)*. IEEE, Los Alamitos (2007)
- [28] Ucan, E., Thompson, N., Gupta, I.: A PiggyBacking Approach to Reduce Overhead in Sensor Network Gossiping. In: *MidSens 2007*, Newport Beach, USA, November 26-30, ACM, New York (2007); Copyright 2007, 978-1-59593-929-6/07/11
- [29] Singh, S., Mghavendra, M.W.C.S.: Power-Aware Routing in Mobile Ad Hoc Networks. In: *Proceedings of Fourth Annual ACM IEEE International Conference on Mobile Computing and Networking*, Dallas, TX, October 1998, pp. 181–190 (1998)
- [30] Feng, M., Cheng, S., Zhang, X., Ding, W.: A Self-Healing Routing Scheme Based on AODV in Ad hoc Networks. In: *The Fourth International Conference on Computer and Information Technology CIT 2004*, September 14-16, pp. 616–620 (2004)
- [31] Gil, H., Yoo, J., Lee, J.: An On-demand Energy-efficient Routing Algorithm for Wireless Ad hoc Networks. In: Chung, C.-W., et al. (eds.) *HSI 2003*. LNCS, vol. 2713, pp. 302–311. Springer, Heidelberg (2003)

- [32] Chenchen, Z., Zhen, Y.: A New EAAODV Routing Protocol based on Mobile Agent. In: International Conference on Systems and Networks Communications, ICSNC 2006 (October 2006)
- [33] Chang, C.Y., Chang, C.T., Tu, S.C., Hsieh, T.T.: Active Route-Maintenance Protocol for Signal-Based Communication Path in Ad-Hoc Networks. In: Proceedings of Ninth IEEE International Conference on Networks, October 10-12, pp. 25–30 (2001)
- [34] Bergamo, P., Giovanardi, A., Travasoni, A., Maniezzo, D., Mazzini, G., Zorzi, M.: Distributed Power Control for Energy Efficient Routing in Ad Hoc Networks. *ACM/Kluwer WINET* 10(1) (January 2004)
- [35] Gong, M.X., Midkiff, S.F., Buehrer, R.M.: A New Piconet Formation Algorithm for UWB Ad Hoc Networks. In: IEEE Conference on Ultra Wideband Systems and Technologies, November 16-19, pp. 180–184 (2003)
- [36] Huang, Y., Zeng, W.: A reliable Routing Protocol for Multi-hop Ad hoc Networks. In: International Conference on Wireless Communications, Networking and Mobile Computing, WiCom, September 21-25, pp. 1656–1659 (2007)
- [37] Nardis, L., Baldi, P., Benedetto, M.G.: UWB Ad-Hoc Networks. In: IEEE Conference on Ultra Wideband Systems and Technologies, Baltimore, Maryland, USA, May 20-23, pp. 219–223 (2002)
- [38] Nardis, L.D., Giancola, G., Benedetto, M.G.D.: A position based routing strategy for UWB networks. *IEEE*, Los Alamitos (2003)
- [39] Nardis, L.D., Giancola, G., Benedetto, M.G.D.: A power-efficient routing metric for UWB wireless mobile networks. In: *Procs. of Vehicular Technology Conference, IEEE VTS Fall VTC 2003*, October 2003, pp. 3105–3109 (2003)
- [40] Horie, W., Sanada, Y.: Novel Routing Schemes Based on Location Information for UWB Ad-Hoc Networks. *Wiley Periodicals, Inc.*, Chichester (2004)
- [41] Home, W., Sanada, Y.: Novel Packet Routing Scheme based on Location Information for UWB Ad-hoc Network. *IEEE*, Los Alamitos (2003)
- [42] Gong, M.X., Midkiff, S.F., Buehrer, R.M.: A Self-Organized Clustering Algorithm for UWB Ad Hoc Networks. In: *WCNC 2004*. *IEEE Communications Society*, Los Alamitos (2004)

An Efficient and Reliable Routing Protocol for Wireless Mesh Networks

Jaydip Sen

Innovation Lab, Tata Consultancy Services Ltd., Bengal Intelligent Park,
Salt Lake Electronics Complex, Kolkata, India
Jaydip.Sen@tcs.com

Abstract. Wireless mesh networks (WMNs) have emerged as a key technology for next generation wireless networks showing rapid progress and inspiring numerous applications. The persistence driving force in the development of WMNs comes from their envisioned advantages including extended coverage, robustness, self-configuration, easy maintenance, and low cost. However, to support real-time applications with stringent quality of support (QoS), WMNs must be equipped with a robust, reliable and extremely efficient routing protocol so that packets can be routed through them with minimum delay. In this paper, we focus on the critical factors in designing a routing protocol for WMNs, and propose an efficient and reliable routing protocol. The protocol is based on a reliable estimation of available bandwidth in a wireless path and end-to-end delay measurements. Simulations carried out on the protocol demonstrate that it is more efficient than some of the current routing protocols.

Keywords: Wireless mesh networks, quality of service, routing, end-to-end delay, bandwidth estimation, selfish nodes.

1 Introduction

Wireless mesh networking has emerged as a promising concept to meet the challenges in next-generation wireless networks such as providing flexible, adaptive, and reconfigurable architecture while offering cost-effective solutions to service providers. WMNs are multi-hop wireless networks formed by mesh routers (which form a wireless mesh backbone) and mesh clients. The mesh routers provide a rich radio mesh connectivity which significantly reduces the up-front deployment cost of the network. Mesh routers are typically stationary and do not have power constraints. However, the clients are mobile and energy-constrained. Some mesh routers are designated as gateway routers which are connected to the Internet through a wired backbone. A gateway router provides access to conventional clients and interconnects ad hoc, sensor, cellular, and other networks to the Internet. A mesh network can provide multi-hop communication paths between wireless clients, thereby serving as a community network, or can provide multi-hop paths between the client and the gateway router, thereby providing broadband Internet access to the clients.

As WMNs become an increasingly popular replacement technology for last-mile connectivity to the home networking, community and neighborhood networking, it is imperative to design an efficient resource management system for these networks. Routing is one of the most challenging issues in resource management for supporting real-time applications with stringent QoS requirements. However, most of the existing routing protocols for WMNs are extensions of protocols originally designed for MANETs and thus they perform sub-optimally.

The paper presents an efficient routing protocol for WMNs that is able to handle stringent QoS requirements of real-time applications. It utilizes two metrics- (i) minimum bandwidth (B_{min}) and (ii) maximum delay (T_{min}) that can be tolerated by the applications that it supports. While issues such as end-to-end route discovery and QoS-aware routing have been proposed for WMNs in [1], the protocol proposed in this paper is more efficient than those schemes as shown in the simulation results. The QoS-awareness in the protocol is achieved by a robust estimation of the available bandwidth of the wireless channel and a proactive discovery of the routing path. However, if the bandwidth estimation process indicates that it is not possible to guarantee the required QoS, the protocol does not admit the flow. In the event of any node or link failure, the protocol adapts itself and re-discovers an alternate path.

The key contributions of the paper are as follows: (i) It provides an accurate estimation of the end-to-end delay in a routing path. This estimated value is then used to check whether the routing can guarantee the application QoS. (ii) It computes a link quality estimator and utilizes it in route selection. (iii) It provides a framework for reliable estimation of available bandwidth in a routing path so that flow admission with guaranteed QoS satisfaction can be made. In addition the proposed protocol helps in identifying and isolating selfish nodes [2].

The rest of this paper is organized as follows. Section 2 describes related work on routing in WMNs. Section 3 describes the network model and the design objectives. Next, Section 4 describes the details of the proposed routing protocol. Simulation results are presented in Section 5. Section 6 concludes the paper.

2 Related Work

Although significant work has been done on routing in wireless networks, very little work has been done for WMNs. Most of the routing protocols for MANETs such as AODV and DSR use hop-count as the routing metric. However, this is not well-suited for WMNs. The basic idea in minimizing the hop-count is that it reduces delay and maximizes the throughput. But the assumption here is that the links in the path are either perfect or do not work at all, and all links are of equal bandwidth. A routing scheme that uses the hop-count metric does not take link quality into consideration. A minimum hop-count path has, on the average, longer links between the nodes present in the path compared to a higher hop-count path. This reduces the signal strength received by the nodes in that path and thereby increases the loss ratio at each link [3]. Hence, it is always possible that a two-hop path with a good link quality provides higher throughput than a one-hop path with a poor link quality. Moreover, wireless links usually have asymmetric loss rate [4]. Hence, new routing metrics based on link

quality are proposed such *expected transmission count* (ETX), *per-hop round-trip time* (RTT), and per-hop packet pair. Different approaches have been suggested by researchers for designing routing protocols for WMNs. In [5], a QoS routing over OLSR protocol has been proposed that takes into account metrics such as bandwidth and delay where the source node proactively changes a flow's next hop in response to the change in available bandwidth on its path. In [6], the authors have proposed a *link quality source routing* (LQSR) protocol. It is based on DSR and uses ETX as the routing metric. A new routing protocol called *multi-radio link quality source routing* (MR-LQSR) is proposed in [7]. The process of neighbor node discovery and propagation of link metric are same as those in DSR. However, assignment of link weight and computation of the path weight is different. A QoS enabling routing algorithm for mesh-based wireless LAN architecture has been proposed in [8] where the wireless users form an ad hoc peer-to-peer network. The authors also have proposed a protocol for MANET called *ad hoc QoS on-demand routing* (AQOR) [9]. However, both these algorithms are incapable of guaranteeing QoS [1]. In [10], the authors have shown that if a *weighted cumulative expected transmission time* [7] is used in a link state routing protocol, it does not satisfy the isotonicity property of the routing protocol and leads to formation of routing loops. To avoid routing loops, an algorithm is proposed that uses *metric of interference and channel switching* (MIC) as the routing metric.

In contrast to most of the above approaches, the proposed protocol performs on-demand route discovery using multiple metrics like bandwidth, delay, and reliability of the links and provides a routing framework that can support application QoS.

3 Network Model and Challenges

This section presents the model of the network and discusses the design key challenges for routing protocols in WMNs.

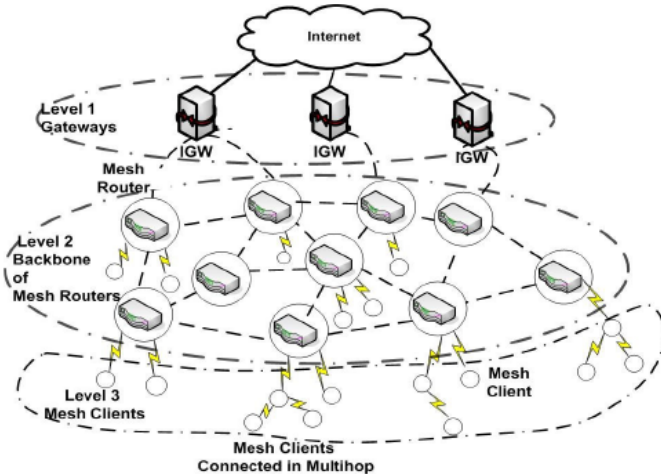


Fig. 1. The hierarchical architecture of a WMN

3.1 Design challenges

Two major design challenges for design of a mesh routing protocol are detecting and avoiding unreliable routes, and accurately estimating the end-to-end delay of a flow.

1. *Measuring link reliability*: It has been observed that in wireless ad hoc networks neighbor sensing with broadcast messages introduces *communication gray zones* [11]. In such zones, data messages cannot be exchanged although the *Hello* messages indicate neighbor reachability. This leads to a systematic mismatch between the route state and the real world connectivity, resulting in disruptive behavior. Since the routing protocols such as AODV and WMR [8] rely on control packets like RREQ, these protocols are highly unreliable for estimating the quality of wireless links. Due to communication gray zone problem, nodes that are able to send and receive bidirectional RREQ packets may not be able to send/receive high bit-rate data packets. These fragile links trigger link repairs resulting in high control overhead for routing protocols. Therefore, designing a robust and reliable link quality estimator is crucial for reducing the control overhead of a routing protocol in WMNs.

2. *End-to-end delay estimation*: An important issue in a routing protocol is end-to-end delay estimation. Current protocols estimate end-to-end delay by measuring the time taken to route RREQ and RREP packets along the given path. However, RREQ and RREPs are quite different from normal data packets, and therefore are unlikely to experience the same levels of traffic delay and loss as data packets. It has been observed through simulation that a RREP-based estimator overestimates while a hop-count-based estimator underestimates the actual delay experienced by the data packets [1]. The reason for the significant deviation of a RREP-based estimator from the actual end-to-end delay is interference of signals. The RREQ packets are flooded in the network resulting in a heavy burst of traffic. This heavy traffic causes inter-flow interference in the paths. The unicast data packets do not cause such events. Moreover, as a stream of packets traverse along a route, due to the broadcast nature of wireless links, different packets in the same flow interfere with each other resulting in per-packet delays. Since the control packets do not experience per-packet delay, the estimate based on control packet delay deviate widely from the actual delay experienced by the data packets. Computing a reliable estimate of the end-to-end delay is, therefore, another challenge in design of a routing protocol for WMNs.

4 The Proposed Routing Protocol

The goal of the proposed routing protocol is to establish a route from the source to the destination that can allow traffic flow to satisfy the QoS of the application, i.e., the flow should be able to deliver packets with a guaranteed minimum bandwidth B_{min} and maximum allowable end-to-end delay T_{max} .

The proposed protocol is a reactive routing protocol, in which during the route discovery phase, each intermediate node uses an admission control scheme to check whether the flow can be accepted or not. If a flow is accepted, an entry is created for the flow in table (called the *flow table*) maintained locally by the node. The admission control scheme used is the one proposed in [9]. The novel features of the proposed protocol are discussed in the remainder of this section.

4.1 Estimating Reliability of Routing Paths

Every node estimates the reliability of each of its wireless links to its one-hop neighbor nodes. For computing the reliability of a link, the number of control packets that a node receives in a given time window is used as a base parameter. An *exponentially weighted moving average* (EWMA) method is used to update the link reliability estimate. If the percentage of control packets received by a node over a link in the last interval of measurement of link reliability is N_t , and if N_{t-1} is the historical value of the link reliability before the last measurement interval, $\alpha = 0.5$ is the weighting parameter, then the updated link reliability (R) is computed as:

$$R = \alpha * N_t + (1 - \alpha) * N_{t-1} \quad (1)$$

Every node maintains the estimates of the reliability of each of its links with its neighbors in a *link reliability table*. The reliability for an end-to-end routing path is computed by taking the average of the reliability values of all the links on the path. Computation of the link reliability values are based on the RREQ packets on the reverse path and the RREP packets on the forward path. The use of routing path with the highest reliability reduces the overhead of route repair and makes the routing process more efficient.

4.2 Use of Network Topological Information in Route Discovery

The proposed protocol makes use of the knowledge of network topology by utilizing selective flooding of control messages in a portion of the network. In this way, broadcasting of control messages is avoided and thus the chances of network congestion and disruption to the flows in the network are reduced. If both the source and the destination are under the control of the same mesh router (Fig. 1), the flooding of the control messages are confined within the portion of the network served by the mesh router only. However, if the source and the destination are under different mesh routers, the control traffic is limited to the two mesh groups.

To further reduce the overhead of control message and enhance the reliability in routing, the nodes accept broadcast control messages from only those neighbors which have the link reliability value greater than 0.5 (i.e., on average 50% of the control packets sent from those nodes have been received by the node). This ensures that path with less reliability values are not discovered and therefore not considered for routing.

4.3 Estimating End-to-End Delay in a Routing Path

For accurate estimation of end-to-end delay in a routing path, an approach similar to the one proposed in [1] has been taken. For addressing the issue of differential delays experienced by the control and the data packets, the proposed protocol makes use of some *probe packets* during the route discovery phase. When a source node receives RREP packets from the destination in response to its RREQ, it stores in a table, the records for all the RREP packets together with the path through which the packets have arrived at it. However, instead of randomly selecting a path to send probe packets to the destination as suggested in [1], the packets are sent along the path from which the RREP messages have arrived at the source first. This ensures that the probe

packets are sent along the path which is likely to induce less end-to-end delay resulting in a better performance of the protocol as observed from the simulation results presented in Section 5. The probe packets are identical to data packets so far as their size, priority and flow rates are concerned. The objective of sending probe packets is to simulate the data flow and observe the delay characteristics in the routing path. Number of probe packets is kept limited to $2H$ for a path consisting of H hops to make a tradeoff between control overhead and measurement accuracy [1].

A destination node sets on a timer after it receives the first probe packet from the source node. The timer duration is based on the estimated time for receiving all the probe packets and is computed statistically. The destination computes the average delay experienced by all the probe packets it has received, and sends the computed value to the source node piggybacking it on a RREP message. If the computed value is within the limit of tolerance of the application QoS, the source selects the route and sends data packets. If the delay exceeds the required limit, the source selects the next best path (based on the arrival of RREP packets) from its table and tries once again. Since the routing path is set up based on the probe packets rather than the naïve RREP packets, the proposed protocol has higher route establishment. The proposed algorithm has higher setup time due to sending of the probe packets and selection of the best path based on the estimated end-to-end delay. However, since the selected paths have high end-to-end reliability, the delay and the control packet overhead are reduced because of minimal subsequent route breaks.

4.4 Estimation of Available Network Bandwidth

In addition to computation of path reliability and end-to-end delay, it is also necessary that the effective bandwidth in a routing path is reliably estimated. This is extremely important to support real-time applications since these applications require a guarantee for a minimum available bandwidth for their flows. In the proposed protocol the available bandwidth for a wireless link is estimated using its end-to-end delay and the loss of packets in it due to congestion. The packet-loss due to congestion in the link is estimated as follows.

In a wireless link packet loss may happen due to two reasons: (i) loss due to faulty wireless links and (ii) loss due to network congestion. The *radio link control* (RLC) layer segments an IP packet into several RLC frames before transmission and reassembles them into an IP packet at the receiver side. An IP packet loss occurs when any RLC frame belonging to an IP packet fails to be delivered. When this happens, the receiver knows the RLC frames reassembly has failed and the IP packet has been lost due to wireless error. Meanwhile, the sender detects *retransmission time out* (RTO) of the frame and discards all the RLC frames belonging to the IP packet. This enables the sender to compute packet drop rate in the wireless links. Moreover, using the sequence numbers of the IP packets received at the receiver, it is possible to differentiate the packet loss due to link error and packet loss due to congestion [Yang paper]. For example, while receiving two incoming packets with sequence number i and $i + 2$, if the receiver finds an IP packet assembly failure in RLC layer, the packet with sequence number $i + 1$ is lost due to wireless channel. Once the packet loss ratio due to congestion ($P_{congestion}$) is estimated, the available bandwidth in the wireless link, *estrat*, is computed as follows [12]:

$$estrat = \frac{PacketSize}{X + Y} \quad (2)$$

where,

$$X = RTT \sqrt{\frac{2P_{congestion}}{3}} \quad (3)$$

$$Y = RTO * \min(1, 3 * \sqrt{\frac{3P_{congestion}}{8}} P_{congestion} (1 + 32P_{congestion}^2)) \quad (4)$$

In (2), RTT is the average round trip time for a control packet. RTO is the retransmission time out for a packet, and is computed using (5).

$$RTO = \overline{RTT} + k * \overline{RTT}_{var} \quad (5)$$

where, \overline{RTT} and \overline{RTT}_{var} are the mean and variance respectively of RTTs and k is set to 4. This bandwidth estimator is employed to dynamically compute the available bandwidth in the wireless links on a routing path so that the guaranteed minimum bandwidth for the flow is always maintained throughout the application life-time.

4.5 Identifying Selfishness in Wireless Nodes

The proposed routing protocol also enforces cooperation among the nodes by identifying the selfish nodes in the network and isolating them. Selfishness is an inherent problem associated with any capacity-constrained multi-hop wireless networks like WMNs. A mesh router can behave selfishly owing to various reasons such as: (i) to obtain more wireless or Internet throughput, or (ii) to avoid path congestion. A selfish mesh router increases the packet delivery latency, and also increases the packet loss rate. A selfish node while utilizing the network resources for routing its own packet, avoids forwarding packets for others to conserve its energy. Identification of selfish nodes is therefore, a vital issue.

Several schemes proposed in the literature to mitigate the selfish behavior of nodes in wireless networks such as credit-based schemes, reputation-based schemes and game theory-based schemes [13]. However, to keep the overhead of computation and communication at the minimum, the proposed protocol employs a simple mechanism to discourage selfish behavior and encourage cooperation among nodes. To punish the selfish nodes, each node forwards packets to its neighbor node for routing only if the link reliability of that node is greater than a threshold (set at 0.5). Since the link reliability of a selfish node is 0, the packets arriving from this node will not be forwarded. Therefore, to keep its link reliability higher than the threshold, each node has to participate and cooperate in routing. The link reliability serves a dual purpose of enhancing reliability and enforcing node cooperation in the network.

4.6 QoS Violation and Recovery

The proposed protocol detects failure to guarantee QoS along a path with the help of *reservation timeouts* in flow tables records maintained in the nodes, and by detection of non-availability of minimum bandwidth as estimated along its outbound wireless

link. Failure to guarantee QoS may occur in three different scenarios. In the first case, a node receives a data packet for which it does not find a corresponding record in its flow table. This implies that a reservation time-out has happened for that flow. The node, therefore, sends a *route error* (RERR) to the source which re-initiates route discovery. In the second scenario, a destination node detects from its flow table records that the data packets received have exceeded the maximum allowable delay (T_{max}). To restore the path, the destination broadcasts a new RREP back to the source, and the source starts rerouting the packets via the same path of the RREP. This approach is similar to the one proposed in [9]. In the third case, an intermediate node on the routing path may find that the estimated bandwidth (using (2)) in its forwarding link is less than the guaranteed minimum (B_{min}) value. In this case, the intermediate node sends a route error (RERR) to the source which re-initiates the route discovery process. The real-time estimation of the bandwidth in the next-hop wireless link at each node on the routing path makes the proposed protocol more robust and reliable compared to other existing routing protocols for WMNs. For example, the protocol presented in [1] does not provide any bandwidth estimation mechanism at intermediate nodes, and therefore cannot provide delivery of all packets for every admitted flow in the system.

5 Simulation Results

The proposed protocol has been implemented in Qualnet network simulator [14]. The WMN topology is shown in Fig. 2. The simulated network consists of 50 nodes of which 5 are mesh routers and remaining nodes are mesh clients. Each mesh router has 9 mesh clients associated with it. To compare the performance of the proposed protocol with the protocol presented in [1] by Kone et al, parameters in the simulation are chosen to be identical as those used in [1]. The simulation parameters are presented in Table 1. The MAC protocol used is 802.11b, which is based on CSMA/CA protocol. No wireless fading model is used scheme to keep it simple, yet realistic.

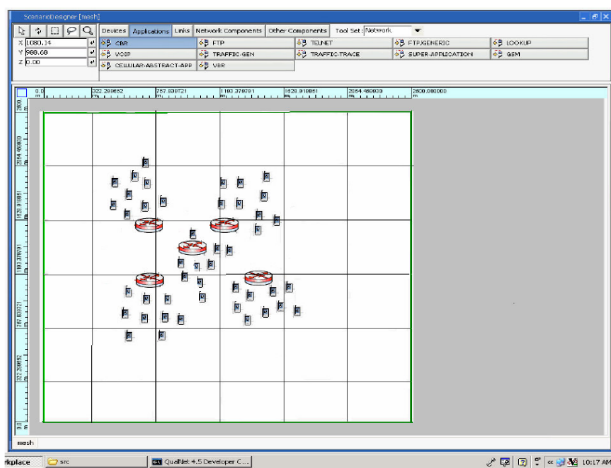


Fig. 2. The network topology of the simulated WMN

Table 1. Simulation parameters

Parameter	Values
Simulation area	1300 m * 1300 m
MAC protocol	802.11
PHY-model	802.11b
Propagation channel frequency	2.4 GHz
Raw channel bandwidth	11 Mbps
Simulation duration	15 s
Traffic type	CBR UDP
Packet size	512 bytes
Max. no. of packets from a source node	10000
Hello packet broadcast interval	200 ms
Wireless fading model	None
Node mobility	None
Propagation limit of each node	-111.dBm
IP queue scheduler	Strict priority

5.1 QoS Routing Behavior

To evaluate the effectiveness of the admission control module of the protocol, the performance of the protocol in a heavy traffic condition is observed. Due to heavy traffic, some flows in the network will not have their bandwidth and end-to-end delay requirement satisfied, and hence will be rejected by the protocol. AODV, on the other hand being non-aware about the QoS parameters, will try to route packets from all the sources. The experimental results have shown that the proposed protocol has rejected admission of new flows after the traffic level in the network has gone above a certain limit. However, for flows that have been admitted, all packets are routed within the required end-to-end delay. AODV, on the other hand has admitted all the flows and routed them with increasing delay as the congestion in the network increased to due more number of control packets. An important point worth mentioning here is that the protocol in [1] is not able to achieve 100% packet delivery ratio for all flows. In contrast, the proposed protocol has been found to achieve 100% packet delivery. This is due to bandwidth estimation at each intermediate node as discussed in Section 4.4.

5.2 Control Overhead

For computing the overhead due to control packets, only the RREQ messages and the probe packets in the proposed algorithm are considered since these are broadcast messages and largely contribute to the control overhead of the protocols.

Fig. 3(a) shows the overhead of RREQ packets in AODV and the proposed protocol for different data rates. The control overhead of the proposed protocol is first evaluated only with the RREQ packets and then with RREQ packets and the probe packets together. This also gives an idea of the additional overhead introduced due to the probe packets. It can be easily observed that the proposed protocol has a very limited overhead even with the probe packets compared to naïve AODV protocol. It is also worth observing that the proposed protocol has about 20% less control overhead than the protocol proposed in [1] due to its robust bandwidth estimation of links.

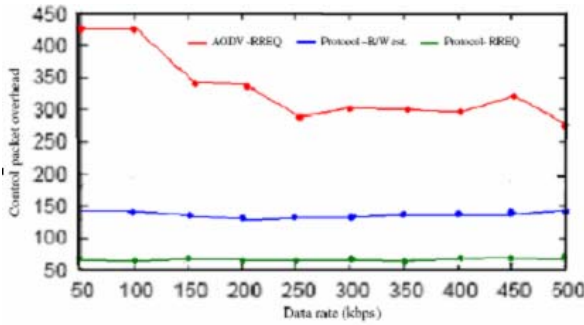


Fig. 3(a). Control overhead vs. data rate

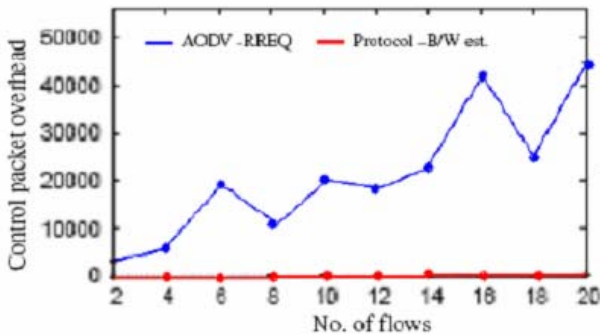


Fig. 3 (b). Control overhead with no. of flows

Fig. 3(b) shows that AODV has a very high overhead of control packets for increasing number of flows in the system. AODV always tries to establish routing paths based on minimum hop-counts. It does not pay any attention to the link reliability. This results in frequent selection of unreliable links and consequent link break and re-discovery of route. This leads to large overhead of control information. In contrast, the proposed protocol has a very limited control overhead since paths with higher link reliability are only used for routing. The algorithm proposed in [1] has similar performance as the proposed protocol in this case.

5.3 End-to-End Delay Estimation

To demonstrate the effectiveness of the end-to-end delay estimation mechanism by probe packets, delays estimated by naïve RREP approach and the probe packet approach are compared with the actual end-to-end delay in the routing paths. The records are observed for different flow rates with each flow having a minimum bandwidth requirement of $B_{min} = 50$ kbps. Fig. 4 shows that the probe packet-based mechanism very accurately estimates the actual delay, and has similar performance as that of the algorithm in [1].

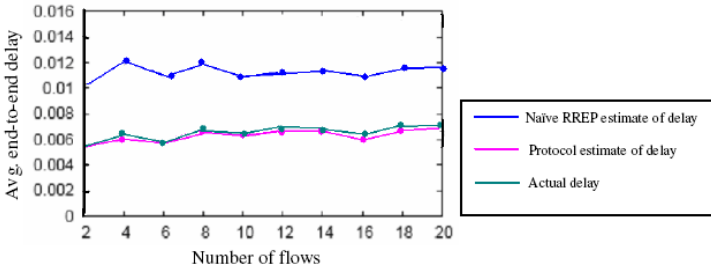


Fig. 4. End-to-end delay estimation for different no. of flows

5.4 Detection of Selfish Nodes

As mentioned in Section 4.4, the proposed protocol has the ability to detect selfish nodes which try to use network resources without contributing to the cooperative framework of routing. To evaluate its detection capability of selfish nodes, two types of flows are distinguished: *selfish* and *honest*. A flow is considered selfish if either its source or destination is a selfish node, otherwise the flow is considered to be honest. Some mesh clients are selected randomly and configured as selfish nodes. In each of the 20 run of the simulation, 10 flows of 50 kbps data rate are generated randomly.

It may be observed from Fig. 5 that AODV cannot restrict the traffic along the selfish flows. The selfish nodes can fully exploit the routing process to have their packets routed in the network. However, the proposed protocol reduces the flows along the selfish paths. In fact, the performance of AODV is not very much affected by presence of selfish nodes, since it never establishes routing path based on hello packets. Since the proposed protocol establishes route based on hello packets received form neighbors, its performance is affected by the presence of selfish nodes. However, its performance is not substantially affected, since most in most cases, these nodes are not allowed to participate in the routing, because of the low values of their link reliability. It may also be mentioned that the proposed protocol is able to maintain, on average, 35% more throughput in presence of selfish nodes when compared with the protocol proposed in [1]. The large difference is due to its ability to detect selfish nodes faster by its effective bandwidth estimation in the route where non-forwarding of packets is treated as packet drops due to congestion.

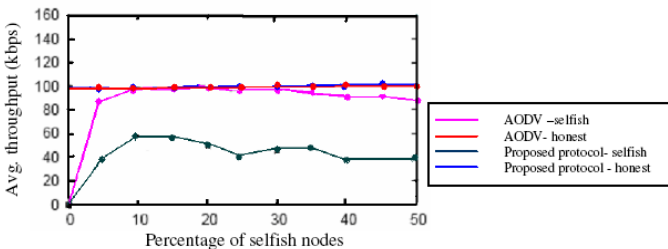


Fig. 5. Avg. throughput with different protocols for varying percentage of selfish nodes

6 Conclusion

WMNs present a promising technology for providing low-cost last mile broadband Internet connectivity through multi-hop communications. Designing an efficient and reliable routing protocol for WMNs is a challenging issue. This paper has presented a QoS-aware routing protocol that is capable of supporting real-time applications. By robust estimation of available bandwidth in the wireless route and a reliable computation of end-to-end delay along a path, the protocol is able to sustain a high level of throughput even in presence of a large proportion of selfish nodes in a WMN. Moreover, it is able to reduce control overhead substantially by exploiting the knowledge of network topology. Simulation results have shown that the proposed protocol is more efficient than some of the current routing protocols for WMNs.

References

1. Kone, V., Das, S., Zhao, B.Y., Zheng, H.: Quorum: Quality of Service in Wireless Mesh Networks. *Journal of Mobile Networks and Applications* 12(5), 358–369 (2007)
2. Mahajan, R., Rodrig, M., Wetherall, D., Zahorjan, J.: Sustaining Cooperation in Multi-Hop Wireless Networks. *Proc. of NSDI 2*, 231–244 (2005)
3. Cuto, D.S.J.D., Aguayo, D., Bricket, J., Morris, R.: A High-Throughput Path Metric for Multi-Hop Wireless Routing. In: *Proc. of ACM MOBICOM*, pp. 134–146 (2003)
4. Aguayo, D., Bicket, J., Biswas, S., Judd, G., Morris, R.: Link-Level Measurements from an 802.11b Mesh Network. In: *Proc. of ACM SIGCOMM*, pp. 121–132 (2004)
5. Badis, H., Gawedzki, I., Al Agha, K.: QoS Routing in Ad Hoc Networks using QOLSR with no need of Explicit Reservation. In: *Proc. of VTC (September 2004)*
6. Draves, R., Padhya, J., Zill, B.: Comparison of Routing Metrics for Static Multi-Hop Wireless Networks. In: *Proc. of ACM SIGCOMM*, pp. 133–144 (2004)
7. Draves, R., Padhye, J., Zill, B.: Routing in Multi-Radio, Multi-Hop Wireless Mesh Networks. In: *Proc. of ACM MOBICOM*, pp. 114–128 (2004)
8. Xue, Q., Ganz, A.: QoS Routing for Mesh-Based Wireless LANs. *International Journal of Wireless Information Networks* 9(3), 179–190 (2002)
9. Xue, Q., Ganz, A.: Ad Hoc QoS On-Demand Routing (AQOR) in Mobile Ad Hoc Networks. *Journal of Parallel and Distributed Computing* 63(2), 154–165 (2003)
10. Yang, Y., Wang, J., Kravets, R.: Interference-Aware Load Balancing for Multi-Hop Wireless Networks. Technical Report UIUCDCS-R-2005-2526, Dept. of Computer Science, University of Illinois at Urbana-Champaign (2005)
11. Lundgren, H., Nordstrom, E., Tschudin, C.: The Gray Zone Problem in IEEE 802.11b Based Ad Hoc Networks. *MC2R* 6(3), 104–105 (2002)
12. Yang, F., Zhang, Q., Zhu, W., Zhang, Y.-Q.: End-to-End TCP-Friendly Streaming Protocol and Bit Allocation for Scalable Video over Wireless Internet. *IEEE Journal on Selected Areas in Communications* 22(22), 777–790 (2004)
13. Santhanam, L., Xie, B., Agrawal, D.: Selfishness in Mesh Networks: Wired Multi-Hop MANETs. *IEEE Journal on Wireless Communication* 15(4), 16–23 (2008)
14. Qualnet Simulator, <http://www.scalable-networks.com>

A Context-Aware Service Model Based on Workflows for u-Agriculture^{*}

Yongyun Cho¹, Jongbae Moon^{2,**}, and Hyun Yoe¹

¹ Information and Communication Engineering, Suncheon National University,
413 Jungangno, Suncheon, Jeonnam 540-742, Korea

{yycho, yhyun}@suncheon.ac.kr

² Supercomputing Center, Korea Institute of Science and Technology Information,
52-11 Eoeun-dong, Yuseong-gu, Daejeon 305-806, Korea
comdoct@shinbiro.com

Abstract. In agricultural processes, many works may need to be automated, because those are very hard labors or time-consuming jobs for farmers. Workflow technologies, which have successfully been a good model for a service automation in various computing environments, can be used as a possible service automation model in agriculture. A workflow for u-agriculture may need various contexts sensed from real sensor networks for service automation. Recently, many researches have applied workflow technologies into the various fields such as smart home, u-health care, u-city, u-port, and u-agriculture. However, many current workflow service have difficulty to control work processes and execute services according to context information in u-agricultural environment. This paper proposes a context-aware service model for u-agriculture, which is based on workflows and is aimed at supporting smart workflow services based on ubiquitous sensor networks in u-agriculture. With the proposed context-aware service model, developers can easily integrate various service demands into a service workflow, and can easily develop a context-aware workflow service for u-agriculture. Therefore, the proposed service model can be greatly helpful in the development of smart applications or the work automation in the fields of u-agriculture.

1 Introduction

Service automation is commonly centered in the service execution without intervention of humans according to the changes of surrounding situation information [1]. Because agricultural environment is hard to work and labor-intensive in most case and areas, various efforts for the work automation and smart cultivation is invaluable. In recent years, because workflow is a very good model for service

^{*} This research was supported by the MKE(Ministry of Knowledge Economy), Korea, under the ITRC(Information Technology Research Center) support program supervised by the IITA(Institute of Information Technology Advancement)” (IITA-2009-(C1090-0902-0047)).

^{**} Corresponding author.

automation in the business computing environments, some studies have tried to use workflow techniques to automate services according to context information or situation information in ubiquitous computing environment or pervasive computing environment [2,3,4]. And, some researches have tried to apply workflow and/or context-aware technologies in the field of agriculture [5,6,7].

A context-aware service is a kind of a smart service which can be automatically processed according to various situation information. A context in a ubiquitous environment means any information that can be used to characterize the situation of an entity [11,2]. For example, in u-agricultural environments, contexts may be the plant's rate of growth, the amount of sunshine, the percentage of humidity and the indoor temperature of greenhouse and so on. These contexts can be an important service execution information for work automation or smart cultivation in u-agriculture environment. Therefore, context-aware service for u-agriculture can consider contexts around agricultural environment with any work automation model.

This paper introduces a new context-aware service model based on workflow as a service model for work automation in u-agriculture environment. The suggested approach uses uWDL [2], which is an academic context-aware workflow language, as a workflow model for work automation. So, the developers can easily describe context information directly into a workflow service scenario, and can correctly execute the best one of the services described into the workflow according to the contexts. Therefore, with the suggested model, context-aware service developers can easily development various context-aware applications in the fields of u-agriculture environments. And, because a developed workflow can be easily reused for other context-aware services, the developers can raise development efficiency and reusability.

2 Related Work

2.1 Context-Aware Workflow Languages

Workflow languages such as XLANG [8], WSFL [9] and WS-BPEL were developed by international IT companies such as Microsoft, IBM, Bea, Oracle etc. in order to support service automation in business processes. The BPEL4WS [10] is a workflow language based on web-service technologies and includes the advantages of both WSFL and XLANG technologies. XLANG is one of the initial workflow languages using a directed graph approach and WSFL is a block-structured workflow language for describing web services. However, They do not have any element in their language schema or grammar to describe context information as transition conditions of services. The workflow languages can use simple return values through XPath, XLink, Xpoint to express current service's transition conditions. Therefore, they are not suitable for a workflow language for context-aware services in ubiquitous environment.

uWDL [2] is a workflow language for ubiquitous computing environment. It includes <constraints>, <contexts>, and <profiles> elements to describe various contexts as service transition conditions in its language schema. And, For

the automatic execution of context-aware workflow services, it has also `flows` element, which supports various service flow patterns such as a single-flow, a multi-flow, and a sub-flow. Even if uWDL is a kind of academic and experimental context-aware workflow language, it can enough describe contexts as transition conditions of services directly in a uWDL workflow scenario. To do that, it offers a structural context model based on RDF [11] triplet-subject, verb, and object.

2.2 Context-Aware Workflow Systems

A workflow system manages and controls flows of subtasks using state transition constraints specified in a workflow language. WorkSco [12] is an situation-adaptable workflow system that can support service demands generated dynamically in a business process. It is based on a micro workflow model, a dynamic evolution and an open-point adaptation techniques to dynamically handle user's requests, which may be generated in various business domains. However, it does not yet give an explicit method to do that. Even though WorkSco considers dynamic handling for user's requests in a workflow system, because it does not consider situation information or contexts as user's requests, it is basically not adequate for ubiquitous computing environments.

CAWE [13] is a framework for the management of context-aware workflow systems based on Web Services. CAWE supports a synthetic and extensible specification of context-sensitive workflows, which can be executed by standard workflow engines. The system explained an initial prototype of the CAWE framework using jBPM, which is a business process management system based on Petri Netmodel implemented in Java. However, even if it proposed a systematic and valuable architecture for context-aware workflow services, it does not offer an explicit method to express the contexts into a workflow scenario as service execution condition, and doesn't enough describe how it can use user's situation information or contexts, which can be dynamically occurred in real ubiquitous environments.

uFlow [14] is an another ubiquitous workflow framework to support a context-aware service based on a uWDL. uFlow is based a workflow scenario, and offers a method to handle the changes of user's demands or situation information which may be dynamically generated during service processing. POESIA [7] is a result of the effort to use a ontology-based workflow to compose Web services in agriculture. The approach described how to compose Web services, using domain-specific multidimensional ontologies in the field of agricultural environment. POESIA is invaluable for introducing a systematic model or method to use ontologies, web-services, and workflow technologies to support context-aware services in agriculture. However, it seems to be mainly aimed to define Web services using ontologies, but it is not a main purpose to describe contexts into a workflow scenario, and to control context-aware service flows according the scenario and real context information for work automation in agriculture. Therefore, it may be inappropriate to intactly use POESIA or uFlow as a service model for automatic and context-aware services in agriculture.

3 A Model of Context-Aware Workflow Services for u-Agriculture

3.1 A Suggested Service Model Architecture

All things considered for a context-aware service in u-agriculture, it is desirable to be autonomic and context-aware. So, a context-aware service model for u-agriculture has to offer a good model to control execution of services for work automation and a method to handle context information for a smart work without human's intervention.

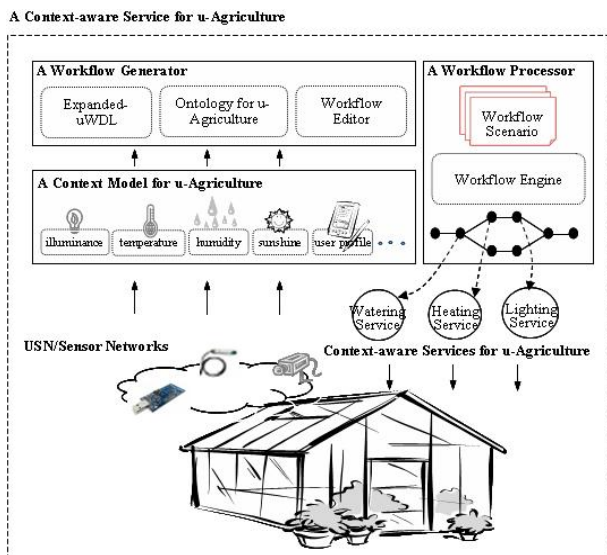


Fig. 1. An architecture for the suggested context-aware service model.

In Figure 1, the suggested service model architecture has a context model, a workflow generator and a workflow processor. The workflow model is based on a RDF-based structural context model. The situation information received from USNs/Sensor Networks the model is transformed as a context data through the workflow model. The workflow generator includes an expanded-uWDL as a automation model for service execution, a OWL-based ontology for u-agricultural environment, and a GUI-based workflow editor for generating a context-aware workflow scenario. With the workflow generator, a developer can explicitly describe the contexts into a workflow scenario as service execution information. The workflow processor is to recognize a workflow scenario and to process services described into the scenario. In this case, the workflow processor determines whether a service defines the sensed context as its execution condition and controls the service flows according to the sensed contexts.

3.2 A Context Model for u-Agriculture

Contexts in u-agriculture may be very various and mainly sensor-based. For example, in u-agricultural environments contexts may vary from the plant’s rate of growth, the amount of sunshine, the percentage of humidity and the indoor temperature of greenhouse to a farmer’s job schedule and his individual profile information. In Figure 1, the module recognizes the contexts sensed from USNs or sensor networks in u-agricultural environment.

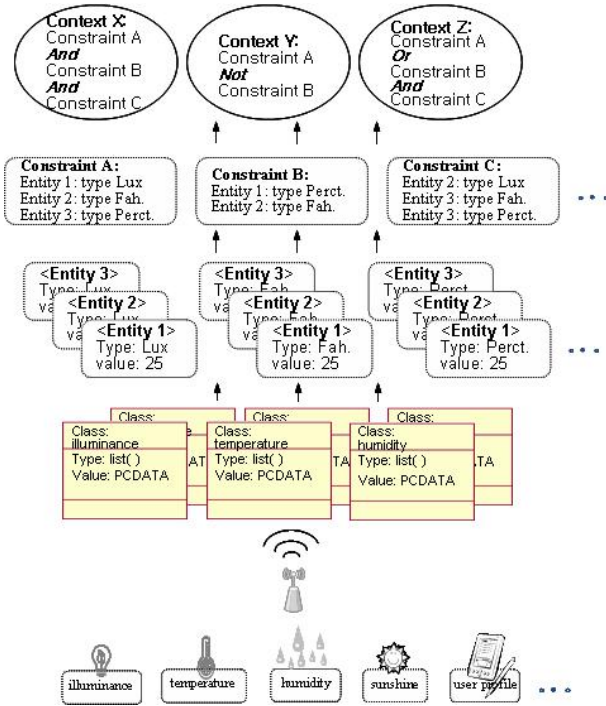


Fig. 2. A context model for u-agriculture

A service in u-agriculture has to be executed according to contexts so that it becomes context-aware. Commonly, A data in real agricultural environment may occur mainly in real-time and is continuously transferred from various kinds of sensors. Therefore, its value and type need to be transferred together for any context-aware service. For example, if the numeric value 60 is transferred from anyone of temperature sensors in a greenhouse, it needs to show that a data type of it is Fahrenheit and a data value is 60. It means that the value is related with just temperature, not anything else. In a context-aware service, the additional semantic information can be meaningful as a data unit to make more conceptual and higher level data. Figure 2 shows a context model in this paper.

In Figure 2, the data set that consists of a value and type is an entity. A few of entities form a constraint, which is based on RDF triplet. And, a context includes several constraints using rule-based operation. For example, let suppose that a numeric value 60 is sensed from the current temperature of a section A in a greenhouse 1. A farmer may wants a service scenario for a heating service, that when the current temperature of a section A in a greenhouse 1 is less than Fahrenheit 50, start a heating system. In that case, a context may needs constraints consisting of 3 entities, $\langle \text{Type:Fahr.}, \text{value: } 60 \rangle$, $\langle \text{Type:greenhouse}, \text{value: no.1} \rangle$, and $\langle \text{Type:Position.}, \text{value: } (x, y) \rangle$. Therefore, the farmer only has to describe the context as execution conditions of the heating service into his workflow scenario. Then, when the workflow processor in Figure 1 finds the context among numerous entities sensed from a real environment, it will start the heating service.

3.3 Workflow Patterns for Context-Aware Services

Figure 3 describes workflow service patterns [15] added into the uWDL schema. The suggested service model uses a workflow model as an work automation model based on uWDL.

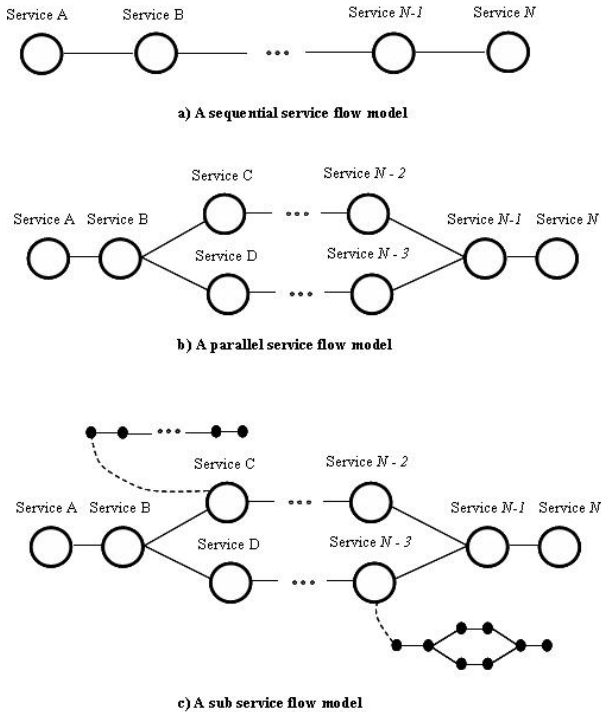


Fig. 3. Workflow patterns based on the extended uWDL

However, uWDL didn't support enough workflow patterns to satisfy various flow pattern of context-aware services. And, uWDL has difficulties in composing a new workflow using some of workflows, because it provide only a sequential flow pattern and a simple workflow pattern. So, this paper expands the schema of uWDL. A workflow based on the expanded uWDL can support various service flows and is reused as a part of a new workflow. This can raise development efficiency and workflow reusability.

In Figure 3, (a) represents a sequential flow model. It means a deterministic flow that various service nodes are connected in a sequence. However, in many cases of real u-agricultural environment, the sequential flow model is not enough to meet various service requirements. Figure 3(b) is a parallel service flow is based on the basic workflow patterns such as the parallel split, synchronization and simple merge pattern. This means that a workflow is deterministic, but it has the flows that split into two or more. Parallel workflow model also is used to build more complicated workflow services. This means that it is able to be used to execute two or more workflow services independently and concurrently. In Figure 3(c), a subflow model is based on both exclusive choice and synchronization pattern. It means that other workflows, which may are developed already, can be called and attached on one of the nodes on flow B. So, with the pattern, a developer can quickly and easily compose new context-aware service flows by reusing pre-developed valuable workflows. In Figure 3(c), the service 3 and $N-1$ represent the workflow reusing with the subflow model.

The flow model (a), (b), and (c) in Figure 3 can be composed for more complicated and large service flows. There are various service domains in ubiquitous computing environments. Thus, service developer must be able to provide multiple workflow services suitable to users by describing many context-aware workflow documents. Each workflow service document can be used independently for a user or be composed in a shape of a new integrated workflow for many users. And service developer is able to express a composite workflow service, which is composed of some of workflows, by reusing a workflow document or more.

4 A Experimental Scenario for Context-Aware Services in u-Agriculture

For an experiment using the suggested service model, we will take experimental scenario for context-aware Services in u-agriculture using the expanded uWDL. To do this, let's suppose a work scenario in u-agricultural environment. The scenario is as like this:

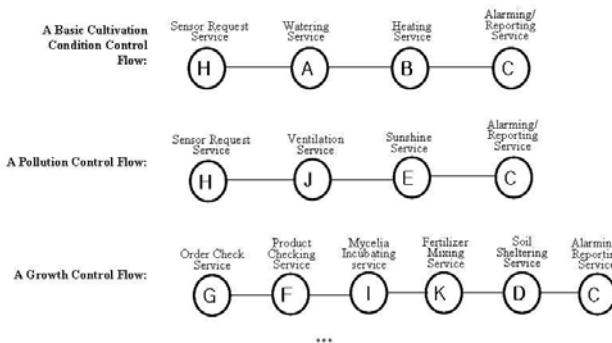
There is a farmer to cultivate of mushrooms in his/her greenhouse. Mushrooms is one of crops that are very sensitive to temperature, humidity, and PH of soil. From morning to midnight everyday, the farmer checks the cultivation factors for his mushrooms. And, the farmer has to take proper works according to the results of the examination. Furthermore, he/she has to incubate the mycelia of mushrooms and shelter the cultivation soil for the continuous cultivation of mushrooms. And, the farmer has to control the speed and the amount of the

works according the market demands and the production of mushrooms. Clearly, it is very labor-intensive jobs, and the farmer will want to reduce the amount of the works which he/she has to do. To do this, let's suppose a next scenario. The scenario is based on a workflow which describes the works as the services. Figure 4 shows the parts of the required services and the workflow scenario.

In Figure 4(a), the services are based on Web services and exist independently. The service H may request the current sensed values of humidity, temperature, PH of soil, a pollution degree of air in the greenhouse, sunshine, and illuminance, to various sensors positioned around the mushroom greenhouse. In Figure 4(b) shows the possible service flows that can be connected with the services of Figure 4(a) in order to achieve a specific work. The service H is needed in the basic cultivation control flow and the pollution control flow. The service C is used in all of the three flows in Figure 4(b). Of course, There can be other possible and reasonable flows in addition to those in Figure 4. And, the procedural orders of the flows in Figure 4 don't have to be absolutely like that.

Watering Service (A)	Heating Service (B)	Alarming/Reporting Service (C)	Soil Sheltering Service (D)	Sunshine Service (E)	Product Checking Service (F)
<type: Hum., value: id_1> <type: Log_op, value: lt> <type: Percr, value: 30> AND <type: Order.return, value: False> ...	<type: Temp., value: id_1> <type: Log_op, value: gt> <type: Fahr., value: 60> AND <type: Order.return, value: False> ...	<type: Rec., value: True > AND <type: Serv.name, value: PCDATA> ...	<type: Ph, value: id_1> <type: Log_op, value: lt> <type: Ph, value: 90> AND <type: Order.return, value: False> ...	<type: Bright, value: id_1> <type: Log_op, value: lt> <type: Lux, value: 50> ...	<type: Prod., value: id_1> <type: Log_op, value: gt> <type: Kg, value: 3000> ...
Order Check Service (G)	Sensor Request Service (H)	Mycelia Incubating service (I)	Ventilation Service (J)	Fertilizer Mixing Service (K)	
<type: Order_idg, value: id_1> <type: Log_op, value: ge> <type: Serv.F.return, value: var_1> ...	<type: Time, value: var_current_time> <type: Log_op, value: eq> <type: Mill_time, value: 08:00> OR <type: Period.return, value: True> ...	<type: Serv.F.return, value: var_1> OR <type: Serv.G.return, value: var_1 > ...	<type: CO ₂ , value: id_1> <type: Log_op, value: gt> <type: PPM, value: 3000> OR <type: Serv.L.return, value: True> ...	<type: Month, value: id_1> <type: Log_op, value: eq> <type: Day, value: 30> OR { <type: Serv.F.return, value: var_1> AND <type: Serv.G.return, value: var_1 > } ...	

a) A part of the demanded services



b) A part of the possible service flows

Fig. 4. The parts of the demanded services and the possible service flows

Now, let's make a context-aware workflow service model by using the three flows with contexts, which may occur in the greenhouse. Figure 5 represents the context-aware workflow service model. Figure 5(a) shows a possible composition of the basic cultivation condition control flow and the pollution control flow in Figure 4. The service H is needed for the basic cultivation condition control flow and the pollution control flow.

The service H selects a next service flow according to contexts that the service received from the requested sensors. For example, in case of a context data that can provoke the basic cultivation services is sensed from the sensors, the service H will move to the service A. And, when the sensed data is related with the pollution control services, the service H will move to the service J. Of cause, if the contexts for the two flows occur all together, the service H will go to the two flows in a type of parallel. However, the ventilation service or the sunshine

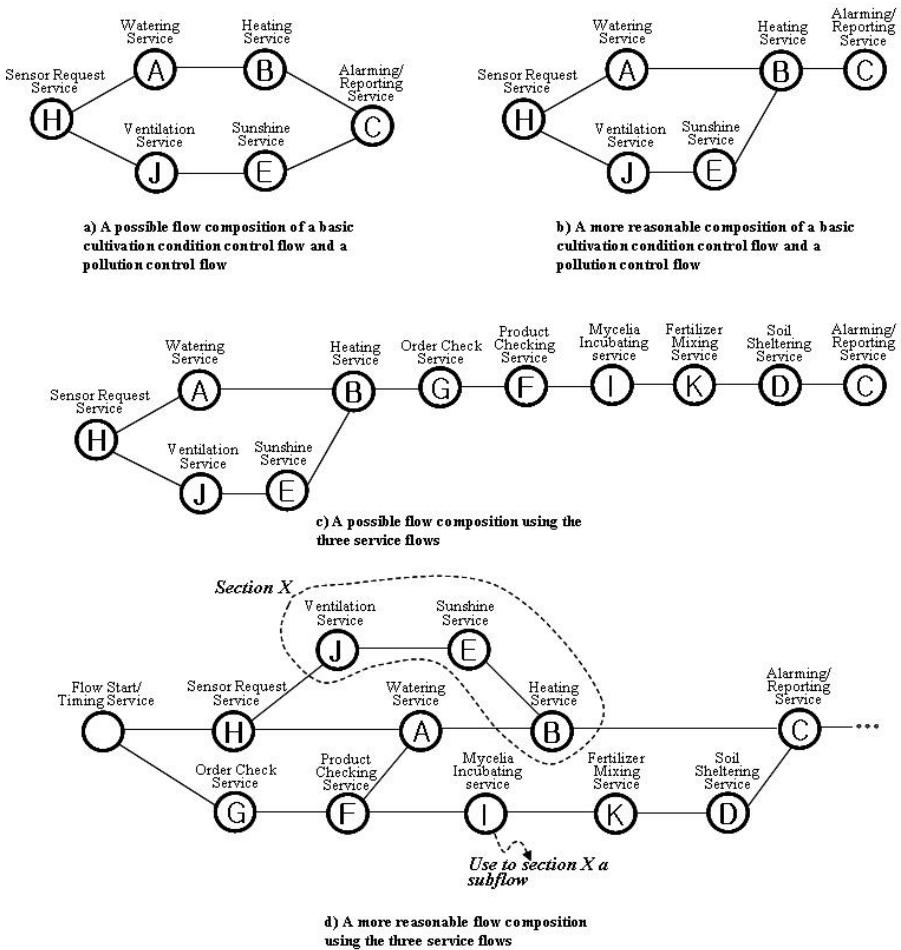


Fig. 5. The possible composition for a context-aware service in u-agriculture

service to purify the air of the greenhouse has to be executed just before the heating service. So, it may be more reasonable that the service E connects to the service B just before. Figure 5(b) shows that. Now, the growth control flow has to be added to the Figure 5(b). First, the sequential flow like Figure 5(c) can be considered simply. However, because the order check service or the product checking service can affect the watering service or heating service, it is desired that the services are started as soon as possible. And, the service I and K may probably generate air pollution of the greenhouse. Therefore, those service can reuse the existing service flow, which consists of J, E, and B, as their subflow. The service J in Figure 5(d) illustrates the process to reuse the section X in Figure 5(c) as a subflow. The expended uWDL can grammatically support the method to reuse any partition of the existing service flows as a subflow.

5 Conclusion

A research for service automation and context-aware services is a valuable subject in u-agricultural environment like in the other fields. However, most existing context-aware service models are not suitable for work automation and smart cultivation in u-agriculture, because they have troubles to use or express contexts to automatically execute a service. This paper introduced a context-aware service model based on workflows for u-agriculture. The suggested service model offers a structural context-model to describe contexts in u-agriculture through a entity, a RDF-based constraint, and rule-based operators. And, it offers flow models of an expanded uWDL as a service automation model. With the proposed context-aware service model, developers can easily develop various context-aware service flows into a service workflow for u-agriculture. And, a developer can reuse the existing service flows in composition of a new service flow with the flow models. Therefore, the proposed service model is expected to reduce the time and effort to develop context-aware applications or the work automation in the fields of u-agriculture.

References

1. Abowd, G.D., Dey, A.K., Brown, P.J., Davies, N., Smith, M., Steggles, P.: Towards a Better Understanding of Context and Context-Awareness. In: HUC 1999, pp. 304–307 (1999)
2. Han, J., Cho, Y., Choi, J.: Context-Aware Workflow Language based on Web Services for Ubiquitous Computing. In: Gervasi, O., Gavrilova, M.L., Kumar, V., Laganá, A., Lee, H.P., Mun, Y., Taniar, D., Tan, C.J.K. (eds.) ICCSA 2005. LNCS, vol. 3481, pp. 1008–1017. Springer, Heidelberg (2005)
3. Aqeel-ur-Rehman, Shaikh, Z.A.: Intelligent Workflow: A Step towards Intelligent Agent based Academic Workflow System. In: 2nd All Pakistan UTECH and ACM Technical Conference (UTECH 2007), Karachi, Pakistan, November 23-24 (2007)
4. Tang, F., Guo, M., Dong, M., Li, M., Guan, H.: Towards Context-Aware Workflow Management for Ubiquitous Computing. In: Proceedings of ICES 2008, pp. 221–228 (2008)

5. Lauser, B., Sini, M., Liang, A., Keizer, J., Katz, S.: From AGROVOC to the Agricultural Ontology Service / Concept Server - An OWL model for creating ontologies in the agricultural domain. In: Proceedings of the OWLED*06 Workshop on OWL: Experiences and Directions, Athens, Georgia (USA), pp. 10–11 (2006)
6. Wark, T., Corke, P., Sikka, P., Klingbeil, L.: Transforming Agriculture through Pervasive Wireless Sensor Networks. *IEEE Pervasive Computing*, 50–57 (2007)
7. Fileto, R., Liu, L., Pu, C., Assad, E.D., Medeiros, C.B.: POESIA: An Ontological Workflow Approach for Composing Web Services in Agriculture. *The VLDB J.* 12(4), 352–367 (2003)
8. Satish, T.: XLANG: Web Services for Business Process Design. Microsoft Corp. (2001)
9. Leymann, F.: Web Services Flow Language (WSFL 1.0), Distinguished Engineer Member IBM Academy of Technology, IBM Software Group (May 2001)
10. Tony, A., Francisco, C., et al.: Business Process Execution Language for Web Services, BEA Systems. Microsoft Corp. IBM Corp., Version 1.1 (2003)
11. W3C: RDF/XML Syntax Specification, W3C Recommendation, 30-39 (2004/2006)
12. Vieira, P., Rito-Silva, A.: Adaptive Workflow Management in WorkSCo. In: 16th International Workshop on Database and Expert Systems Applications (DEXA 2005), pp. 640–645 (2005)
13. Ardissono, L., Furnari, R., Goy, A., Petrone, G., Segnan, M.: Context-aware workflow management. In: Baresi, L., Fraternali, P., Houben, G.-J. (eds.) ICWE 2007. LNCS, vol. 4607, pp. 47–52. Springer, Heidelberg (2007)
14. Han, J., Cho, Y., Kim, E., Choi, J.: A Ubiquitous Workflow Service Framework. In: Proceedings of the 2006 International Conference on Computational Science and its Application (2006)
15. van der Aalst, W.M.P., ter Hofstede, A.H.M., Kiepuszewski, B., Barros, A.P.: Workflow Patterns. *Distributed and Parallel Databases* 14(3), 5–51 (2003)

A History-Based Scheduler for Dynamic Load Balancing on Distributed VOD Server Environments

Jongbae Moon¹, Hyun-joo Moon², and Yongyun Cho^{3,*}

¹ Supercomputing Center, Korea Institute of Science and Technology Information,
52-11, Eoeun-dong, Yuseong-gu, Daejeon, 305-806, Korea
jymoon@kisti.re.kr

² Dept. of Cultural Contents, Hankuk University of Foreign Studies,
270 Imun 2-dong, Dongdaemun-gu, Seoul, 130-082, Korea
hyunjoomoon@gmail.com,

³ Information and Communication Engineering, Suncheon National University,
413 Jungangno, Suncheon, Jeonnam 540-742, Korea
sslabycho@hotmail.com

Abstract. As computer and network technology advance, multimedia data can be transferred in real time on the Internet. The increasing user demands for various multimedia data make VOD (Video-on-Demand) services to be developed. VOD services are being used in lots of fields, such as entertainment, distance learning, home shopping, and interactive news. In comparison to the existing HTTP services, these VOD services have special features that the time for their services is longer and the services request more disks and network bandwidths. Therefore, compared to HTTP service environments, VOD services have some different workload patterns. In these VOD service environments, the existing load balancing algorithms researched before are not proper. In this paper, we propose a new load balancing algorithm that is based on the history of past user access patterns to make server loads even on hierarchically distributed VOD system environments. This algorithm uses a dynamic genetic algorithm. The proposed distributed VOD system environment consists of a number of VOD servers, which are distributed geographically, and control servers that manage each group of VOD servers. User requests are distributed to prevent convergence by distributing VOD servers geographically. We use a genetic algorithm based on history data to distribute user requests in a local service area. The information of user requests and services is stored and referred in a database as history data. By applying these history data to the fitness function of genetic algorithms, we implemented the genetic algorithm and operations for VOD systems. The load balancing algorithm proposed in this paper can distributed workloads by predicting workloads precisely on VOD environments.

Keywords: VOD, distributed system, load balancing, genetic algorithm.

* Corresponding Author.

1 Introduction

Recently, the Internet is growing increasingly due to the development of computer networks. In addition, there were many attempts to transfer audio and video data through the Internet, and it can be possible to transfer movies via the World Wide Web and high-speed Internet. Video-on-Demand (VOD) service is activating due to these attempts and technologies. VOD allows users to select and watch/listen to video or audio contents on demand over high speed networks, while sitting at home, and have almost as good a control over the viewing of the video as when using a conventional VCR [1]. VOD services have some different characteristics from HTTP service as well access patterns. VOD services have long service time more than HTTP services, and request more network bandwidth and hardware disk. Moreover, VOD is more QoS-sensitive service than HTTP service, so users cancel their requests when the service is delayed. Therefore, the scheduling algorithms that used in distributed environments are not suitable for VOD system environments.

There are many researches on scheduling algorithm for distributed system [2][3][4]. We have to consider dynamic loads information about servers in dynamic service environments. Moreover, analysis of users request pattern must be considered to predict system loads and balance loads of VOD system. For examples, when we know the access pattern of users for a day, we are able to get information about the time of day that resources are most frequently used. CDN (Content Distribution Network) service uses such access pattern to store contents that are frequently used into cache memory. If we record histories of user requests and then analyze the history, the user request pattern is recognized by using some data mining techniques. Such users request pattern operates as an important part of load balancing for VOD system environments. In this paper, we suggest a scheduling algorithm for hierarchical VOD system environment to provide high QoS and load balancing. The proposed scheduling algorithm uses users request pattern, and selects appropriate servers to assign various requests by using genetic algorithm. We developed a fitness function based on history of user requests for the genetic algorithm.

2 System Model

VOD systems divided into three categories according to how to manage VOD servers: central, distributed, and hierarchical server model. The central server model consists of many VOD servers and a central server which controls users request to the VOD servers. This model is easy to set up and provides fast response time due to very simple scheduling algorithm. However, the model results in lower quality when users far from the central server. Moreover, the front-end server may cause a Single Point of Failure (SPOF) because the central server does not work when user requests increase rapidly.

Distributed server model distributes VOD servers managing user requests by themselves and has no central managing server. This model provides fast response time due to assigning user requests to the nearest VOD server. When the VOD server gets high loads, user requests are forwarded to another VOD server nearest. Therefore, there is no SPOF and provides scalability easily. However, VOD server has to keep track of other servers' loads information to forward user requests to another server

which has low load and located near. The number of distributed servers increase, the amount of information to keep track of as well as network overhead also increase.

Hierarchical distributed server model combines the advantages of central and distributed server model. This model distributes VOD servers geographically, and control servers that located in each site manage user requests and VOD servers. This model can reduce control servers' loads by assigning user requests to the nearest server group. Therefore, there are no SPOF and scalability is provided by adding VOD server in each site.

In this paper, we construct a hierarchical distributed VOD system model to develop history based scheduling algorithm. Fig. 1 shows the system model which consists of a global DNS, VOD servers, control servers, and database server. Each server group can be composed of many VOD servers and a control server which manage the servers in the group. And the control server performs load balancing by distributing user requests.

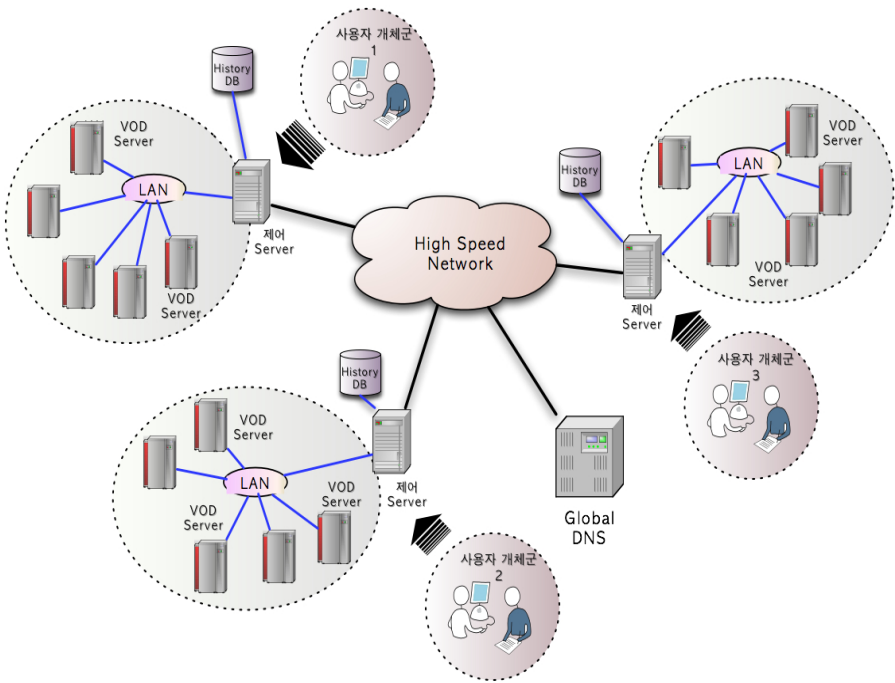


Fig. 1. Hierarchical Distributed VOD System Model

3 Proposed Scheduling Algorithm

In this section, we suggest a two level load balancing algorithm for hierarchical distributed VOD system. The first step is a global load distribution by using global DNS, and the second is local load distribution. Fig. 2 shows the global load balancing algorithm which keeps track of best VOD server list and is performed by the global DNS [5]. We assumed that there are N servers which are registered in the global.

```

while (1) {
  for(i=1 to n)
  {
     $T_{i, Load} = \alpha \times L_{i, CPU} + \beta \times L_{i, MEM} + \gamma \times L_{i, NT}$ 
    if (( $S_i$  is in DNS List) and ( $T_{i, Load} > TH$ ))
      then remove  $S_i$  from DNS list
    else if (( $S_i$  is not in DNS list) and ( $T_{i, Load} < TH$ ))
      then add  $S_i$  to DNS list
  }
}
 $T_{i, Load}$  = total loads of  $i_{th}$  control server
 $\alpha, \beta, \gamma$  : weighting value ( $\alpha + \beta + \gamma = 100$ ),  $TH$  : threshold value
 $S_i = i_{th}$  control server
 $L_{i, CPU}$  = CPU loads of  $i_{th}$  Server,  $L_{i, MEM}$  = Memory loads of  $i_{th}$  server
 $L_{i, NT}$  = Network loads of  $i_{th}$  Server

```

Fig. 2. Global load balancing algorithm

First of all, we need to calculate total load of i th control server S_i by estimating CPU, memory, and network utilization rate of each server. A monitoring module running on the control server collects each server's CPU, memory, and network loads. We assign a different weighting value according to the service provided by distributed VOD system. Some services that transfer large amounts of data such as online broadcasting or VOD service have more increase in network overhead than CPU and memory, while some services that request large amount of computing power such as online game have more increase in CPU and memory utilization than network. Therefore, we put more weight on the resources that have a big increase according to the type of service.

When a control server S_i is in the global DNS list and total load is larger than the threshold value TH , the S_i is removed from the list. When the S_i is not in the list and total load is less than the threshold value, S_i is inserted in the list again. Then the DNS manage the list with best servers. The threshold value is determined by using heuristics and is assigned by administrator manually. We provide the following four scheduling policies based on load information that is collected by monitoring module running on control servers, and can modify the policy dynamically.

- CPU based : $\alpha = 100, \beta = 0, \gamma = 0$, when a service is CPU intensive
- Memory based : $\alpha = 100, \beta = 0, \gamma = 0$, when a service is Memory intensive
- Network based : $\alpha = 0, \beta = 0, \gamma = 100$, when a service is Network intensive
- Combined based : $\alpha + \beta + \gamma = 100$, when a service uses all of the resources, the weight is assigned by administrator

Local load distribution is a method to distribute load between VOD servers managed by a control server. We adopt a genetic algorithm to local distribution because there is no general solution to find best servers under the various types of service and lots of requests. When the requests come into the queue, we need to sort the requests as the service type. After sorting, we generate a population of candidate solutions. We create candidate solutions from 10% of all $m \times n$ cases randomly. A fitness function is a particular type of objective function that prescribes the optimality of a solution in a genetic algorithm so that that particular chromosome may be ranked against all the other chromosomes. Optimal chromosomes, or at least chromosomes which are more optimal, are allowed to breed and mix their datasets by any of several techniques. And the following is a fitness function and related expressions.

$$F(x) = \sum_{i=1}^m T_i$$

T_i = current loads on $P_i \times$ working time on $P_i \times$ Load History Value

(P_i : i th process, T_i : expected working time)

(Loads on P_i = CPU Utilization \times Memory Utilization \times Network Utilization on P_i)

(Load History Value = Average load of jobs that executed before)

$$\text{Fitness value} = \frac{1}{F(x)}$$

4 Performance Evaluation

We implemented a simulator which is based on OPNET Modeler [6] and conducted some comparison experiments with GA (genetic algorithm), RR (round-robin), and Random scheduling algorithm when users increase. The first experiment is about changes of average response time when users increase from 1 to 3,000, and Fig. 3 shows the result. As you can see from Fig. 3, the average response time is pretty much the same to some extent which is about 800. However, there is a growing gap among scheduling algorithms. Random scheduling shows the largest increase because highly loaded servers are getting user requests continuously. In this experiment, the average response time of GA, the proposed scheduling algorithm, is under 40 seconds until users are 3,000.

Fig. 4 shows the cancellation rate when users increase. There is a little probability that users cancel popular VODs and there is strong probability that users cancel non-popular VODs. The cancellation rate also increases when the response time is getting

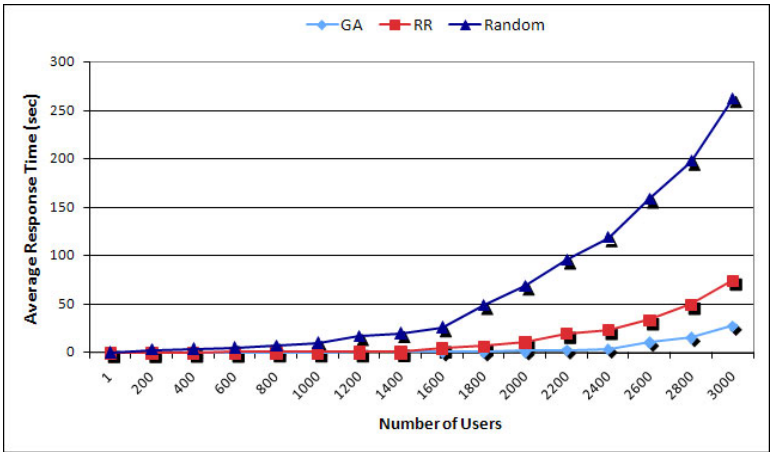


Fig. 3. Average Response Time when users increase

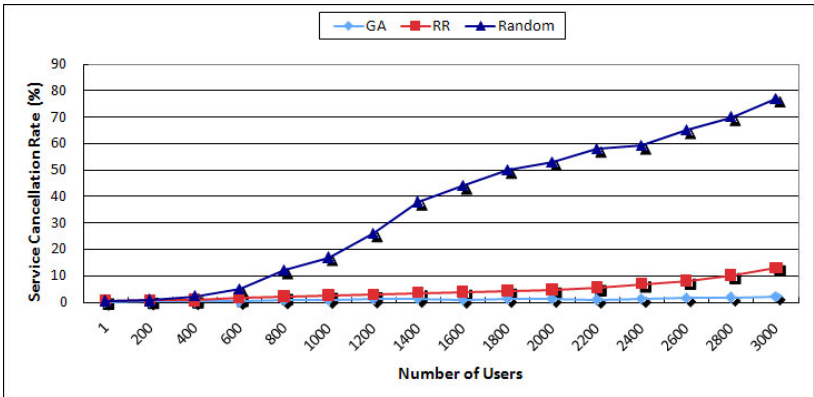


Fig. 4. Service Cancellation Rate when users increase

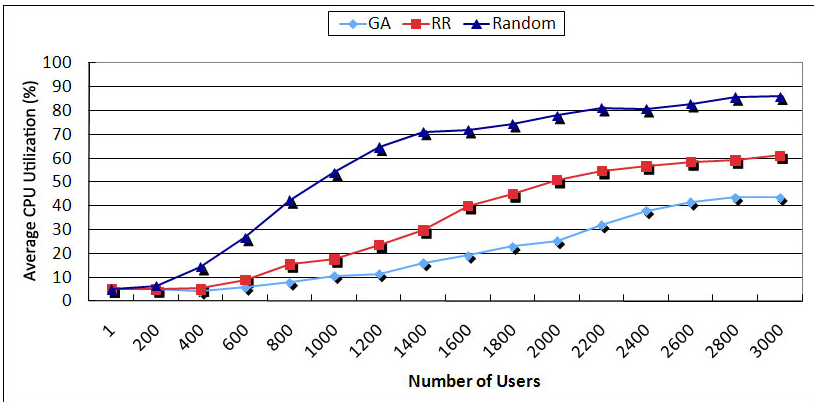


Fig. 5. Average CPU Utilization when users increase

too long. Compared to response time, it looks like Random algorithm has a high cancellation rate. This means that there are many users who cancel the VOD due to long response time because the response time had a lot of ups and downs. There is no big change of the average response times of RR and GA. However, the gap is growing over 2,600 when using RR. In this experiment, Random algorithm shows 77%, RR shows 13%, and GA shows 2.2% of cancellation rate. The result shows that there is a little probability that users cancel the service using GA excepting when the user cancel the service due to lack of interest.

The next experiment is about average CPU utilization when the users increase. Fig. 5 shows the result that the CPU utilization rate is growing when users increase. The CPU utilization rate increase sharply when we use Random scheduling because the requests may be assigned into highly loaded server. In response to this we can see the response time increase in Fig. 3. The average CPU utilization of RR and GA scheduling increases slowly and steadily. GA scheduling shows the very gently increase because lots of requests are evenly distributed and the gap between servers' loads is getting close.

We also conducted an experiment using user request pattern and history during a day. The users' pattern can get from the service history. In this experiment, we assumed the user requests increase dramatically at 6 P.M. and the numbers of users increase from 1 to 1,000.

Fig. 6 shows the result of the changes of response time as the user request pattern. The x-axis refers to the time from 10 A.M. to 10 P.M. The response time is getting longer when the rush hours of about 6 P.M. Random and RR scheduling shows the response time is getting longer sharply. Random scheduling has a large gap between response times when the system has not many requests, and moreover the response time is longer when the requests increase rapidly. However, GA scheduling predicts the loads of servers and distributes the loads evenly among servers, so the response time increases slowly. Therefore, we can see that the proposed scheduling algorithm provides more enhanced QoS.

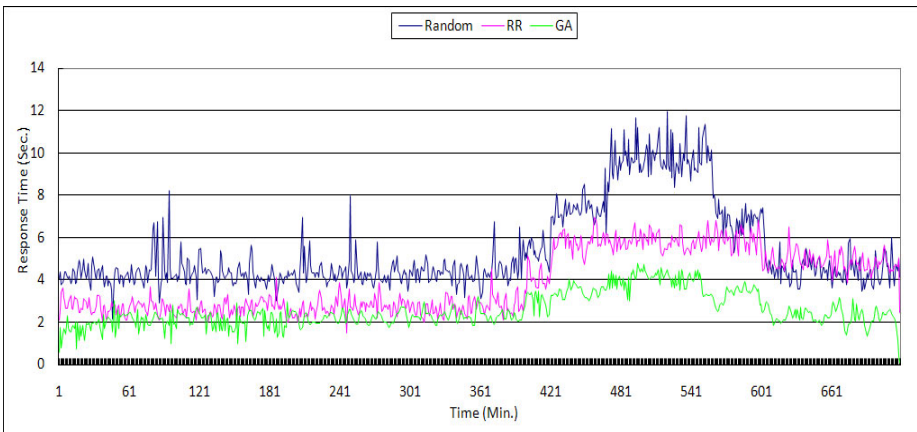


Fig. 6. Response time as time passed

5 Conclusions and Future Works

In this paper, we suggested a dynamic load balancing algorithm in hierarchical distributed VOD system environment. The hierarchical distributed VOD system model consists of a global DNS, VOD servers, control servers, and database server. Each server group can be composed of many VOD servers and a control server which manage the servers in the group. And the control server performs load balancing by distributing user requests. The dynamic load balancing algorithm is based on genetic algorithm and history information of user request pattern. We implemented a simulator based on OPNET. And we conducted some comparison experiments with GA, RR and Random scheduling method. According to the results, we can see that the proposed algorithm provides more enhanced QoS.

References

1. Ma, K., Shin, K.G.: Multicast Video-On-Demand Services. *ACM SIGCOMM Computer Communication Review* 32(1), 31–43 (2002)
2. Loukopoulos, T., Ahmad, I.: Static and Adaptive Distributed Data Replication Using Genetic Algorithms. *J. of Parallel and Distributed Computing* 64(11), 1270–1285 (2004)
3. Rothlauf, F.: Representations for Genetic and Evolutionary Algorithms. In: *Studies on Fuzziness and Soft Computing*, vol. 104 (2002)
4. Mundur, P., Simon, R., Sood, A.K.: End-to-End Analysis of Distributed Video-on-Demand Systems. *IEEE Transactions on Multimedia* 6(1), 129–141 (2004)
5. Moon, J.-B., Kim, M.-H.: Dynamic Load Balancing Method Based on DNS for Distributed Web Systems. In: Bauknecht, K., Pröll, B., Werthner, H. (eds.) *EC-Web 2005*. LNCS, vol. 3590, pp. 238–247. Springer, Heidelberg (2005)
6. OPNET Technologies Inc., <http://www.opnet.com/>

A Secure Routing Protocol for Wireless Sensor Networks

Jaydip Sen and Arijit Ukil

Innovation Lab, Tata Consultancy Services Ltd., Bengal Intelligent Park,
Salt Lake Electronics Complex, Kolkata, India
{Jaydip.Sen,Arijit.Ukil}@tcs.com

Abstract. Wireless sensor networking has been a subject of extensive research efforts in the recent years, and has been well recognized as a ubiquitous and general approach for some emerging applications such as a real-time traffic monitoring, ecosystem and battlefield surveillance. Since these networks deal with sensitive data, it is imperative that they are made secure against various types of attacks such as node capture, physical tampering, eavesdropping, denial of service etc. This paper presents a secure routing mechanism for wireless sensor networks. The protocol is resilient in the presence of malicious nodes that may launch selective packet dropping attack on the routing path. The scheme employs single-path routing, and therefore, it is energy-efficient. While the packets are forwarded towards the base station, if any node fails to forward a packet, it is isolated immediately. Packets are then routed around the node. Simulation conducted on the scheme clearly demonstrates that it is more efficient than some of the existing similar schemes.

Keywords: Wireless sensor networks, routing, malicious packet dropping, neighborhood watch, energy-efficiency, reliability.

1 Introduction

Recent convergence of technological and application trends has resulted in an exceptional level of interest in wireless ad hoc networks, and in particular, wireless sensor networks (WSNs). The push was provided by rapid progress in computation and communication technology as well as the emerging field of low cost, reliable, MEMS-based sensors, while the pull was provided by numerous applications that can be summarized under the umbrella of computational worlds, where the physical world can be observed and influenced through the Internet and WSN infrastructures. Typically, WSNs contain hundreds or thousands of sensor nodes that have the ability to communicate among them and also to a base station (BS) [1]. While the sensor nodes have limited sensing region, processing power and energy, networking a large number of such nodes makes a robust, reliable and accurate sensor network covering a wider region. The WSNs are envisioned to play an important role in a wide variety of applications ranging from critical military surveillance to forest fire monitoring and the building security monitoring. Since these networks are mostly deployed in harsh and often hostile operating environments, security and reliability are two important aspects for their operations. Design of a robust and secure routing protocol for WSNs is

a particularly challenging problem due to limited computational resources of the nodes and unreliable wireless links. In addition, data transmission through these networks is prone to many types of failures and attacks such as node capture, physical tampering, eavesdropping, denial of service, etc [2, 3, 4]. Some of the previous research efforts in secure routing for WSNs have focused on the design of single-path routing using low-rate and periodic flooding of events to avoid failed nodes [5]. A fault-tolerant cluster-based routing scheme has been proposed in [6]. In [24], a multi-path routing protocol has been presented that uses three paths from the source node to the sink to provide robustness.

To address the problem of security and efficiency in routing in WSNs, this paper proposes a scheme that reliably identifies compromised (or faulty) nodes and utilizes a routing path that avoids these nodes. Essentially, it utilizes a single-path routing concept and thereby saves energy-consumption. If a malicious node is detected on the next-hop on the routing path, the node is efficiently bypassed and the packets are routed around the node to the base station still in a single-path. The proposed protocol is a modification of the routing scheme proposed in [7]. However, it is more energy-efficient and less delay-inducing as observed from the simulation results.

The rest of the paper is organized as follows. Section 2 discusses some related work in the area of secure routing in WSNs. Section 3 presents the WSN model and possible types of packet dropping attacks on these networks. It also gives a brief introduction to the LEAP protocol for key management used in the proposed scheme. Section 4 gives a detailed description of the proposed protocol. Section 5 presents the simulation results and Section 6 concludes the paper.

2 Related Work

A considerable amount of research effort has been put in design of secure routing protocols for WSNs. While some of these works have focused on defense against outsider attacks on the key management protocols being used in the network [8, 9, 10], the others have attempted to secure node-to-node communication [10, 11]. However, all these schemes fail even if a single node is compromised.

An insider attack on a WSN is even more dangerous, since an attacker in this case can launch various types of attacks such as dropping of legitimate packets, injecting of bogus packets and sensing reports, advertising false routing messages, eavesdropping on the network communication etc. Parno et al have proposed a distributed detection mechanism [12]. Newsome et al [13] have presented some techniques for preventing an adversary from launching Sybil attack by arbitrarily creating new identities.

Deng et al have proposed an *intrusion tolerant routing protocol* in wireless sensor networks (INSENS) that adopts a routing-based approach to security in WSNs [14]. It constructs routing tables in each node, bypassing malicious nodes in the network. The protocol cannot totally prevent attack on nodes, but it minimizes the damage caused to the network due to an attack. The scheme has reduced computation, communication, storage, and bandwidth requirements at the sensor nodes at the cost of greater computation and communication overhead at the base station. Tanachaiwiwat et al have proposed a novel secure routing protocol- *trust routing for location-aware*

sensor networks (TRANS) [15]. It makes use of a loose-time synchronization asymmetric key cryptographic scheme to ensure message confidentiality. Zhu et al have proposed a *localized encryption and authentication protocol* (LEAP) to prevent insider attacks--particularly message eavesdropping and message fabrication attacks [16]. Compromised or faulty nodes may also drop legitimate packets by launching selective forwarding attacks [2]. In order to defend against such type of attacks, most of the existing secure routing protocols in sensor networks are based on multi-path forwarding scheme [17, 18], or interleaved mesh forwarding scheme [19].

In contrast to multi-path approach, this paper presents a single-path packet-forwarding scheme in sensor networks where some of the nodes may be malicious and drop any packets that they receive. The protocol is an extension of the protocol presented in [7] and is more energy-efficient than the base protocol. The main contribution of the paper is its energy-efficiency and increased reliability, since the packets in the network are always routed in single-path avoiding costly multi-path approach. In the next section, the sensor network model and its associated security vulnerabilities are discussed.

3 Network Threats and Key Management

The proposed protocol is suitable for a large and dense sensor network where one or multiple base stations and a large number of tiny, inexpensive, static and resource-constrained sensor nodes are deployed over a wide area. Data gathered by the sensor nodes are sent along multi-hop routes to the base station. It is assumed that each sensor node has the same transmission range and all links are bi-directional.

It is also assumed that a key establishment scheme like LEAP [16] is already deployed in the nodes. A short description of LEAP is given in Section 3.2.

3.1 Threat Model

The goal of the proposed scheme is to defend against packet-dropping attacks that may be launched by compromised nodes in a WSN. As mentioned in Section 2, the insider attacks may have devastating effects on the performance of a WSN and suitable defense mechanism should be deployed to counter such attacks.

A node in a WSN may drop packets because of two reasons: (i) it may be faulty due to some problems in its hardware and/or software, or (ii) it may be malicious and intentionally may not forward any packets. To ensure reliable packet delivery, some authors [20, 21] have proposed use of acknowledgements (ACK). However, in presence of malicious nodes, an ACK received from the next-hop node cannot guarantee that the packet is forwarded by the next-hop node; it only implies that the packet has been received by the next-hop node. It will not be possible to ensure that the next-hop node really forwards the packet just by implementing an ACK mechanism.

Due to the broadcast nature of the wireless links in the sensor network, it is possible for a node to overhear communications to and from its neighbor nodes. Therefore, a node u after forwarding a packet to its next-hop v can listen on v 's communication to check whether v really forwards the packet further. The simple monitoring mechanism, however, cannot also guarantee reliable packet delivery either. A malicious

node may forward the packet but it can make its intended next-hop id such that there is no node in the neighborhood which has the matching id. Therefore, even if the packet is forwarded by the malicious node it is ultimately lost in the network. This is called *routing disruption attack*, and it is impossible to detect such attack by a naïve neighborhood monitoring scheme.

The objective of the proposed scheme is to address all types of packet-dropping attacks ranging from the most naïve type (e.g., detection of a faulty node) to the most malicious type (e.g., routing disruption attack).

3.2 Key Establishment Scheme in LEAP

The *localized encryption and authentication protocol* (LEAP) proposed by Zhu et al [16] is a key management protocol for WSNs based on symmetric key algorithms. It uses different keying mechanisms for different types of packets depending on their security requirements. Four types of keys are established for each node: (i) an *individual key* shared with the base station (pre-distributed), (ii) a *group key* shared by all the nodes in the network (pre-distributed), (iii) *pair-wise keys* shared with its neighbor node, and (iv) a *cluster key* shared with multiple neighbor nodes. The pair-wise keys shared with immediate neighbor nodes are used to protect peer-to-peer communication and the cluster key is used for local broadcast.

It is assumed that the time required for an attacker to attack a node, T_{min} , is greater than the time interval during which a sensor node can detect all its intermediate neighbors, T_{est} . This is a reasonable and practical assumption which has also been made in some of the other research contributions [22].

A newly joined node u in the network executes the following four steps to establish a pair-wise key with each of its neighbors. The steps are as follows:

1. *Key Pre-distribution*: each node u is loaded with a common initial key K_I and derives its master key $K_u = f_{K_I}(u)$, where f_K is a pseudo-random function. The master key depends on the common key and the unique identifier of the node.
2. *Neighbor discovery*: node u sets up a timer to fire after time T_{min} , broadcasts its id, and waits for each neighbor v 's ACK. The ACK from v is authenticated using the master key K_v , of node v . Since node u knows K_I , it can derive $K_v = f_{K_I}(v)$.

$$\begin{aligned} u &\rightarrow * : u, R_u \\ v &\rightarrow u : v, MAC(K_v, R_u | v) \end{aligned}$$

3. *Pair-wise key establishment*: node u computes its pair-wise key (K_{uv}) with v , as $K_{uv} = F_{K_v}(u)$. Node v also computes K_{uv} in the same way. K_{uv} serves as their pair-wise key between the nodes u and v .
4. *Key revocation*: when its timer expires, node u erases K_I and all the master keys of its neighbors. Every node, however, keeps its own master key, in order to establish pair-wise keys with nodes which have joined later.

Once K_I is erased, a node will not be able to establish a pair-wise key with any other nodes that have also erased K_I . Without the knowledge of K_I , an attacker who has been able to compromise a node after T_{min} , fails to establish pair-wise keys with any nodes except the neighbors of the compromised node. In such a way, the protocol is able to localize any possible attack on the network.

The cluster key is established by a node after the establishment of the pair-wise key establishment. A node generates a cluster key and sends it to each of its neighbors after encrypting it using the pair-wise shared key. The group key is pre-loaded but is updated once any compromised node is detected in the network. This could be done either by base station's sending the new group key to each node using its individual key, or a hop-by-hop basis using cluster keys.

4 The Proposed Routing Protocol

This section presents the proposed secure routing protocol for WSNs. The protocol is based on a robust *neighborhood monitoring system* (NMS) and is an extension of the scheme proposed by Lee and Choi [7]. NMS works on promiscuous monitoring of the neighborhood by a node and detection of any possible malicious packet dropping attack by a cooperative algorithm using *neighbor list checking* to be described in Section 4.1.

The proposed scheme ensures reliable hop-by-hop delivery of packets in a WSN even in presence of malicious nodes that may launch packet-dropping attack in the routing path. To defend against packet-dropping attack, most of the proposed algorithms in the literature have exploited the concept of multi-path routing, where a single packet is routed through multiple paths from the source to the sink. While this approach ensures reliable packet delivery, it consumes a high amount of energy for delivery of each packet. To avoid this problem, the proposed algorithm uses a single-path routing mechanism. If a malicious node is encountered, the node is avoided and the packet is routed around it in an efficient manner (discussed later in this section) still in a single-path mode to the base station. The selection of the new path is based on some broadcast signaling in the neighborhood of the malicious node.

The proposed algorithm is more efficient than the algorithm proposed in [7], since it always uses single-path routing in contrast to the occasional multi-path approach used in [7]. The use of single-path routing makes it more energy-efficient and less delay-inducing. The algorithm is described below:

1. *Neighbor list checking*: during the neighbor discovery phase, each node exchanges Hello messages with its neighbor nodes to know its 1-hop and 2-hop neighbors (i.e., neighbors of each of its neighboring nodes). The neighborhood information is subsequently verified by exchange of *neighbor list checking* messages as described in Section 4.1.
2. *One-hop packet forwarding*: when a node u sends a packet to its neighbor, it first keeps a copy of the packet in its buffer, and then forwards it to its next-hop node v before encrypting it with the cluster key of the node u . Since the cluster key is shared between the node and all its neighbors, the packet encrypted and sent by node u to node v can be overheard by all the neighbors of node u .
3. *Monitoring nodes selection*: as the packet is being forwarded from node u to node v , the neighbors of node u that are also neighbors of node v receive the packet and store it in their buffers. These nodes are designated as the secondary monitoring nodes. For example, in Fig. 1, nodes w and y are the secondary monitoring nodes for node v . The node u is the primary monitoring node in this case. The nodes that are not neighbors of node v but have received the packet

because they are neighbors of node u , discard the packet. The primary node knows the secondary monitoring nodes, since every node has the knowledge of its 1-hop and 2-hop neighbors.

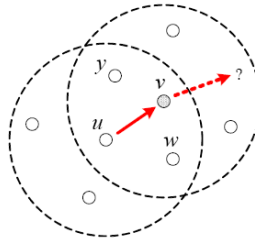


Fig. 1. Neighbor monitor system (secondary nodes w , y and primary node u)

4. *Role of secondary monitoring nodes:* secondary monitoring nodes w and y monitor the traffic from node v and compare the outbound packets from node v with the packets stored in their buffer. The next-hop address of each packet is also verified to check whether the packet's intended next-hop is a really a neighbor of node v , by cross-checking the neighbor list of node v . If both these checks yield positive results, the secondary monitoring nodes remove the packet from their buffer and their role of monitoring is complete for that packet. If any packet is found to remain in the buffer of a secondary monitoring node for more than a threshold period of time, it first sends a broadcast signal in its neighborhood to inform all its neighbors that it is going to forward the packet to its designated next-hop so that other neighbors do not forward the same packet. The secondary monitoring node now forwards the packet to its designated next-hop after encrypting the packet with the cluster key. The role of the secondary node now becomes that of the primary node and its neighbors become the secondary node. This is in contrast to the scheme proposed in [7] in which all the secondary nodes forward the packet in a multi-path mode. The use of a single-path instead of multi-path makes the proposed scheme more energy-efficient and less delay-inducing compared to the one proposed in [7].
5. *Role of primary monitoring node:* The role of a primary monitoring node (node u) is identical to that of secondary monitoring nodes (nodes w and y); the only difference is that it listens not only on the traffic from node v , but also on the traffic from the nodes w and y . If the packet is correctly forwarded by any one of the nodes v , w , y , the node u removes the packet from its buffer. The role of node u as the primary monitoring node is now complete. If time out occurs for a packet, the primary monitoring node u forwards the packet (encrypted with its cluster key) to its next-hop other than node v .

As the packet is routed along a path towards the sink, the above steps of NMS algorithm except the *neighbor list checking* are executed at each hop so that reliable packet delivery can happen through a single path. This is in contrast to the previous schemes proposed in [19, 21, 22]. In these schemes, a node broadcasts a packet without specifying a designated next-hop, and all neighboring nodes with smaller costs

(the cost at a node is the minimum energy required to forward a packet from the node to the base station) or within a specific geographic region continue forwarding the packet to the base station. If nodes v , w , and y have smaller costs than node u in Fig. 1, then each of them will forward packets received from node u following the existing approaches. However, in the proposed scheme, nodes w and y only observe the packet forwarding activities of node v , instead of actively forwarding the packets. In the event of no packet drop, the routing to the base station happens in a single-path thereby making the process highly energy-efficient.

Even in the event of a packet drop, the proposed algorithm works in a single-path mode. This makes it more efficient than the one proposed in [7]. If the node v in Fig. 1 does not forward the packet it has received from node u , then one of the secondary monitoring nodes w and y would forward the packet to its next-hop nodes. The node (either w or y) that forwards the packet to its next-hop neighbors will first send a broadcast message in its neighborhood so that its other neighbors would not forward the same packet. This single-path routing approach even in the event of packet dropping makes the proposed scheme extremely energy-efficient and very less delay-inducing. Fig. 2 shows an example of the application of the scheme, where two malicious (or faulty) nodes are bypassed as the packet is routed to the base station in a single-path.

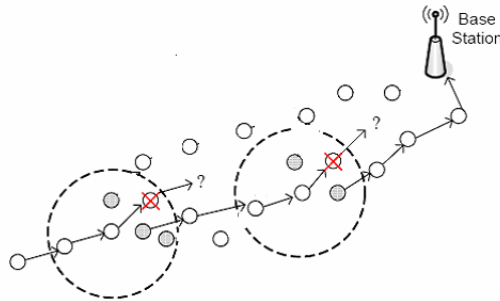


Fig. 2. Two malicious nodes identified by secondary monitoring nodes

For the scheme to work, each packet should be encrypted with a cluster key of the forwarding node so that all the neighbors of the forwarding node can decrypt and overhear it. If a link-level encryption was applied between each pair of nodes in the routing path, the scheme would have been more robust, since a compromised node could decrypt only the packets which were destined to it. However, it would have made the scheme less resilient to packet dropping attack. Since encryption with a cluster key provides a reasonable level of robustness to a node compromise, [11] and also supports local broadcast (i.e. resiliency against packet-dropping) it makes the algorithm optimum in its performance [11].

To make the scheme robust to routing disruption attack, where a node intentionally forwards the packets to a spurious address of the next-hop so that the packet is lost in routing, it is necessary that each node should prove that it really has the claimed neighbors. Apparently, a node has the knowledge of its direct neighbors by neighbor discovery and pair-wise key establishment phases discussed earlier. However, in the

case of two-hop neighbors, a malicious node v can inform its neighbor u that it also has neighbor node x (any possible id in the network) which in fact is not a neighbor of node v (Fig. 1). Apparently, there is no way node u can detect these false claim of v since x is not in the neighborhood of u . Therefore, it is absolutely necessary to devise a scheme by which a node can verify the neighbors of each of its neighboring nodes. This scheme is explained in the next subsection.

4.1 Neighbor List Checking

To verify the 2-hop neighbors, i.e., the neighbors of the neighbors of each node, the proposed scheme utilizes a concept of *neighbor list checking* (NLC). NLC uses the pair-wise keys established in LEAP. During neighbor discovery, NLC uses *three-way handshaking* in order to identify not only the communicating nodes but also their respective neighbors. There are two distinct types of neighbor discovery in NLC. In one case, both the nodes involved in the neighbor discovery process are still within the initial T_{min} – both the nodes are referred to as new nodes. In the other case, the neighbor discovery happens between a newly-deployed node within the initial T_{min} and an existing node over the initial T_{min} – the latter node is referred to as an old node.

(i) Neighbor discovery between new nodes: If a new node broadcasts its neighbor list before the expiry of its initial T_{min} , it is verifiable. Therefore, it is essential that both nodes broadcast their respective neighbor lists before their respective T_{min} . The three-way handshaking neighbor discovery between pure node u and v may be represented as follows [7]:

$$\begin{aligned}
 & u \rightarrow * : u, R_u \\
 v \rightarrow u : & v, T_v, R_v, \text{MAC}(K_v, R_u \mid K_u \mid (u, T_v, R_v)) \\
 u \rightarrow v : & u, T_u, \text{MAC}(K_{uv}, R_v \mid (u, T_u))
 \end{aligned}$$

T_v and T_u are the time remaining to reach T_{min} of v and T_{min} of u , respectively. R_u is a random number generated by node u . $\text{MAC}(K, M_1 \mid M_2)$ denotes message authentication code of concatenated message of M_1 and M_2 using key K . Node u broadcasts its id, and waits for each neighbor v 's ACK. The ACK from every neighbor v is authenticated using the master key K_v of node v . Since node u knows K_I , it can derive $K_v = f_{KI}(v)$. The ACK from node v contains T_v . If T_v is non-zero, node v claims to be a new node. K_u in MAC proves node v to be a new node, since new node v should know K_I and derive $K_u = f_{KI}(u)$. Node u records T^*_v (T_v added to the current time of node u) in the entry for node v in the *neighbor information table* (NIT). Node u computes its pair-wise key with v , $K_{uv} = f_{Kv}(u)$. Node u also generates $\text{MAC}(K_v, v \mid u)$ (which means that v certifies u as an immediate neighbor), and stores it as a *certificate*.

The ACK from node u also contains T_u . This ACK is authenticated using their pair-wise key K_{uv} , which proves node u a new node and u 's identity. Node v then records T^*_u (T_u added to the current time of v) in the entry for u in its NIT. It also generates $\text{MAC}(K_{uv}, u \mid v)$ and stores it as a certificate. This completes the three-way handshaking process.

Every new node u broadcasts its neighbor list just prior to T_{min} of u . Each receiving neighbor v checks whether the receiving time at v is prior to T_u^\wedge in the NIT. If yes, the neighbor list of u is now certified by each neighbor v .

(ii) Discovery between a new and an old node: After the bootstrapping phase, new nodes can join the network. This situation is depicted in Fig. 3, where the new node x has joined the network. The nodes u, v and w are old nodes. In this scenario, the messages exchanged during the three-way handshaking are different. The message exchange between the new node x and the old node u is as follows [7]:

$$\begin{aligned}
 & x \rightarrow * : x, R_x \\
 u \rightarrow x : & u, T_u, R_u, v, \text{MAC}(K_v, v | u), w, \text{MAC}(K_w, w | u), \text{MAC}(K_u, R_x | (u, T_u, v, \\
 & \text{MAC}(K_v, v | u), w, \text{MAC}(K_w, w | u))) \\
 x \rightarrow u : & x, T_x, \text{MAC}(K_x, x | u), v, \text{MAC}(K_v, x | u), w, \text{MAC}(K_w, x | u), \text{MAC}(K_{xu}, \\
 & R_u | (x, T_x, \text{MAC}(K_x, x | u), v, \text{MAC}(K_v, x | u), w, \text{MAC}(K_w, x | u)))
 \end{aligned}$$

The node x broadcasts its id, and waits for the ACKs from its neighbors. The ACK from neighbor u is authenticated using the master key K_u of u . Since x knows K_I , it computes $K_u = f_{KI}(u)$. The ACK from u contains T_u , the time remaining to reach T_{min} of u . As u is an old node in Fig. 3, T_u is zero. Node u sends its neighbor list (v and w) to x including their respective certificates in the ACK. Node x can verify the list since it knows K_I . If the verification is positive, x computes pair-wise key $K_{xu} = f_{Ku}(x)$, and also generates $\text{MAC}(K_u, u | x)$ and stores it as a certificate.

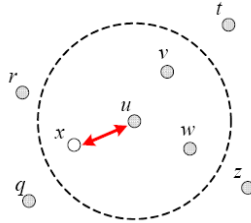


Fig. 3. Neighbor discover between a new node x and an old node u

The ACK from x also contains T_x , the time remaining to reach T_{min} of x . This ACK is authenticated using K_{xu} , to prove x 's identity as a new node. Node u computes T_x^\wedge (T_x added to the current time of u) and makes an entry for x in the neighbor table. Node x also sends the certificate for u , $\text{MAC}(K_x, x | u)$ in the ACK. It also provides one-time certificates for each of u 's neighbors (v and w) and completes the three-way handshaking process.

Node u broadcasts the one-time certificates received from x to its neighbors w and y so that they are able to verify the authenticity of x . The contents of the one-time certificates are as follows [7]:

$$u \rightarrow * : u, x, v, \text{MAC}(K_v, x | u), w, \text{MAC}(K_w, x | u), K_{xu}^A, \text{MAC}(K_{xu}^c, (u, x, v, \text{MAC}(K_v, x | u), w, \text{MAC}(K_w, x | u), K_{xu}^A))$$

K_u^A is a local broadcast authentication key in u 's one-way key chain, K_u^c is the cluster key of u . Each neighbor v of u verifies x by checking its one-time certificate, MAC ($K_v, x | u$) and x is certified by each neighbor v of u .

Before the expiry of T_{min} of x , it broadcasts its neighbor list. Each neighbor u of x checks whether the receiving time at u is before T_x^A in the neighbor table. If yes, the neighbor list of x is certified by u .

Through the above neighbor list checking process, every node knows the certified neighbors of each of its neighbors. With this knowledge, the monitoring nodes check the validity of the claims of the next-hop of the target forwarding node.

4.2 Neighbor Information Table

The information obtained through neighbor list checking is stored in the *neighbor information table* (NIT) of each node. Table 1 shows the NIT of node u as depicted in Fig. 3.

Table 1. Neighbor information table of node u

Neighbor ID	Certificate	Neighbor list
v	$MAC(K_v, v u)$	u, w, t
w	$MAC(K_w, w u)$	u, v, z
x	$MAC(K_x, x u)$	u, r, q

The NIT is used by the monitoring nodes. For example, if node u overhears a packet sent from w to v , it first checks its NIT and finds that both w and v are its neighbors. Node u starts listening on v 's traffic as a secondary monitoring node, and if v does not forward the packet to one of its verified neighbors, u will first broadcast a message in its neighborhood and then forward the packet to its next-hop other than v . Node u then starts acting as a primary monitoring node.

5 Simulation Results

To evaluate the performance of the proposed protocol, simulations have been carried out on network simulator *ns-2*. To compare the performance of the scheme with that proposed in [7] identical parameters are chosen. The simulation parameters are presented in Table 2.

The malicious nodes are configured in such a way that they do not forward any packets. These nodes are strategically distributed in a square area of 200 m side (from 150m to 350m of each side of the 500 m \times 500 m target area), so that they are half-way between the source and the base station. TinyOs beaconing [23] has been used as the base routing protocol in the simulation to maintain uniformity with [7]. Each simulation is repeated with 100 different topologies, and the average value is taken as the final result.

Fig. 4 and Fig. 5 show the *packet delivery ratio* (the ratio of the number of packets successfully reaching the base station to the total number packets sent by the source) in presence of varying number of malicious nodes (x) with $N = 300$ and $N = 600$

Table 2. Simulation parameters

Parameter	Values
No. of sensor nodes (N)	300, 600
Area of simulation	500 m * 500 m
Transmission range of each sensor	30 m
Degree of each node	10 ($N=300$), 20 ($N=600$)
Position of source	(50, 50)
Position of BS	(450, 450)
Avg. no of hops from source to BS	18

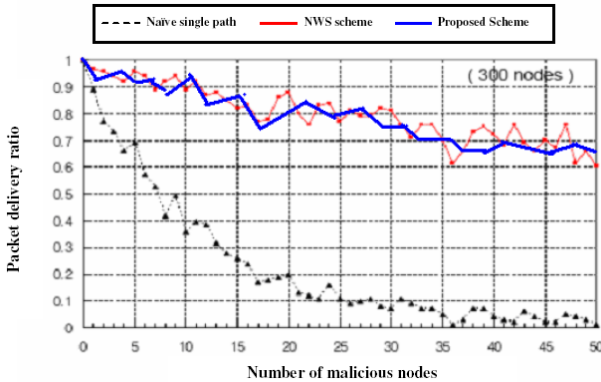


Fig. 4. Packet delivery ratio vs. no. of malicious nodes ($N = 300, x = 0 \sim 50$)

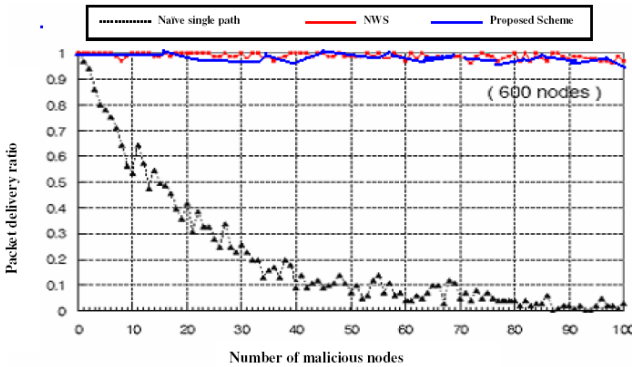


Fig. 5. Packet delivery ratio vs. no. of malicious nodes ($N = 600, x = 0 \sim 100$)

respectively. It is observed that with $N = 300$, the packet delivery ratio falls with the increase in number of malicious nodes. The performance of the proposed protocol is found to be almost identical to that of NWS proposed in [7]. With both these protocols, the packet delivery ratio is on the average more than 80%. Since in NWS some packets are routed in multi-path, it has a marginally higher of packet delivery ratio in

Fig. 4. As can be observed from Fig. 5, with $N = 600$, the proposed protocol has same performance as that of NWS. Even when $x = 100$, the delivery ratio is as high as 99%. As has been pointed out in [7], this is because of the higher degree of sensor nodes ($d = 20$), more number of malicious nodes can be bypassed.

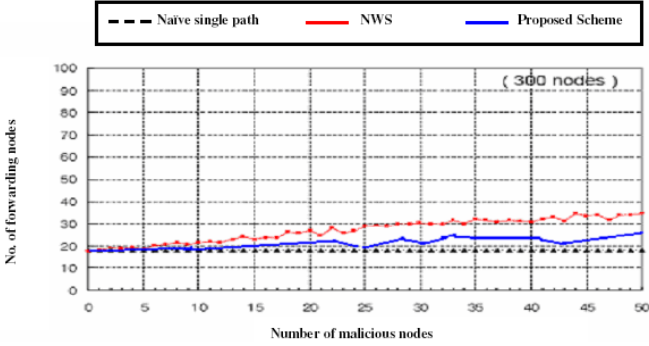


Fig. 6. No. of forwarding nodes vs. no of malicious nodes ($N = 300, x = 0 \sim 50$)

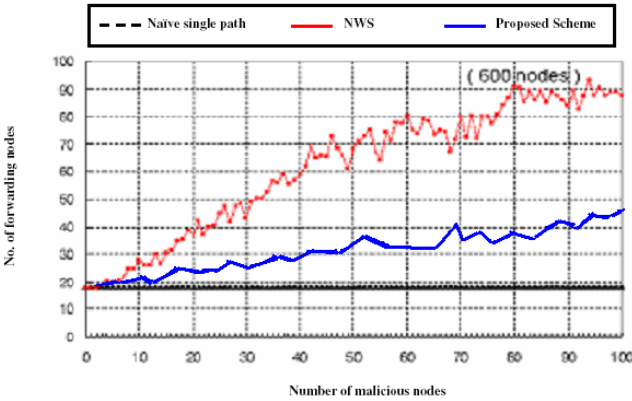


Fig. 7. No. of forwarding nodes vs. no of malicious nodes ($N = 600, x = 0 \sim 100$)

Fig. 6 and Fig. 7 show the number of nodes that has participated in packet forwarding for $N = 300$ and $N = 600$ respectively. Since the source is 18 hops away from the base station, the number of relaying nodes in single-path routing case remains constant at 18. The number of forwarding nodes was found to be increasing with x , since more number of nodes needs to be bypassed. However, in the proposed protocol, the number of forwarding nodes was found to be much smaller than in the NWS. In NWS, with $N = 600$, the number of forwarding nodes increases very fast with the increase in x , since more number of secondary monitoring nodes participate in routing. While this ensures reliability, it has an adverse effect on the energy-efficiency of the network. As the same packet is being routed in multi-path, NWS protocol has large energy consumption on the average. The proposed protocol while maintaining

the same level of packet delivery ratio (i.e. reliability), involves less number of forwarding nodes. It, therefore, consumes much less energy than NWS and thereby enhances the network life-time.

6 Conclusion

This paper has presented a reliable and energy-efficient single-path routing protocol for WSNs. The proposed protocol works reliably and efficiently even in presence of malicious nodes on the routing path that drop packets without forwarding them. The scheme is designed for hop-by-hop reliable forwarding of packets from the source to the base station while routing every packets in single-path. The simulation results have shown that the protocol is more energy-efficient than some of the existing routing protocols for WSNs while maintaining the same level of reliability.

References

1. Akyildiz, I.F., Su, W., Sankarasubramaniam, Y., Cayirci, E.: A Survey on Sensor Networks. *IEEE Communications Magazine* 40(8), 102–114 (2002)
2. Karlof, C., Wagner, D.: Secure Routing in Wireless Sensor Networks: Attacks and Countermeasures. In: *Proc. of 1st IEEE International Workshop on Sensor Network Protocols and Applications*, pp. 113–127 (2005)
3. Perrig, A., Stankovic, J., Wagner, D.: Security in Wireless Sensor Networks. *Communications of the ACM* 47(6), 53–57 (2004)
4. Wood, A., Stankovic, J.: Denial of Service in Sensor Networks. *IEEE Computer* 35, 54–62 (2002)
5. Intanagonwiwat, C., Govindan, R., Estrin, D., Heidemann, J., Silva, F.: Directed Diffusion for Wireless Sensor Networking. *ACM/IEEE Transactions on Networking* 11(1), 2–16 (2003)
6. Gupta, G., Younis, M.: Fault-Tolerant Clustering of Wireless Sensor Networks. In: *Proc. of IEEE Wireless Communications and Network Conference (WCNC)*, pp. 1579–1584 (2003)
7. Lee, S.-B., Choi, Y.-H.: A Resilient Packet-Forwarding Scheme against Maliciously Packet-Dropping Nodes in Sensor Networks. In: *Proc. of the 4th ACM Workshop on Security of Ad Hoc and Sensor Networks*, pp. 59–70 (2006)
8. Chan, H., Perrig, A., Song, D.: Random Key Pre-Distribution Schemes for Sensor Networks. In: *Proc. of IEEE Symposium on Security and Privacy*, pp. 197–213 (2003)
9. Eschenauer, L., Gligor, V.D.: A Key Management Scheme for Distributed Sensor Networks. In: *Proc. of the 9th ACM Conference on Computer and Communication Security (CCS)*, pp. 41–47 (2002)
10. Perrig, A., Szewczyk, R., Wen, V., Culler, D., Tygar, J.: SPINS: Security Protocols for Sensor Networks. In: *Proc. of ACM MobiCom*, pp. 189–199 (2001)
11. Karlof, C., Sastry, N., Wagner, D.: TinySec: A Link Layer Security Architecture for Wireless Sensor Networks. In: *Proc. of ACM SensSys*, pp. 162–175 (2004)
12. Parno, B., Perrig, A., Gligor, V.D.: Distributed Detection of Node Replication Attacks in Sensor Networks. In: *Proc. of IEEE Symposium on Security and Privacy*, pp. 49–63 (2005)

13. Newsome, J., Shi, E., Song, D., Perrig, A.: The Sybil Attack in Sensor Networks: Analysis and Defenses. In: Proc. of IEEE IPSN, pp. 259–268 (2004)
14. Deng, J., Han, R., Mishra, S.: INSENS: Intrusion-Tolerant Routing in Wireless Sensor Networks. Technical Report CU-CS-939-02, Dept. of Computer Science, University of Colorado at Boulder (2002)
15. Tanachaiwiwat, S., Dave, P., Bhindwale, R., Helmy, A.: Routing on Trust and Isolating Compromised Sensors in Location-Aware Sensor Networks. In: Proc. of 1st International Conference on Embedded Networked Sensor Systems, pp. 324–325 (2003)
16. Zhu, S., Setia, S., Jajodia, S.: LEAP: Efficient Security Mechanisms for Large-Scale Distributed Sensor Networks. In: Proc. of the 10th ACM Conference on Computer and Communications Security (CCS), pp. 62–72 (2003)
17. Ganesan, D., Govindan, R., Shenker, S., Estrin, D.: Highly Resilient, Energy-Efficient Multipath Routing in Wireless Sensor Networks. *Computing and Communications Review (MC2R)* 1, 11–25 (2002)
18. Deng, J., Han, R., Mishra, S.: A Performance Evaluation of Intrusion-Tolerant Routing in Wireless Sensor Networks. In: Proc. of the 2nd International Workshop on Information Processing in Sensor Networks (IPSN), pp. 349–364 (2003)
19. Ye, F., Zhong, G., Lu, S., Zhang, L.: GRAdient Broadcast: A Robust Data Delivery Protocol for Large Scale Sensor Networks. *ACM Wireless Networks, WINET* (2005)
20. Carbunar, B., Ioannidis, I., Nita-Rotaru, C.: JANUS: Towards Robust and Malicious Resilient Routing in Hybrid Wireless Networks. In: Proc. of ACM Workshop on Wireless Security, WiSe (2004)
21. Morcos, H., Matta, I., Bestavros, A.: M2RC: Multiplicative-Increase/Additive-Decrease Multipath Routing Control for Wireless Sensor Networks. *ACM SIGBED Review* 2 (2005)
22. Yang, H., Ye, F., Yuan, Y., Lu, S., Arbough, W.: Towards Resilient Security in Wireless Sensor Networks. In: Proc. of ACM MobiHoc, pp. 34–45 (2005)
23. Hill, J., Szewczyk, R., Woo, A., Hollar, S., Culler, D., Pister, K.: Systems Architecture Directions for Networked Sensors. In: *ACU ASPLOS IX* (2000)
24. Ssu, K.-F., Wu, T.-T., Huang, S.-K., Jiau, H.C.: A Dependable Routing Protocol in Wireless Sensor Networks. In: Proc. of ICWMC, pp. 1–6 (2006)

Efficient Pairwise Key Establishment Scheme Based on Random Pre-distribution Keys in WSN*

Hao Wang, Jian Yang, Ping Wang, and Pu Tu

Chongqing University of Posts and Telecommunications,
400065, Chongqing, China
wanghaonet@163.com, yangjian_666@sina.com,
wangping@cqupt.edu.cn, tupu0803@126.com

Abstract. Key establishment and management are the cores of many security mechanisms and services in wireless sensor networks which are commonly deployed in a hostile environment. Pre-distributing keys into sensor nodes is a practical option for key management in wireless sensor environment. The negotiation and establishment of pairwise keys are much more important after the nodes are deployed, so that the whole network will work in expectation and safely. In this paper, we proposed an efficient pairwise key establishment scheme used in cluster-based random key pre-distribution protocol, which can provide perfect connectivity and security for large scale wireless sensor networks. Simulation results and performance analysis indicates that it also has other strongpoint, lower communication and computational overhead, much more scalable and flexible for the network size, a better balance between the node storage capacity and communication overhead.

Keywords: Wireless sensor networks, security, key pre-distribution, pairwise key.

1 Introduction

Wireless Sensor Networks (WSN), the main task is sending data which is collected from the sensor nodes to the base station in time, consists of a large number of sensor nodes. And the sensor node, which is powered by the built-in battery and equipped with a tiny sensor device, has limited capability of data processing and storage, and can carry out a certain distance wireless communication. WSN is a new research direction in information technology due to its large-scale of applications, such as in environmental monitoring, military, homeland security, traffic control, community security, forest fire prevention, targeting areas and so on. As the sensor nodes are mostly deployed in unattended or hostile areas, once the sensor network deployed in enemy controlled areas, the enemy will tap the network communication, palm off the right nodes, send malicious information to the network, so that they can destroy the network, and make it loss of normal function.

* The work reported in this paper is supported in part by the national science and technology major project No. 2009ZX03006-001.

Therefore, Sensor network security problem is especially prominent. In order to ensure the sensor networks works in safe way, when sensor nodes communicate with each other, the message transmitted between each other should be encrypted, important data must also be certified. To achieve this goal, we must first solve the key distribution problem, establishment of session key. To ensure that communications in node-to-node or nodes-to-header (such as the cluster head, the base station) is safe, we should establish a session key, namely, the pairwise key, between two nodes effectively and safely. Focus of this paper is to build the network model and give the establishment of the pairwise key option.

The main contributions of this paper are summarized in the following:

- 1) We propose our cluster-based model of the network topology, which is greatly based on practical application. This model achieves greatly strengthened security with excellent scalability and increases the connected probability.
- 2) We present the pairwise key establishment scheme, both with common keys directly and using different-intermediate-nodes path, which extremely guarantee the key security and secure communications in the future.
- 3) We give our analysis of selecting keys randomly and establishing pairwise keys via intermediate nodes, which excellently balances the storage and the communications.

The remainder of this paper is organized as follows. We briefly describe and discuss current key pre-distribution schemes in Section 2. We then give a characterization of our network topology and model in Section 3. A detailed introduction of our proposed preferable pairwise key establishment scheme is presented in Section 4. Section 5 gives analysis and evaluation of our proposed scheme. We conclude the paper and summarize our results in section 5.

2 Related Work

To ensure secure communications, the sensor nodes need to consult keys; according to the characteristics of sensor networks and the practical application, secret keys are pre-loaded into each sensor node by using random deployment method before they are deployed. There are several key pre-distribution models, the simplest way is that all nodes share a master key, but this method has poor security; or any two of nodes share the pairwise key, the security of such programs is good, but poor network expansion is not suitable for large-scale node network.

E-G program [9], the basic scheme, with the consideration of security and the space of storage in nodes, gives the basic random key pre-distribution scheme. This program is suit for sensor networks and complete the establishment of the keys in three phases, but did not give a clear way how to determine the pairwise key, or how to establish the pairwise key by the search path. Chan et al. [3] add the minimum requirements that any two of nodes share at least q common keys based the scheme [9]. This way increased security, but with the increase of the compromised nodes, the performance become relatively poor, meanwhile, it has the same deficiency with scheme [9].

Ren et al. [4] improve the structure of key ring in the basic scheme, give the method of sub-key pool allocation; Du et al. [5] build a multi-key space based on matrix theory; Liu et al. [6] use the polynomial instead of the key pool and construct the application of the safety threshold in the polynomial; Du et al. [7] and Liu et al.

[8] further extended the basic scheme using network pre-deployment knowledge. All of these have a clear way how to determine the pairwise key.

Blom et al. [11] proposed the symmetric matrix operations; in order to the network has the λ -secure property. Blundo et al. [12] proposed a polynomial-based key pre-distribution scheme. Scheme [13] [14] also use the symmetric polynomial $f(x, y) = f(y, x)$ to establish the pairwise keys, either for single base-station or multiple base-station. All these need much computing and have a poor scalability.

Our goal is to develop and improve the current key pre-distribution schemes. We give the detailed description how to establish the pairwise keys which provide a perfect security and connectivity. The scheme we proposed is not only reduced the keys stored in nodes but also improve the network's security.

3 Cluster-Based Deployment Model

3.1 Notation

For clarity, we list the symbols used in the paper below:

P key pool, set of keys randomly chosen from the setup server

$|P|$ size of the key pool, number of keys in the P set

M sub key pool

$|M|$ size of the sub key pool

m number of keys in a node's key ring

n network size, in nodes

t number of clusters

N_i node whose identifier is N_i

r the expected number of one cluster, $r=n/t$

p probability

B_u the broadcasting message of node u

k_{uv} pairwise key shared by u and v

N_u random value generated by u

$k_{ID_{uv}w_1}$ common key's ID between u and v , w_1 is the number

$E_{k_{uc}}\{M\}$ ciphertext of message M , encrypted by the key k_{uc}

3.2 Network Topology

Sensor network are usually with a large number of nodes, extensive coverage, but a single node only has a limited short-range radio communication. Except of communicating with the cluster head or base station, ordinary nodes usually communicate directly with the neighbor nodes. If the entire network is divided into several smaller regions, we call clusters in this paper; all nodes in a cluster neighboring each other can communicate with each other directly. At the same time, the key pool is also divided into the sub-key pools which correspond with the clusters. The sensor nodes from the corresponding sub-key pool select the keys, and then follow the division of the region to deploy the nodes. This way not only increases the probability of shard

keys within the same cluster of nodes, but also can reduce the number of keys stored in nodes. Through the overlapping sub-key pool, the nodes in adjacent cluster can be up to a certain probability, then, when using cluster head mechanism, inter-clusters communication can be completed through the cluster head. In this model, we use the head management mechanisms, cluster head products after achieving connectivity in the cluster; the specific method is not our focus of this paper. After sensor network is divided into several clusters, when the base station sends control commands, we can broadcast news according to the requirements of different regions, this has more pertinence. When the whole network has requires or control commands to, it broadcasts to the whole network. If this, it can reduce the numbers of broadcasting information.

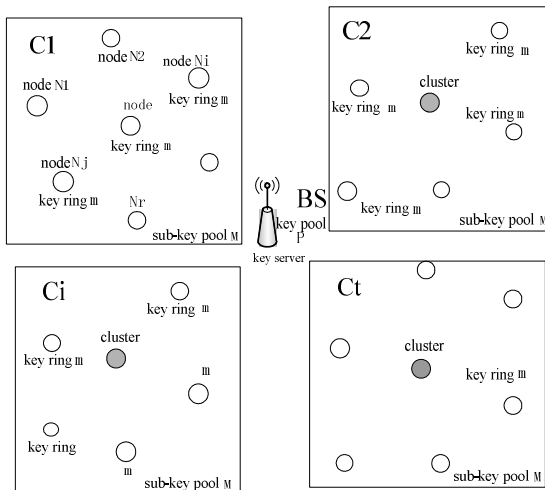


Fig. 1. Network topology of our model

The entire network topology model is shown in Fig. 1. Key server (base station) establishes key pool P by generating random numbers as keys, the size of key pool is $|P|$. Suppose that the total number of nodes in the network is n . It has the node clustering when deploying the networks. Assuming that it is divided into t , the number of nodes in each cluster is r . Before the deployment of the network, the sub-key pool selected randomly by the key server for each cluster is M , which has a number of $|M|$; in each cluster, for each node we take randomly select m keys from its sub-key pool M .

When the network deployed, it has a certain probability that any two or more of nodes share the same key within a cluster. When the two nodes have the same key, we can choose one or more as the pairwise key. As the network takes the cluster head as manager for each region, each cluster can consider a cluster head using some methods and establish the pairwise keys with the adjacent clusters. Obviously the different cluster still has a certain probability to get the same key.

3.3 Probability Model

According to this model of network topology, node connectivity has two kinds of situation, connectivity between two nodes in one cluster, namely, the probability of

existence of the same key in the different clusters. For the same cluster, in the network of Fig. 1, take C1 for example, the probability that Ni and Nj share at least one key is p. Ni selects m keys from |M| has a total of $C_{|M|}^m$ ways, Nj has also $C_{|M|}^m$ ways. Suppose that nodes Ni, Nj shared w keys, selecting w shared keys has $C_{|M|}^w$ methods, Ni and Nj still have 2(m-w) distinct keys, which has a kind of $C_{|M|-w}^{2(m-w)}$ selection ways. And to partition 2(m-w) different keys to the two nodes equally, we have $C_{2(m-w)}^{m-w}$ ways. So the probability that two nodes share w keys is

$$p(w) = \frac{C_{|M|}^w C_{|M|-w}^{2(m-w)} C_{2(m-w)}^{m-w}}{(C_{|M|}^m)^2} \quad (1)$$

The probability that two nodes share at least one key is

$$p = \sum_{i=1}^w p(i) \quad (w \leq m) \quad (2)$$

According to the probability of the complementary event, the result is

$$p_1 = 1 - p(0) = 1 - \frac{C_{|M|}^m C_{|M|-m}^m}{(C_{|M|}^m)^2} \quad (3)$$

Similarly, according to the above reason, the probability that two nodes in the different clusters shared at least one key is

$$p_2 = \left(1 - \frac{C_{|P|}^{|M|} C_{|P|-|M|}^{|M|}}{(C_{|P|}^{|M|})^2}\right) \left(1 - \frac{C_{|M|}^m C_{|M|-m}^m}{(C_{|M|}^m)^2}\right) \quad (4)$$

4 Efficient Pairwise Key Establishment and Management

The security connectivity of wireless sensor network is on the basis of communications connectivity; to ensure that any two nodes within a cluster are able to communicate, it requires that it can establish a secure path at least between any two nodes after the deployment of networks. Based on random key pre-distribution scheme, the key server loads the pre-shared keys to each node, that means that it can establish the communication link after self-organized, and the nodes guarantee the session key by consultation according to the predetermined rules prior to the establishment of the link, that is the pairwise key between any arbitrary nodes. Described in detail in this section about the building process of the pairwise key between the nodes within the cluster, the entire process of building the key is similar to the literature [9], but improve and perfect it, and make it closer to practical application.

4.1 Establishment of the Pairwise Keys Directly

When the deployment of a whole network is completed, each cluster enters bootstrapping phase. For each cluster, the node find the same key by broadcasting consultation

mechanism, and select one of the common keys (if there were multiple) as the pairwise key through the rules scheduled, the process of looking for shared key is as follows.

The node u broadcasts its own message that contains the node's ID u , random number N_u and m keys' IDs,

$$u \rightarrow * : B_u = \{u, N_u\} \cup \{k_{ID1}, k_{ID2}, \dots, k_{IDm}\}$$

When other nodes of the cluster, such as the node v , received the message B_u , it will compare all the ID with its own stored in the keys' ring, then reply to u with the same ID which is in order from small to large. Due to the rules predetermined by the node u and node v : they take the key of the smallest identifier as their pairwise key, and when this key is compromised or leaks, the key by followed will be as pairwise key. Node v sends a message to the node u as follows,

$$v \rightarrow u : M = \{v, N_u\} \cup \{k_{ID_{uv1}}, k_{ID_{uv2}}, \dots, k_{ID_{uvw}}\}$$

When u received all the nodes' messages of common key in the cluster, the pairwise keys will be stored in a form just as Table 1. ($1 \leq w_i \leq w_{i+1} \leq m$ ($1 \leq i \leq n'$))

The other nodes broadcast their own key messages similarly, and build the pairwise key tables with neighbor nodes. If network needs, such as unified management in cluster, data integration and so on, we can use hierarchical, that's to say, the establishment of cluster head for each cluster. The communication between clusters will be achieved by the cluster head who already had the pairwise keys with the cluster head of the adjacent clusters. The method establishing the pairwise key between the cluster heads is the same as the above.

Table 1. Pairwise keys stored in node u

Neighbor node	Pairwise key	common keys
v	K_{uv1}	$K_{uv2}, \dots, K_{uvw_1}$
a	K_{ua1}	$K_{ua2}, \dots, K_{uaw_2}$
b	K_{ub1}	$K_{ub2}, \dots, K_{ubw_3}$
c	K_{uc1}	$K_{uc2}, \dots, K_{ucw_4}$
.....

4.2 Establishment of the Pairwise Keys via Intermediate Nodes

In section 4.1, the pairwise key is established when there exist the common keys between two nodes, so it just needs one session. When the direct establishment of the

pairwise key is completed, the node-to-node within the cluster not only has the pairwise key, but also forms a table of shared keys. So does with the cluster head-to-head between the clusters. The table is composed of the neighbor nodes and their security links which have established their pairwise keys. The physically adjacent nodes which are not able to find common keys can create it via intermediate nodes. This section will give an approach about indirectly establishing the pairwise key.

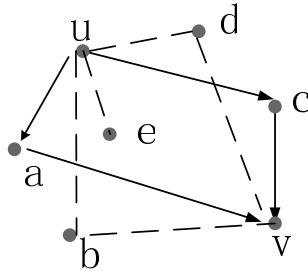


Fig. 2. Indirectly establishing the pairwise key

When there is no common key between two nodes, such as the node u and v, but it has formed the table of the pairwise key and the shared keys with other neighboring nodes. Now the node u and v have to establish the pairwise key, and it must be safe, not known to other nodes in the establishing process. The principle of establishment is: we choose the node which exists the least shared keys with the sending node as the forward transmitting node. Shown in Fig. 2, let u stored the shared key with a, b, c, d, e, it needs to choose two nodes from them, making the both nodes and v have the maximum probability of the common keys. Assume that a, c have the least numbers of shared keys in u's pairwise keys' table, then u chooses a, c as the forwarding node. The reason is that there is no shared key between u and v, yet a, c has the least shared keys with u, then the probability of existing shared keys with v will be the largest. The specific description of establishing process is as follows.

When node v receives the broadcast message of node u, found that there is no common key, and then returns the following message,

$$M = \{v, N_u, N_v\}.$$

Node u receives the message M, stores N_v , then selects a, c from the pairwise keys table to forwarding the key's materials in two paths according to above principle.

$$u \rightarrow a: M = \{u, E_{k_{ua}} \{v | N_{u1}\}\}$$

$$u \rightarrow c: M = \{u, E_{k_{uc}} \{v | N_{u2}\}\}$$

When node a, c receive their message M, if they have the pairwise keys with v, transmit them to v. As the node v receives the message from node a and c, calculates the pairwise key:

$$k_{uv} = N_v \oplus N_{u1} \oplus N_{u2},$$

and node u generates pairwise key using the same method.

When either a or c doesn't have pairwise key with v, u receives the response; it will choose other nodes in the pairwise key table in the same principle except a and c, repeats the above steps until successfully establish pairwise key with v. Other nodes use the same method to establish pairwise keys with the nodes they do not have shared keys in the cluster, and add the pairwise keys into their own pairwise keys table. If there is no shared key between the adjacent cluster-heads, they establish the pairwise key through the base station. In the whole network, each node established pairwise key among the nodes in the cluster and between the cluster-heads, we achieved the whole network connectivity.

5 Performance Analysis

In this section, we analysis and evaluate the security and performance of our scheme in connectivity, security, scalability and node storage. Connectivity: it is on behalf of the capability of the network security connectivity. It defined as the probability of establishing secure link when two nodes are within the scope of communication; for this article, it means that the probability of the two nodes in a cluster storing at least one common key, as well as the two nodes in different clusters. When nodes store fewer keys, the probability should be as high as possible. Security: the nodes in the network should ensure that expected security, which should store the least key and cost a minimum communication overhead. Node storage: since the node's limited storage resource, when the network reaches certain connectivity, the nodes should store fewer keys.

5.1 Network Connectivity

In the section 3.2, we have got the probability of at least one shared key existed in two nodes. Suppose the number of key pool is $|P|=10000$, when $|M|=(0.01,0.02, 0.03,0.04)|P|$, this probability is equation (3), the relationship between the probability and the keys stored in node is shown in Fig. 3.

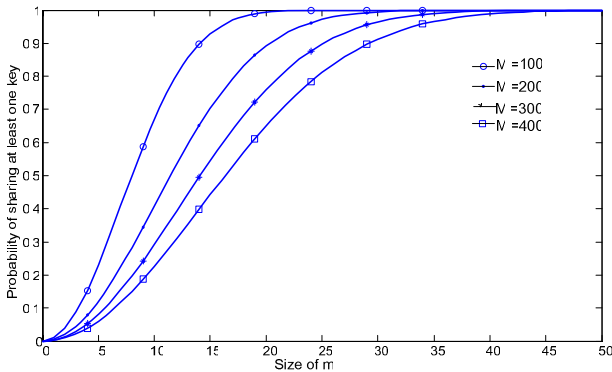


Fig. 3. Probability of sharing at least one key in two nodes in cluster

Table 2. In a certain probability, the number of keys stored in node in cluster

$ M \backslash P$	0.92~0.95	0.95~0.98	0.98~0.99
100	15	16,17	18,19
200	22,23	24-26	27,28
300	27,28	29-32	33,34
400	31-33	34-37	38-40

The Fig. 3 shows, when the sub-key pool is constant, the probability increases as the number of keys stored in nodes increases; nodes store the same amount of keys, the smaller the sub-key pool, the greater the probability; that's to say the possibility of finding pairwise keys will be larger. Even if the two nodes does not exist any same key, we can establish pairwise keys which is described in section 4.2. Use the method in section 4.2, we can guarantee the network wide connectivity. Table 2 shows when asked a certain probability, a single node, should store the number of keys. Compared to E-G scheme, the new scheme significantly reduces the number of keys.

For the adjacent two clusters, the probability of two nodes, equation (4), simplified as,

$$p_2 = (1 - f_1(M, P))(1 - f_2(m, M))$$

$$f_1(M, P) = \frac{(P - M)(P - M - 1) \cdots (P - 2M + 1)}{P(P - 1) \cdots (P - M + 1)}$$

$$f_2(m, M) = \frac{(M - m)(M - m - 1) \cdots (M - 2m + 1)}{M(M - 1) \cdots (M - m + 1)}$$

Fig. 4 shows, when sub-key pool is constant, the relationship between the probability and the number of keys stored in node.

Fig. 4 and Table 3 shows that sub-key pool is small, such as 100,200, the probability of finding common keys between clusters will be very small; when the sub-key pool increases, the probability significantly increases while the node storage has not increased much; such as when sub-key pool is 400, the amount of keys stored in different clusters is approximately as 1.4 times as the number of inner cluster. According to practical application, when we need a large amount of inner-cluster communication, we may properly increase both the number of sub-key pool and key ring stored in nodes. Then any node who will be the cluster head has a big probability with others'. Even if there is no shared key, we can establish pairwise key through the methods in 4.2. In actual application, the communication between the nodes in different clusters is relatively small, if need, we can transmit through the cluster head.

In all, we can see that the scheme has a good performance. When ensure connectivity, we largely reduce the number of key stored in node; and we can pre-configured key flexibility according to the specific application requirement. When cannot find shared keys between nodes, we can create pairwise key through a small amount of communications.

Table 3. In a certain probability, the number of keys stored in node in different clusters

$\frac{p}{ M }$	0.92~0.95	0.95~0.98	0.98~0.99	Maximum p
100	--	--	--	0.3895
200	--	--	--	0.8696
300	38-41	42-49	50-72	0.9900
400	43-46	47-52	53-57	0.9997
1000	69-74	75-85	86-92	1.0000
2000	99-106	107-122	123-132	1.0000
3000	121-131	132-150	151-162	1.0000
4000	140-152	153-173	174-188	1.0000

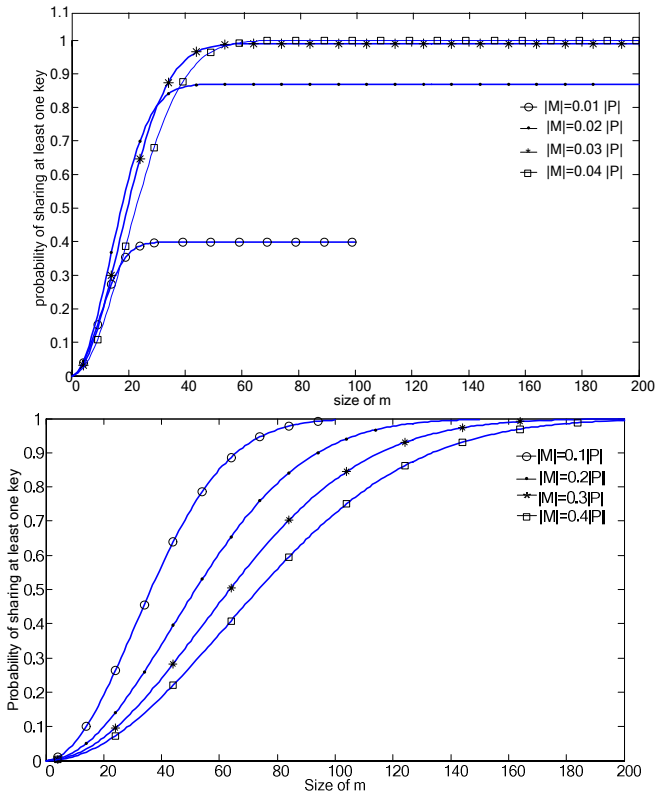


Fig. 4. Probability of sharing at least one key in two nodes in different clusters

5.2 Security Analysis

For the random key pre-distribution scheme, the network must establish all the communication keys in shortest time after the distribution of the network nodes. If there is neither leak of any key nor captured of nodes by attacker, the network is relatively supposed to be secure and all keys are successfully achieved by each node in the future communication. Otherwise, the message encrypted by the pairwise key will lose confidentiality, and one node compromised by adversary will pose a threat to others even the network.

The goal of this scheme is to allow sensor nodes to find a common secret key safely with each of their neighbors after deployment. So every node broadcasts its key's identifiers at the network's initial time to avoid disclose keys. We'll analyze the network's resiliency when one or more nodes are captured and disclose all their keys. Let the number of captured nodes be x . Since each node contains m keys, the probability that a given key has not been compromised is

$$\left(1 - \frac{m}{|S|}\right)^x.$$

Therefore, the expected fraction [6,16] of total keys being compromised can be estimated as:

$$p_3 = 1 - \left(1 - \frac{m}{|S|}\right)^x \quad (5)$$

When $|M|=400$, $|P|=10000$, $m=35$ and $m=55$, the results and comparison with existing key pre-distribution schemes are depicted in Fig. 5, which contains the in-cluster and global resilience of our scheme, E-G scheme[9] that is same as our global resilience and Chan-Perrig-Song scheme[3].The Fig. 5 show that our scheme substantially lowers the fraction of compromised communication after x nodes are compromised compared with the Chan-Perrig-Song scheme in network-wide. We get that if the x nodes were compromised in one cluster, the resilience of the network in that local area (in the cluster) is lower than that of the entire network. We also realize that in reality, these x nodes might be randomly distributed in the entire region, so the network will have a higher probability of resilience.

The result is that our scheme is secure and efficient. Even the adversary eavesdropped the message, it only get the key's ID instead of the keys themselves. For another, other pairwise keys are calculated by the nodes themselves. Because the key materials are transmitted through different intermediate nodes, the forwarding nodes only know one of the materials, random value in this method; they can't evaluate the pairwise key. So, the key established in this method is absolutely secure compared to using the secure channels in previous scheme; it is only known the two nodes who'll communicate with each other in future instead of any path-node or other third node.

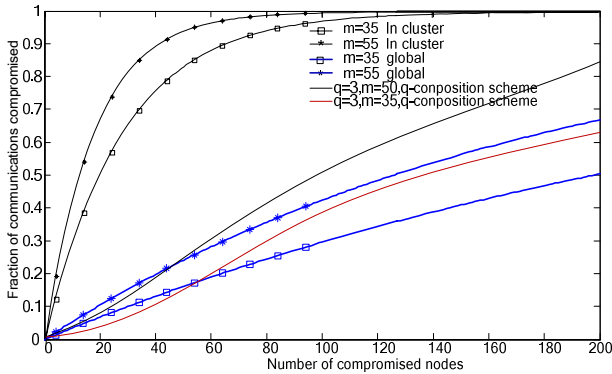


Fig. 5. Network Resilience: Comparisons with existing schemes

5.3 Communication and Computational Overhead

After the deployment of all the nodes in the network, to discover the common keys, each node of the schemes [9] [10] broadcast all of the keys or the key identifiers to their neighbor nodes. In particular, when the network size is large, the node will store a large number of keys in its subset. Then key information exchange process, not only consumes a lot of energy but also wastes considerable radio bandwidth, resulting in the link load overhead, collision opportunities to grow. Our method, according to specific requirements, has the strongpoint that nodes only need to store the appropriate number of keys in cluster, broadcast key identifier, to establish a certain number of pairwise key. The remaining nodes that do not share common keys only take a small amount of communication and calculation to establish the pairwise keys, which can reduce the network traffic load and energy costs as much as possible. With our method, the node transmits its key materials by using its neighbors who have the least common keys to the destination, which means he/she will find the right intermediate nodes in the largest probability. This path-materials computation and pairwise key establishment approach extremely has less CPU power and time consumption than previous two schemes. We can almost guarantee any two neighbor nodes to establish a pairwise key and have a safe communication.

In our scheme, they use the very key whose identity is the smallest as pairwise key in all key's IDs when the two nodes find their common keys. Otherwise, they only need four communications to get their key materials to calculate the key via XOR. This extremely reduces the computation overload for resource constrained sensor nodes especially their limited computing capability. Meanwhile, if they leak their pairwise key, they can easily delete it and take the key of the second smallest ID in their key tables as their later pairwise key while the two nodes have more than one common key. This is much flexible and secure for updating the pairwise key. Therefore its communication and computational overhead are much lower than previous schemes.

6 Conclusion

In this paper, we describe our pre-distribution key management model in WSN and propose a cluster-based key management scheme depending on the application

requirements. Since the main communication in the network is point-to-point, establishing pairwise keys between each two nodes is especially important. Therefore we present a new efficient pairwise key establishment scheme for large scale wireless sensor networks. Compared to the E-G scheme and the Chan-Perrig-Song scheme, our scheme has a number of appealing properties, our scheme can provide the full network connectivity and the substantially best resiliency, as well as the maximum supported network size for fixed key pool size, which is much more scalable and flexible. Lower computational and communication overhead are also achieved in our proposed scheme. Simulation results and performance analysis shows that this network model not only enhanced communication security but also greatly reduced the number of keys stored in nodes, and it also achieved a better balance between the node storage capacity and communication overhead.

References

1. Wang, Y., Attebury, G., Ramamurthy, B.: A survey of security issues in wireless sensor networks. *IEEE Communications Surveys and Tutorials* 8(2) (2006)
2. Su, Z., Lin, C., Feng, F.-J., Ren, F.-Y.: Key Management Schemes and Protocols for Wireless Sensor Networks. *Journal of Software* 18(5), 1218–1231 (2007)
3. Chan, H., Perrig, A., Song, D.: Random key pre-distribution schemes for sensor networks. In: *Proceedings of Symposium on Research in Security and Privacy*, Oakland, USA, pp. 197–213 (2003)
4. Ren, K., Zeng, K., Lou, W.: A new approach for random key pre-distribution in large-scale wireless sensor networks. *Journal of Wireless Communication and Mobile Computing (Special Issue on wireless Networks Security)* 3(6), 307–318 (2006)
5. Du, W., Deng, J., Han, Y.: A pairwise key pre-distribution scheme for wireless sensor networks. In: *Proceedings of the 10th ACM Conference on Computer and Communications Security*, Washington, DC, USA, pp. 42–51 (2003)
6. Liu, D., Ning, P.: Establishing pairwise keys in distributed sensor networks. In: *Proceedings of the 10th ACM Conference on Computer and Communications Security*, Washington, DC, USA, pp. 52–61 (2003)
7. Liu, D., Ning, P.: Location-based pairwise key establishments for static sensor networks. In: *Proceedings of SASN 2003*, Fairfax, USA, pp. 72–82 (2003)
8. Du, W., Deng, J., Han, Y.: A key management scheme for wireless sensor networks using deployment knowledge. In: *Proceedings of IEEE INFOCOM 2004*, Hong Kong, China, pp. 586–597 (2004)
9. Eschenauer, L., Gligor, V.D.: A key-management scheme for distributed sensor networks. In: *Proc. CCS 2002*, pp. 41–47. ACM Press, New York (2002)
10. Bollobas, B., Fulton, W., Katok, A., et al.: *Rand Graphs*, 2nd edn., pp. 160–200. Cambridge University Press, Cambridge (2001)
11. Blom, R.: An optimal class of symmetric key generation systems. In: Beth, T., Cot, N., Ingemarsson, I. (eds.) *EUROCRYPT 1984*. LNCS, vol. 209, pp. 335–338. Springer, Heidelberg (1985)
12. Blundo, C., Santis, A.D., Herzberg, A., Kutten, S., Vaccaro, U., Yung, M.: Perfectly-secure key distribution for dynamic conferences. In: Brickell, E.F. (ed.) *CRYPTO 1992*. LNCS, vol. 740, pp. 471–486. Springer, Heidelberg (1993)

13. Wang, Y., Ramamurthy, B., Xue, Y.: A Key Management Protocol for Wireless Sensor Networks with Multiple Base Stations. In: IEEE International Conference on In Communications, ICC 2008, pp. 1625–1629 (2008)
14. Wang, Y., Ramamurthy, B., Zou, X.: KeyRev: An efficient key revocation scheme for wireless sensor networks. In: Proc. ICC 2007, Glasgow, Scotland, UK (June 2007)
15. Traynor, P., Choi, H., Cao, G., Zhu, S., Porta, T.L.: Establishing pairwise keys in heterogeneous sensor networks. In: Proc. INFOCOM 2006 (2006)
16. Ma, C.-g., Shang, Z.-g., Wang, H.-q.: Domain-based key management for heterogeneous wireless sensor networks. *Journal on Communications* 01.30(5) (2009)
17. Chan, A.C.-F.: Distributed symmetric key management for mobile Ad Hoc networks. In: Proceedings of the INFOCOM 2004, Twenty-third Annual Joint Conference of the IEEE Computer and Communications Societies, vol. 4, pp. 2414–2424 (2004)
18. Liu, Z.-H., Ma, J.-F., Huang, Q.-P.: Domain-Based Key Management for Wireless Sensor Networks. *Chinese Journal of Computers* 29(9) (September 2006)

Agent Based Approach of Routing Protocol Minimizing the Number of Hops and Maintaining Connectivity of Mobile Terminals Which Move One Area to the Other

Kohei Arai¹ and Lipur Sugiyanta^{1,2}

¹ Department of Information Science,
Faculty of Science and Engineering, Saga University, Japan

² Department of Electrical Engineering,
Faculty of Engineering, State University of Jakarta
arai@is.saga-u.ac.jp, lipurs@gmail.com

Abstract. An ad-hoc network is a special kind of network, where all of the nodes move in time. So the topology of the network changes as the nodes are in the proximity of each other. Ad-hoc networks are generally self-configuring, because of the concept of no stable infrastructure takes a place. The research is focused in Wireless (Ad-Hoc) Mesh Networks. In this network, each node should help relaying packets of neighboring nodes using multi-hop routing mechanism. This deployment is needed to reach far destination nodes to solve problem of dead communication. This multiple traffic "hops" within a wireless mesh network caused dilemma. It can extend the coverage area, but the repeatedly relayed traffic will exhaust the radio resource. Wireless mesh networks that contain multiple hops become increasingly vulnerable to problems such as energy degradation and rapid increasing of overhead packets. Minimizing hop or relay is expected to solve this dilemma.

Keywords: wireless ad-hoc mesh network; multi-hop routing; energy, overhead packet.

1 Introduction

Wireless networks in today meaning are usually set up with a centralized access point for each area of connectivity. The access point has knowledge of all devices in its area and routing to them is done in a table driven manner. The future trend of mobility, it becomes important to give users the possibility of finding service they need anywhere in the network. Users would like to obtain access to services automatically, without reconfiguring their system. Especially with mobile devices dynamic discovery of services in other than home or corporate network and automatic network configuration will be very useful.

An ad-hoc network is a special kind of network, where all of the nodes move in time. So the topology of the network changes as the nodes are in the proximity of each other. Ad-hoc networks are generally self-configuring, because of the concept of no stable infrastructure takes a place. Each node in the network must be able to take

care of routing of the data. So the concept of centralized network with pre-defined central routing tables could be applied only with the limited nodes changing and coverage area. The coverage and connectivity of wireless mesh networks decreases as the end-user becomes further from the other nodes. Multi-hop communications among mesh nodes can enable long distance communication with acceptable performance degradation. Multi-hop wireless network is network where communication between two nodes may go through multiple consecutive wireless links without relying on any existing, pre-configured network infrastructure or centralized control. These networks are constrained by the interference of wireless communications, the finite battery energy and the mobility of nodes.

However, wireless network throughput and multi-hop coverage extension are two contradictory goals. On one hand, the multi-hop communications can extend the coverage area. On the other hand, as the number of hops increases, the repeatedly relayed traffic will exhaust the radio resource. The throughput will sharply degrade due to the increase of collisions from a large number of users. Therefore, it becomes an important issue to design a wireless network, so that the coverage of a wireless network can be extended without sacrificing the overall throughput.

Beside throughput, the ability to transport data from a source to a destination is a fundamental ability of wireless network. Delay is characterized by their lack of connectivity, resulting in a lack of instantaneous end-to-end paths. This is due to protocol routing trying to first establish a complete route and then, after the route has been established, forward the actual data. Also, wireless clients are mobile nodes and commonly run on batteries. Thus, power usage of mesh clients should be limited. These natures of wireless ad hoc networks made them require many improvements compared to wireless managed networks, before it ready to be implemented.

2 Wireless Mesh Network Architecture

A wireless mesh network is a communications network made up of radio nodes organized in a mesh topology. It can be seen as a special type of wireless ad-hoc network. In the IEEE 802.11s network, there is WLAN mesh and is defined as a set of mesh points interconnected via wireless links with the capabilities of automatic topology learning and dynamic path selection. Fig. 1 shows an example of IEEE 802.11s WLAN mesh. In the figure, there are two classes of wireless nodes. The mesh points (MPs) are the nodes supporting wireless mesh services, such as mesh routing selection and forwarding, while the non mesh nodes are the pure client STAs. In addition to mesh services, the mesh access point (MAP) also provides wireless access services. The pure client STAs do not participate in the WLAN mesh, but they can associate with the mesh APs to connect to the mesh networks. The WLAN mesh can connect to other networks by the mesh portals (MPPs). Multiple WLAN meshes can also be connected by the MPP.

The research used that network model with modified scheme. The architecture of the network is as follow. Nodes are mesh points with capability of routing selection and forwarding included in. Mesh access points that provide wireless services do not exist. Services are generated among nodes. Other nodes and new nodes that do not participate in the mesh network previously, can be associated to the networks by using a reached node which already joined the network as relay. The wireless mesh

can connect to other networks by the mesh portals (MPPs), but this scheme is not covered in this research. Multiple wireless meshes can also be connected by the MPP. There will be several wireless areas, each with the nodes inside. Nodes are always active moving with direction, speed, and distance change randomly. It is possible that nodes able to move inside wireless area and leaving its home area. During this activity, nodes always reconfigure their routing path allowing dynamically changing of route selection. Interaction among nodes became rapidly increased.

When moving, nodes also allowed entering new wireless area, leaving home/original area. In order to maintain connectivity along this movement, individual components will need to discover and maintain awareness of surroundings nodes and also to configure and adapt themselves in response to changing node's condition.

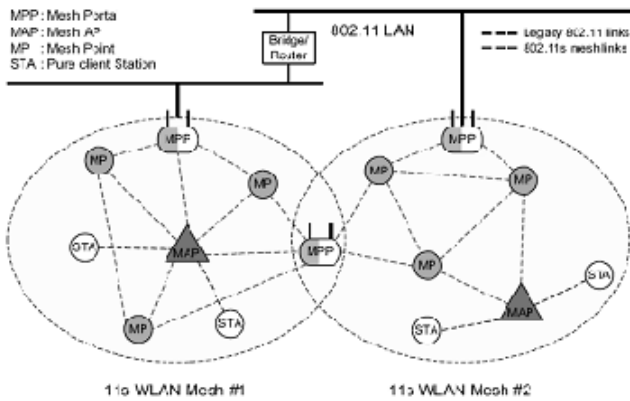


Fig. 1. The network architecture for the IEEE 802.11s WLAN mesh network

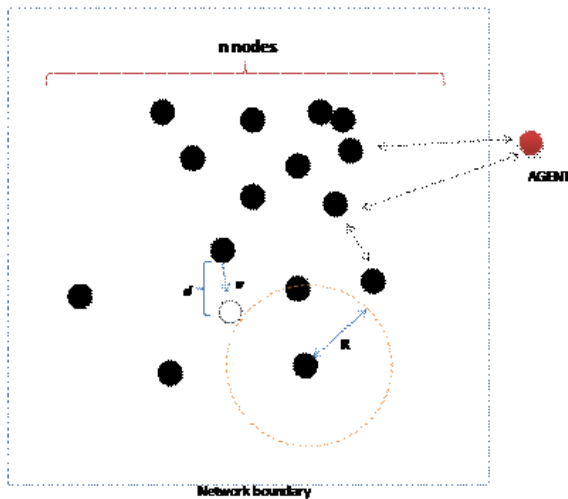


Fig. 2. Network model used in this research

There are special nodes called MPP. A MPP provides the ability to bridge wireless networks over a mesh. MPP has capacity to find optimum path with minimum hops to other nodes. Therefore MPP should support traffic propagation to both wireless areas at where it located. Nodes at both areas able to sent and receive data outside its home/original mesh and vice versa by using this MPP. It has consequences that nodes must aware about both their neighbors and also MPP. As long as communication path between source node and destination node is built inside wireless area, then MPP is not involved. In this paper, this MPP node is called Agent. The agent node reorganizes the network considering the available nodes.

Agents are located in such place to facilitate communication between wireless areas and minimize the number of hops to achieve optimum throughput of mesh clients which communicate each other. Through minimizing the number of hops between source and destination, it is hoped that the optimum throughput between source node and destination can be reached. The optimal tradeoff between network capacity and throughput for wireless mesh networks need to be investigated. Simulation programming approach is built to determine the optimal deployment parameters. This simulation runs with large number of nodes, with high dynamic parameters and considering mobile devices resources (battery power, bandwidth) and wireless environment constraints.

Simulation is split into two steps. First has to cover movement among nodes inside area, not move to another area. Second, simulation covers node's activities that include communication among nodes which also move to other areas. During the first step, route discovery play important role. Source nodes that do route discovery flood the network with a route request. Its neighbors propagate this route request in the same manner until it reaches the destination node or a node that already has a route to the destination. The data analyzed in this paper is intended to this first step.

3 Algorithm

3.1 Establishing the Initial Topology

Initialization topology. Nodes discover themselves and use their maximum transmission power to build the initial topology. We consider the case where all nodes in the network are similar, i.e., assuming an homogeneous infrastructure. Let us assume a square area of side l , in which n nodes are deployed uniformly, distributed at random position. This deployment produces a connected topology under some assumptions; sometimes a completely connected topology is built. From the simulation, a large connected component is built quickly using a communications radius considerably smaller than the radius needed to have the entire network connected.

3.2 Deriving the Dynamic Minimum Hop Tree

The concept used was backbone-based. This backbone approach is to find a subset of nodes that will be a gateway to guarantee connectivity and communication coverage while allowing every other node in the network to reach at least one node on this backbone in a direct way. These nodes provide complete network coverage. Fig. 3 shows an example of this approach in which 35 nodes out of 200 nodes in the network were selected to build the backbone. (nodes with red color are initiator nodes)

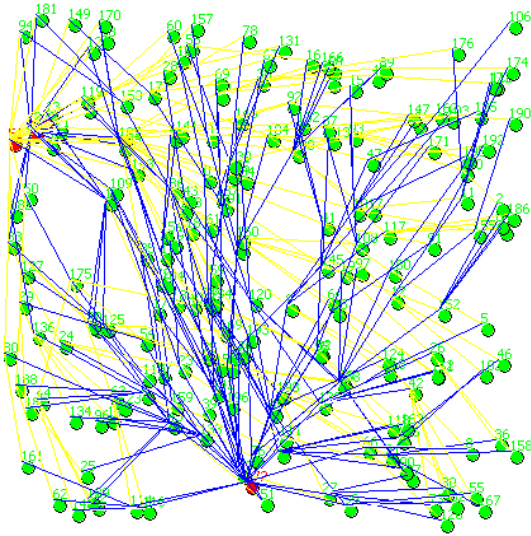


Fig. 3. 200 nodes distributed randomly in the area of 12, some nodes become backbone for others. Captured at initial time ($t=0$).

Topology is dynamically changed along as node moving. An easy way to illustrate the tree construction process is by the minimum spanning tree of a graph. The process works from the set of nodes that are initiator nodes at time t . Between t and $t + 1$, based on certain parameters, the nodes in the set evaluate all adjacent nodes in order to extend the tree. The process continues until all nodes on the graph are evaluated. Not every node will be selected to be part of the tree, and those which were not selected will have a chance until a new tree is requested. In general, the process usually starts on an initiator node, but if there are more than one initiator nodes, then several trees may be built in parallel.

The growing a tree algorithm builds connected graph considering the remaining energy in the nodes and the distance between them. The tree is built using three types of messages: Hello Message, Parent Reply Message, and Time Out Message. The building process is started by a predefined node that might be the initiator, right after the nodes are deployed. The initiator, node 72 in Fig. 3, starts the protocol by sending an initial Hello Message. This message will allow the neighbors of 72 to know their “parent”. This message also includes a selection metric that is calculated based on the signal strength of the received Hello Message and the remaining energy in the node. The metric will be used later by the parent node to sort the candidates. The parent node waits a certain amount of time to receive the answers from its neighbors. Once this timeout expires, the parent node sorts the list in decreasing order according to the selection metric. The neighbors of 72 that have already received Hello messages, after send Reply messages, and then prepare to send Hello messages to their neighbors. The process continues to all reachable nodes. If the neighbors have been already covered by another node, it ignores the Hello Message. The method that a node needs to call on in order to send a message is broadcast. The method that a node needs to send a packet is broadcast.

3.3 Routing Mechanism

Routing algorithms intended to maintain routes to all reachable destinations for sending data. Reactive ad-hoc routing algorithms use on-demand routing, so the routes are only requested when they are needed. This approach stores active routes in memory and no maintaining information for unused routes. Proactive ad-hoc routing algorithms (table-driven) maintain routes to all nodes in the network, so they rely according to tables. Nodes use reactive routing algorithm, whereas agents use proactive routing algorithm.

Main concepts of reactive routing algorithms are route discovery, route caching and route maintenance. Nodes used this method to discover route path. The first data (this means that no data has been sent to that destination before) initialize the route discovery from the source to the destination. A route reply is then sent back the same way the route request traveled. The route reply contains a list of all the intermediate nodes between source and destination. When the route reply reaches back to the source, the source can send the data and cache the new route for future use. Whenever a node discovers that the link to one of its neighbor nodes has failed, it will send a route error packet to the source of any cached routes that use the failed link. This route error will propagate through all intermediate nodes that also have cached this route. Upon receipt of it, node will update and purge the failed route.

On the other hand, agent nodes have enough knowledge to compute minimal route (min. hop distance) to any other destination nodes in both areas. In this simulation, agent nodes create source tree and computes route to the destination by computing the minimal-hop path. Each regular node keeps a path to agents.

3.4 Energy

Energy is power kept in each node. It was assumed that the radio dissipates $E_{elec} = 50$ nJ/bit to run the transmitter or receiver circuitry and $\epsilon_{amp} = 100$ pJ/bit/m² for the transmit amplifier. Thus, to transmit a k-bit message a distance d using this radio model, the radio expends:

$$E_{Tx}(k, d) = E_{elect} * k + \epsilon_{amp} * k * d^2 \quad (1)$$

and to receive this message, the radio expends:

$$E_{Rx}(k) = E_{elect} * k \quad (2)$$

The energy model included in simulation was based on the following formulas, taken from [9]:

$$E_{TXBit} = E_{elect} + (\epsilon_{amp} * (\pi r^2)) \quad (3)$$

$$E_{RXBit} = E_{elect} \quad (4)$$

The energy behaviors of nodes are defined as:

- During the idle time, a node does not spend energy. Even though this assumption has been proven untrue because being idle might be as costly as receiving data, this is still an assumption that can be done in most experiments, since the most important factor is the overhead in terms of message exchange and its associated cost.

- The nodes are assumed to have one radio for general messages. The main radio is used in all operations when the node is in active mode, and to send and receive control packets. When this radio is turned off, then no messages will be received and no energy will be used.
- Energy distribution among nodes can either be constant value, normally distributed, Poissonly distributed, or uniformly distributed.

4 Properties and Evaluation

Simulation was deployed on the abstract area, which is a rectangle where the nodes of the network are located. In order to create the deployment area, its width and height are manually defined.

The wireless mesh network model used in this simulation is based on the following assumptions:

- The communication range of the nodes is a perfect symmetric unit disk. If $d_{x,y} \leq r_x \rightarrow x$ and y can see each other.
- A constant bit error rate (BER) is defined for the whole network. It is a simply implementation of an error model. Whenever a packet is going to be sent, a random number is generated and compared to the message error rate (that depends on the size of the message). If the random number is greater, the message is received, otherwise it is lost. The default value for the BER is 0, which means there is no packet loss.
- To increase the randomization of the simulation process, simulation introduces some delay on some common processes in the network, like message transmission delay, processing delay, time out, etc. so each instance of a simulation, would produce different results.
- Simulation manages homogeneous networks, in which all nodes have the same characteristics. The parameters that can be defined in a homogeneous family of nodes are number of nodes, communication radius, and size of the area of deployment, position distribution, and energy distribution.

4.1 Initiator Nodes

Initiator node is node that initiates transmission of packet. It can be route discovery or data transmission. Like other nodes, initiator is always moving with random direction, speed, and distance. Along moving, initiator node is always sensing its neighbor to maintain connectivity.

4.2 Node's Energy Performance

Nodes uses radio channel that is symmetric such that the energy required transmitting a message from node A to node B is the same as the energy required transmitting a message from node B to node A for a given SNR. For our experiments, we also assume that all nodes are sensing the environment at a fixed rate. Energy was used whenever node sends, receives, and processes message.

The average of remaining node's energy at the end of simulation (with 200 cycles) for all nodes is shown in Fig 4.

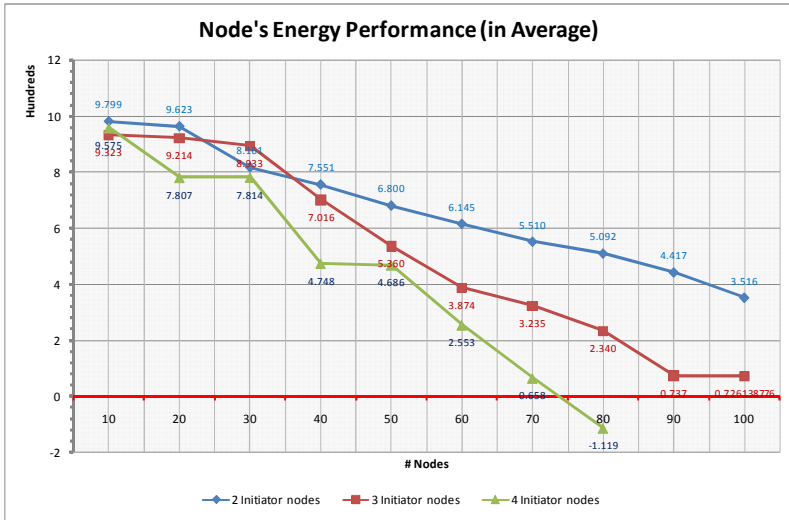
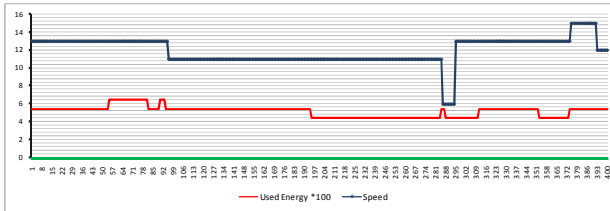
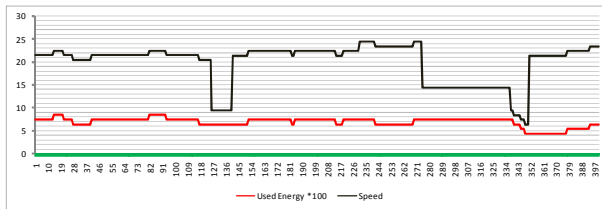


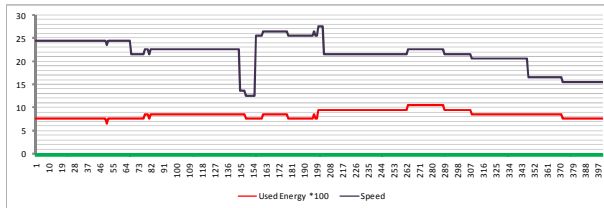
Fig. 4. Node's energy performance in relationship with number of nodes



(a)



(b)



(c)

Fig. 5. The faster node is moving, the more Energy is be used

Each node was given energy of 1000 at the beginning of simulation. Simulation executed repeatedly with number of initiator nodes increment by 10. As number of initiator node increase, the total energy remained on nodes is decreased. Simulation with 4 initiator nodes gave result negative average value at # Nodes= 80. It was to show that some nodes were already dead/OFF.

With node always moving, simulation shows that energy used during cycles is affected with its speed. More energy is used by node when node moves faster. Fig. 5 below illustrated the effect of speed to energy used by initiator node.

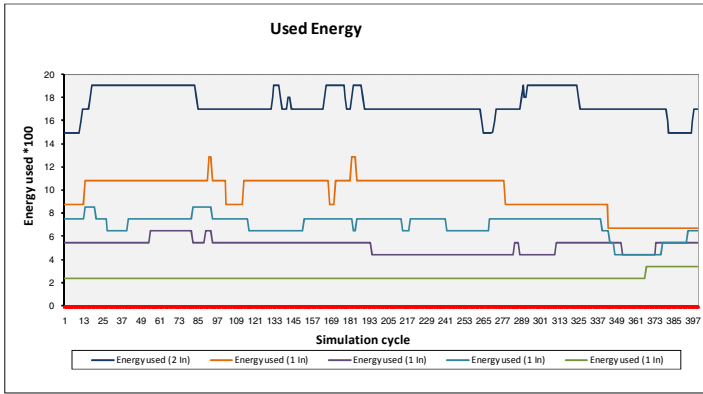


Fig. 6. Comparison for energy used to number of initiator node

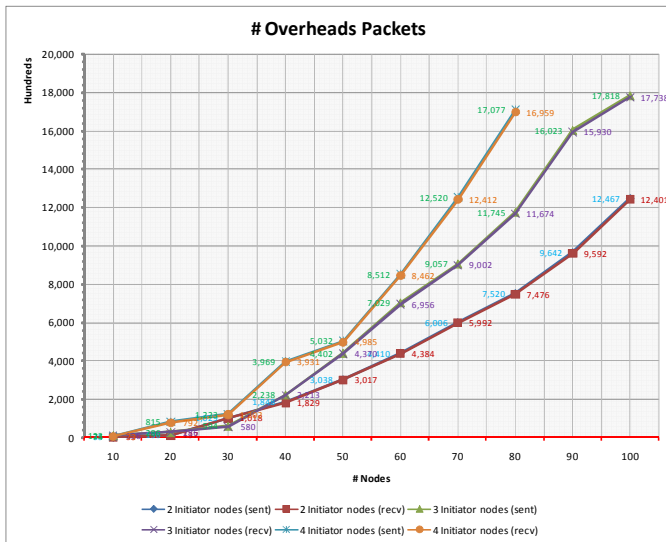


Fig. 7. Overhead packets increase rapidly as more number of initiator nodes

It is clear from Fig. 6, the increasing number of initiator nodes caused nodes to route a large number of messages. Thus these nodes will die out quickly, causing the energy required to get the remaining messages to increase and more nodes to die. This will create a cascading effect that will shorten system lifetime. In addition, as nodes close to the initiator die, then the environment is no longer being monitored. To prove this point, we ran simulations using the random 10 – 100 node shown in Fig. 5. After completed cycle, several nodes were considered dead.

Using Equations 3 and 4 and the random node from 10 until 100, we simulated transmission of first data from initiator to every node (communication range 100 m). Fig. 4 showed the total energy remained on nodes. Topology was changed all the time cycle. With transmission energy is on the same order as receive energy, it can be seen that the energy-efficient to use depends on the network topology.

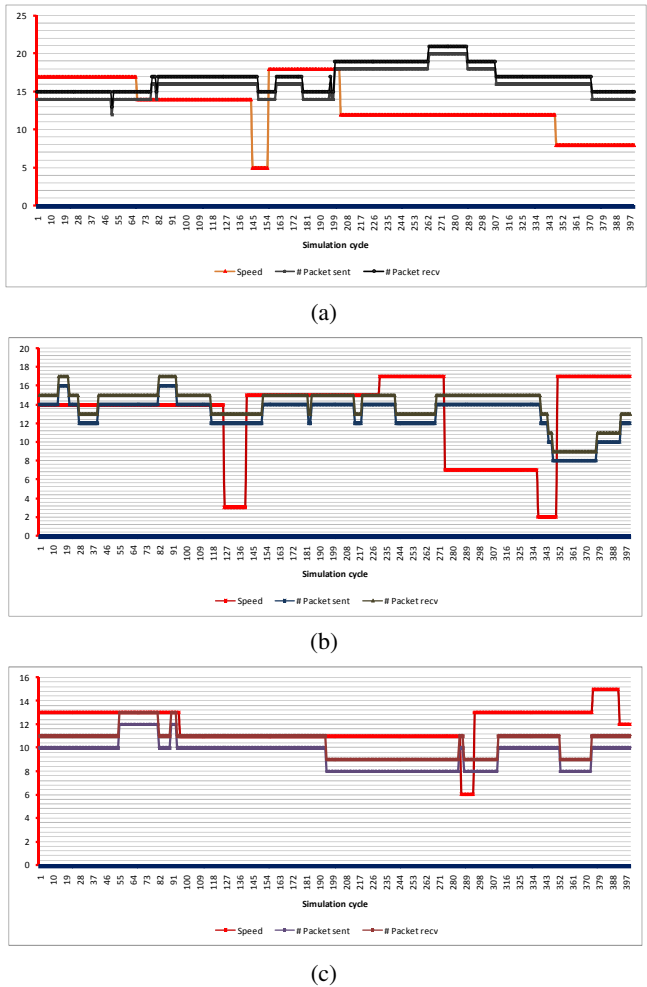


Fig. 8. Number of overhead packets is affected with speed of node

4.3 Overhead of Topology Discovery

Overhead packet refers to the extra bytes it takes to transmit data on a packet-switched network. Each packet requires overhead that is format information stored in the packet header, which, combined with the assembly and disassembly of packets, reduces the overall transmission speed of the raw data.

In this simulation, TCP/IP will be assumed to be the protocol being used for packet data. It was because much of the world's telecommunication networks are forwarding IP packets to and from at this moment. During the course of an IP packet's journey from a node to others, it is encapsulated and de-encapsulated into and out of framing headers and trailers that define how the packet will make its way to the next hop in the path. This Fig. 8 below is to explain what happens when to overhead packets when number of nodes increase and so as number of initiator nodes.

The overhead is incurred on a per-packet basis, so the smaller the packets, the more of them there are, the more overhead. The bigger the frames, the fewer there are, the less overhead is required to send them. This holds true for any technology that supports variable length frames, including wireless network. Beside that, multi-hop routing imply high packet overhead, then more nodes in the network means more hop available. From [10] shows that the packet overhead of the multi-hop routing is extremely high compared to single path routing since many nodes near the shortest path participate in packet forwarding. This additional overhead caused by moving node can cause congestion in the network. Thus Fig. 8 is simply understood. Overhead increases almost linearly with speed of node. (red line is for speed data, whereas other lines are for overhead being sent and received).

4.4 Impact of Multi-hop Transmission

In this section, we validate multi-hop transmission of overhead packet with simulation result.

- With direct communication protocol, each source node sends its data directly to the destination node. If the destination is far away, direct communication will require a large amount of transmit power from each node (since d in Equation 1 is large). Thus nodes route data through intermediate nodes. In this simulation, the intermediate nodes are chosen such that the transmit amplifier energy ϵ_{amp} is minimized. Each message must go through n transmits and n receives. Depending on the relative costs of the transmit amplifier and the radio electronics, the total energy expended in the system might actually be greater than direct transmission.
- In multi-hop environment, the nodes closest to the destination will be used to route a large number of messages to destination. Thus these nodes will die out quickly. This will create a cascading effect that will shorten system lifetime. Fig. 5 might prove this point.
- Fig. 9 and Fig. 8 show the packet overhead with various initiator node speeds. The packet overhead of the multi-hop routing is extremely high compared to single path routing since many nodes near the shortest path participate in packet forwarding.

The accumulated transmission delay and increasing number of hops among nodes that lowered the throughput can be seen from the previous research results. In [7], variables of evaluation are the number of hops, length of data transmission, and round trip time between nodes. The data obtained is shown in Fig. 8 below.

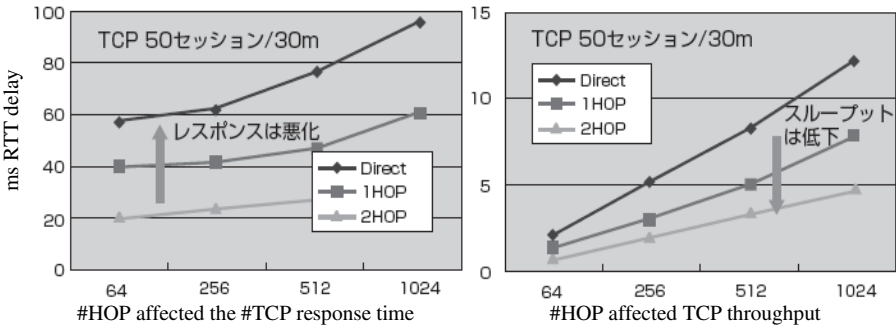


Fig. 9. The effect of multi-hop transmission in Wireless Mesh Network [7]

Those results agreed with Fig. 5 of the simulation results.

The research examined various products from wireless vendors with different architectures, which is set for the same measurement condition and were measured separately. One group is wireless LAN devices equipped with ad hoc mode technology and mesh network routing capability. Meanwhile, other product group is a wireless LAN product which expected to function as access point.

5 Related Works

Many researchers have evaluated general problems of performance of multi-hop wireless mesh network.

One branch of related work is concerned with the evaluation of wireless mesh product considering their applicability into the enterprise network. A second branch of work considers the evaluation of node behavior that composes multi-hop radio network.

Nemoto Nozomu [7] considered to introduce the technical review of wireless mesh network products IEEE802.11 standard with fixed wireless mesh network node. In terms of performance, the research will evaluate the overwhelming use of outdoors Muni-WiFi with a view to applying the same technology in the corporate network.

There are many products and technology for wireless LAN. The difference technology among products includes route, control technology and architecture. Wireless LAN is not competing with legacy network, otherwise it can be considered as complement. In this research, wireless LAN and wireless mesh network is compared to determine the level of performance. Table 1 shows the main contents of evaluation.

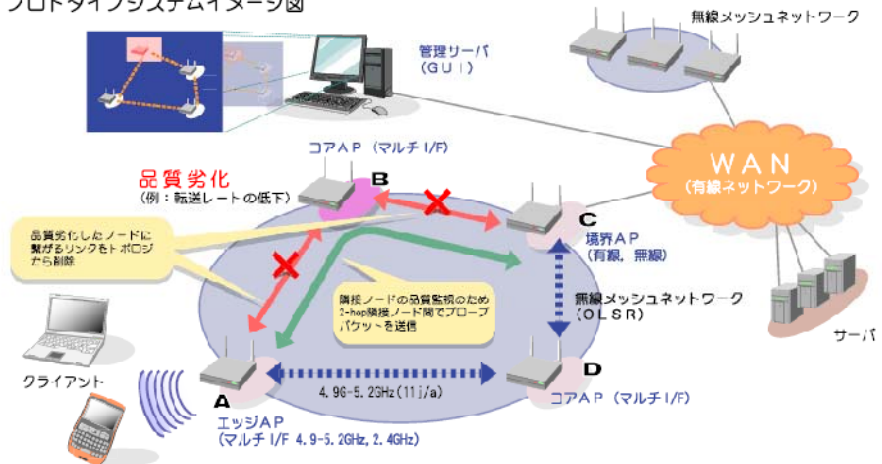
Table 1. Content of evaluation

Item	Test
Installation	Ease of installation
	Ease of scale expansion
Basic performance	Basic transmission performance
	Performance of voice and data transmission
	QoS control function
	Cracking resistance performance
Switch failure	Management functions
	Switching time when disability
Stability	Switching behavior
	Load of stability of network behavior

The second research [12] evaluates an autonomous and distributed control architecture based on the node behavior that composes the multi-hop radio network and the QoS evaluation of the traffic. The developed multi-hop radio network that uses this architecture maximize the autonomous node using wireless as a communication media, monitors the behavior each neighbor nodes and also the traffic in the node that composes the network. It aims at the establishment the architecture of technology for the network construction of autonomous decentralized node with high reliability.

The relay node that composed the multi-hop radio network cooperatively executed and control the autonomous network developed. In a past technology, the behavior of node was determined so that each relay node reports to a certain manager. This kind of technique was suitable if there was correct operation of each relay node, and there was no problem of breakdown and attack. In this research, the behavior of the node is evaluated by its adjacent nodes. As a result, it becomes possible to evaluate the behavior of each relay node accurately when the attack and the un-anticipated trouble. The stabilized operation of the entire network is attempted by excluding the relay node that not meets a standard. The prototype is like the figure below.

プロトタイプシステムイメージ図

**Fig. 10.** The prototype of wireless mesh project [12]

6 Consideration for Dissemination to Multi Areas

Wireless mesh network and wireless technology have emerged in recent years. What type of nodes moving, either inside the home (indoors) electronics appliances or outdoors ITS (Intelligent Transport Systems) is a major attraction point. The problem of routing and performance characteristics such as node capacity and power, expected to be applied to a wide range of subject.

Several other papers provide analytical formulas for calculating performance of wireless network. The main difference with [7], [12] and related papers is that we are here focusing on the performance of wireless network (e.g. #hop, throughput, delay, energy) with nodes which dynamically active. Some results from this simulation tend to agree with previous research results. To the best of our knowledge, in both of previous researches, the nodes are static/not moving. The simulation result can be enhanced through expanding the simulation with multiple areas without violating the continuous connectivity requirement.

6.1 Agents

The agent nodes are special nodes that are included to receive the information from all active nodes in both neighborhood wireless areas. They served as bridges and used to transport the sensed data to its final destination somewhere between the wireless network areas. They are also in charge of initiating, executing, and/or controlling the topology construction and maintenance protocols, routing protocols, etc.

Despite the fact that in real life the hardware configuration of the agent nodes may different compared to a regular wireless device, simulation uses the same data structure to model both nodes, with the main difference that agent nodes are considered to have more amount of energy. Agents in the simulation have two main functions: as service agents and delivery agents.

1. Service agents: takes a place relay between areas which is advertised with service request from source nodes. They could create paths and provide optimum hop/path to destination node.
2. Delivery agents: collectors of information received from nodes in both areas.

6.2 Power Efficiency

Power efficiency is defined as the ratio of the total transmission power over the total maximum transmission power at all nodes and denoted as η . In the worst case that each node still uses its maximum transmission power after the execution of the topology control algorithm, the power efficiency is the minimum value, $\eta = 0$. The higher the power efficiency, the more power saved on transmitting in the network. Theoretically, the upper bound of η is 1. However, practically, η can never be 1 as the transmission power can not be 0 in order to ensure the network connectivity.

7 Conclusion

In this paper, we evaluate topology of wireless network with multi-hop routing. One of the main processes for building hop tree is energy-efficient. The algorithm is done

in not two-hop information queries, which reduces the amount of overhead. Transmission energy is set same as receive energy and topology always changed all the time, and then the energy to use depends on the network topology. Topology is always changed, so as the initiator nodes. Increasing number of initiator nodes caused nodes to route a large number of messages. Thus these nodes will die out quickly, causing the energy required to get the remaining messages to increase and more nodes to die. This will create a cascading effect that will shorten system lifetime. Multi-hop routing, however, imply high packet overhead, then more nodes in the network means more hop available. The packet overhead of the multi-hop routing is extremely high compared to single path routing since many nodes near the shortest path participate in packet forwarding. This additional overhead caused by moving node can cause congestion in the network. This paper doesn't include agent yet. However, from the existing result, several parameters must be paid attention. Energy on each node decreased as number of initiator node increase. It caused by high intensity of nodes to sense neighbors as number of nodes increases. As hop increases and there exist agents as bridge, it is estimated that the energy will decrease rapidly. In future work, we plan to implement agents to minimize hops to achieve optimum transmission while maintaining connectivity among areas. Further investigation need to be conducted on dynamic routing advantages and factors which affect routing mode, e.g., flow type, delay, and etc. We will also further explore the throughput/delay/reliability tradeoffs between wireless network areas that deploy agents and without agents.

Acknowledgments. The authors would like to thank the anonymous reviewers for the helpful comments and suggestions.

References

1. Ettus, M.: System Capacity, Latency, and Power Consumption in Multihop-routed SS-CDMA Wireless Networks. In: Radio and Wireless Conference (RAWCON 1998), August 1998, pp. 55–58 (1998)
2. Rodoplu, V., Meng, T.H.: Minimum energy mobile wireless networks. *IEEE Journal on Selected Areas in Communications* 17(8), 1333–1344 (1999)
3. Lin, X., Stojmenovic, I.: Power-Aware Routing in Ad Hoc Wireless Networks. In: SITE, University of Ottawa, TR-98-11 (Dec. 1998)
4. Meng, T., Volkan, R.: Distributed Network Protocols for Wireless Communication. In: Proc. IEEE ISCAS (May 1998)
5. Manfredi, V., Hancock, R., Kurose, J.: Robust routing for dynamicmanets. In: ACITA 2008 (2008)
6. Ye, Z., Krishnamurthy, S., Tripathi, S.: A framework for reliable routing immobile ad hoc networks. In: Proceedings of IEEE INFOCOM (2003)
7. Nozomu, N.: Consideration and Evaluation of Wireless Mesh Network. Nomura Research Institute (NRI) Pacific Advanced Technologies Eng., 70–85 (2006)
8. Jarábek, B.: Advanced Network Approaches for Wireless Environment. In: Bielíková, M. (ed.) IIT. SRC 2005, April 27, pp. 108–115 (2005)
9. Heinzelman, W., Chandrakasan, A., Balakrishnan, H.: Energy-efficient communication protocol for wireless microsensor networks. In: Proceedings of the 33rd International Conference on System Sciences (HICSS), pp. 1–10 (2000)

10. Oh, S.Y., Gerla, M.: Robust MANET Routing using Adaptive Path Redundancy and Coding. In: COMSNET 2009, Bangalore, India, January 10 (2009)
11. The Dynamic Source Routing Protocol for Mobile Ad Hoc Networks (DSR), RFC Draft (July 2004),
<http://www.ietf.org/internet-drafts/draft-ietf-manet-dsr-10.txt>
12. Kyushu University, Kyushu Institute of Technology.: Developing Multi-Hop Wireless Communication Technology (WiRem), Press Release Project, March 17 (2008)

Ensuring Basic Security and Preventing Replay Attack in a Query Processing Application Domain in WSN

Amrita Ghosal¹, Subir Halder¹, Sanjib Sur², Avishek Dan², and Sipra DasBit²

¹ Dept. of Comp. Sc. & Engg, Dr. B. C. Roy Engineering College, Durgapur, India
ghosal_amrita@yahoo.com, subir_ece@rediffmail.com

² Dept. of Comp. Sc. & Tech., Bengal Engineering and Science University, Shibpur, India
{sanjib.sur11, avishekdan}@gmail.com, siprad@hotmail.com

Abstract. Nodes in a wireless sensor network are susceptible to various attacks primarily due to their nature of deployment. Therefore, providing security to the network becomes a big challenge. We propose a scheme considering cluster architecture based on LEACH protocol to build a security mechanism in a query-processing paradigm within wireless sensor network. The scheme is capable of thwarting replay attack while ensuring essential properties of security such as authentication, data integrity and data freshness. Our scheme is lightweight as it employs symmetric key cryptography with very short-length key. We simulate our scheme to show its efficacy of providing basic security to the network as well as detecting replay attack in the sensor network. Further we compare our scheme with one of the existing schemes taking packet loss and packet rejection ratio as performance metrics.

Keywords: replay attack, denial of service attack, authentication, data integrity, data freshness.

1 Introduction

The advent of efficient short range radio communication and advances in miniaturization of computing devices have given rise to strong interest in wireless sensor networks (WSNs) [1]. A wireless sensor network (WSN) consists of a large number of small, battery-powered [2], [3] wireless sensor nodes. The nodes are distributed in ad-hoc fashion and process sensing tasks. A sensor network may be deployed in hostile environments where there can be malicious attackers. On the other hand, the WSNs are used in critical application domain e.g. defense, medical instrument monitoring. Therefore, securing the activities of WSNs is of utmost importance.

These WSNs are prone to attacks [4], [5] due to their nature of deployment in an unattended environment and also due to the broadcast nature of the medium. One such attack is Denial of service (DoS) attack [6] that possesses a great threat to wireless sensor network environment. One form of DoS attack attempts to disrupt the network service, may be by blocking communication between nodes. The other form of DoS attack is the path-based DoS (PDoS) attack where the network is flooded with bogus packets along a multi-hop data delivery path [7]. The packets are continuously

replayed back to the sink node. This attack is known as replay attack. However, the objective of both types of such attacks is to eliminate or diminish the network performance and thereby hamper the working of the whole system. In replay attack, because of the broadcast nature of sensor networks, an adversary eavesdrops on the traffic, injects new messages and replays/changes old messages.

Many works are so far reported towards the solution of various forms of DoS attack. Deng et al. [7] have considered a type of DoS attack along a multihop data delivery path. As WSNs are generally tree structured, so an attack on the nodes of a path also affects the branches connected to that path. A one way hash chain mechanism has been proposed by the authors to prevent such path based DoS attacks and protect end to end communication along a multihop data delivery path. Here an OHC (One way Hash Chain) number is added with each message packet leading to an extra overhead of 8 bytes per packet. These 8 bytes of additional overhead is a major constraint for a resource constrained sensor node.

Authors in [8] have proposed an authentication protocol capable of resistance against replay attacks and other DoS attacks. But the scheme uses symmetric keys where keys are shared by sensor nodes and therefore a compromised node can send forged messages which may possess a great security threat to the network.

Perrig et al. [9] have used RC5 algorithm for encrypting messages which means a lot of computations have to be done. Here the μ TESLA protocol is used for secure broadcast of messages and an 8-byte Message Authentication Code (MAC) is used by each node for verification along the communication path. However, exchange of huge authentication information is a real bottleneck for the resource-constrained WSN. Moreover a node must have prior knowledge of all the nodes along its path.

The authors in [10] have devised a mechanism to prevent replay attacks. Here it is considered that the packets used for time-synchronization of any two nodes are replayed. A receiver-receiver model for time synchronization has been considered where a reference node broadcasts beacon message to all its neighbours. Based on the arrival time of the beacon message the receiving nodes adjust their clocks. A time offset is calculated based on the difference between the recording times of the beacon message by these receiving nodes. The time offsets are exchanged between the receiving nodes to calculate a threshold value which is the difference between the two time offsets of two nodes. But the arrival of a beacon message at a node can be delayed by certain factors leading to gross errors while deriving the time required for synchronization between the nodes. Moreover, by not using global time synchronization model a large overhead has been introduced as huge number of time offsets need to be computed by the nodes.

Dong et al. [11] have proposed the use of hash chains in their work where each node combines the hash value with its own node-id and forwards this combined value to its next higher hop count node. The receiving node is able to detect replay attacks by observing the combined value of node-id and hash value. But the computation of all these values by nodes takes significant amount of time.

Soroush et al. [12] have developed a scheme to defend replay attacks where a monotonically increasing counter is maintained to keep track of old replayed messages. But here each node maintains a counter to store the timing information of all other nodes which requires a large amount of memory leading to a major bottleneck for memory- constrained sensor nodes.

In this paper we propose a secured query processing scheme in WSN with a target to build a security mechanism from within a query-driven application environment. The proposed security mechanism gives a solution of replay attack while ensuring authentication, data integrity and data freshness.

The rest of the paper is organized as follows. In section 2 system model along with a brief description of the present work is given. A detailed description of the scheme is given in section 3. Section 4 gives the performance evaluation of the proposed scheme and also the comparative study between the proposed scheme and another scheme. Concluding remarks and future scope has been stated in section 5.

2 System Model

An attack [6] is defined as an attempt to gain unauthorized access to a service, resource, or information, or the attempt to compromise integrity, availability, or confidentiality. DoS is one type of attack which hinders communication among the sensor nodes. The present work considers replay attack, which is one form of path based DoS (PDoS) attack. If there is a PDoS attack, the network gets flooded with bogus packets along a multi-hop delivery path [7]. The occurrence of bogus packets in the network is due to the replaying of packets by adversaries leading to replay attack. So PDoS attacks in the network may result in replay attacks.

The system model in the present work considers clustered network architecture based on LEACH protocol [13]. In this architecture nodes organize themselves into local clusters with one node acting as the cluster head (CH). LEACH performs local data fusion to reduce the amount of data sent from the clusters to the base station/sink. Once all the nodes are organized into clusters, each CH assigns a time slot (TDMA schedule) for the member nodes in its cluster. The member nodes sense data and transmit data to the cluster head nodes that are located at one hop distance away from them. The cluster head nodes transmit the same to the base station after receiving data from their member nodes.

The present work considers a query-driven application platform where base station generates query messages and broadcasts the query. Cluster heads receive the broadcasts query and start the registration process to authenticate their respective member nodes. Once the registration phase is over, the CHs forward the query messages to their members. Depending on the nature of queries, specific member nodes send response-packets.

The objective of the proposed scheme named as Secured Query Processing Scheme (SQPS) is to ensure basic security in general and prevent replay attack in particular while working for a query-driven application. The essential properties [14] of a WSN required for maintaining basic security within the network are:

- Authentication- A sensor network must be able to differentiate between data coming from reliable sources and those coming from adversaries.
- Data Integrity- The received data should not be tampered during communication between the sender and the receiver.
- Data Freshness- The data received at a particular time should be current data and not old data which may be replayed by adversaries.

3 The Scheme

In this secured query processing scheme as shown in figure 1, periodically queries are broadcast from base station. The process begins as soon as the cluster heads receive queries from base station. The scheme has two phases-

- Registration Phase
- Query Response Phase

3.1 Registration Phase

This phase is used for registering member nodes by the respective cluster head nodes. The objective of this phase is to register only the authenticated nodes. Moreover the phase provides a mechanism that will help a CH to ensure data integrity and data freshness during query processing phase. Upon receiving a query from the base station, cluster heads initially broadcast registration packet.

The registration packet contains 8 bits including '0' as MSB. This MSB differentiates between a registration packet and a query packet sent by the cluster heads. As mentioned in section 2, each member node is allotted a particular time slot which is used for sending a registration response packet to the cluster head corresponding to the registration packet and also for receiving identification (node-id) for the member node from the cluster head for continuing the query processing session. On receiving the registration packet, the respective member nodes decide the number of bits to be shifted and accordingly left shift the bits of registration packet. Then the nodes generate registration response packet including left-shifted registration packet (8 bit), number of bits left shifted (3 bit), and present time-stamp (12 bit). The key used here is the number of left shifted bits. Symmetric key cryptography is used as the same key is used for encryption and decryption of the data packets. Symmetric key cryptography is beneficial for sensor networks as less computation has to be done and memory requirement is also minimized. Here, though the key is being sent along with the registration packet it will not be possible for any adversary who captures or eavesdrops the contents of the packet to decipher which bits in the packet refer to the key. Moreover in traditional networks the adversary performs some computations to obtain the key used in the network; this is not possible in case of sensor networks as the adversary nodes are also equipped with less computational power. The present time-stamp indicates the time when a member node sends a reply packet in response to the registration packet. The time-stamp is represented in minutes and seconds.

Immediately after receiving a registration response packet from one of its member node, the CH right shifts the reply packet designated number of times. The number of bits to be right shifted is same as the number of bits the registration packet is left shifted and it is provided as control data in the registration response packet itself. If the CH can retrieve the original registration packet after right shifting the received registration response packet, the member node is considered as an authenticated and registered node. Upon authenticating a member node, the CH generates a node-id and sends it to the member node at the same time slot of authenticating and registering the member node. Once the registration of a member node is successful, the CH stores this information.

3.1.1 Encryption at Registration Phase

Let us consider an 8-bit binary registration packet as $S = m_7 \dots m_0$, where MSB is m_7 and LSB is m_0 . The registration packet is broadcast by cluster head at a particular time slot. Time slots are assigned to member nodes according to TDMA schedule. On receiving the registration packet S , member nodes randomly choose the shift bit B for left shifting. If B is 100, S is encrypted as $m_3m_2m_1m_0m_7m_6m_5m_4$ (S'). Accordingly a 3-tuple registration response (RR) packet is formed by the member node as $(S'', B, T_{MN_present})$, where $T_{MN_present}$ is present time-stamp.

3.1.2 Decryption at Registration Phase

Upon receiving the encrypted registration response packet $(S'', B, T_{MN_present})$ from a member node, the cluster head performs the reverse process of encryption to get back the original registration packet. So S' is right shifted B times and converted to S'' . After decrypting, if S'' matches with S , the cluster head accepts the registration packet and authenticates the corresponding member nodes.

3.1.3 Data Stored by Cluster Head

Once the cluster head authenticates a member node, it generates an identification (node-id) for the member node and stores a 6-tuple data (node-id, S , S' , B , $T_{MN_present}$, $T_{MN_previous}$) for the member node where S' , B , $T_{MN_present}$ is RR (packet sent by the node designated by node-id), S is the original registration packet and $T_{MN_previous}$ is a time-stamp which is set to null initially.

3.1.4 Data Stored by Member Node

Once a member node gets its node-id from the cluster head during registration, the member node stores 3-tuple (node-id, B , $T_{MN_present}$) data.

3.1.5 Algorithm for Registration

Begin

// Action executed by cluster heads (CHs)

1: CH broadcasts 8-bit registration packet (S) at time t

// Actions executed by member nodes (MNs)

2: On receiving S at time $(t+1)$

3: for $i=t+1, i \leq t+n, i++$

4: left shift S to obtain S'_i /* i^{th} member node (MN_i) encrypts S by left shifting it by B_i -bit which is arbitrarily chosen by MN_i to obtain S'_i */

5: Generate RR_i /* MN_i generates 3-tuple registration response packet $RR(S'_i, B_i, T_{MN_i_present})$ */

6: Send RR_i to CH

// Actions executed by cluster head (CH) on receiving RR_i from MN_i

```

7:   right shift  $S'_i$  to obtain  $S''_i$            /* CH decrypts  $S'_i$  by right shifting it by
                                            $B_i$ -bit to obtain  $S''_i$  */
8:   if  $S''_i = S$  then
9:        $RR_i$  is authentic                    /* registration response packet  $RR_i$  is
                                           authentic */
10:      accept  $RR_i$                           /* CH accepts registration response packet
                                            $RR_i$  */
11:      generate node-id for  $MN_i$ 
12:      send node-id to  $MN_i$ 
13:      store 6-tuple data                    /* 6-tuple data: ( node-id $_i$ ,  $S$ ,  $S'_i$ ,  $B_i$ ,
                                            $T_{MN_i\_present}$ ,  $T_{MN_i\_previous}$  )*/
// Action executed by  $MN_i$  on receiving node-id $_i$ 
14:      store 3-tuple data                    /*3-tuple data: node-id $_i$ ,  $B_i$ ,  $T_{MN_i\_present}$  */
15:   else
16:       reject  $RR_i$ 
17:   end if
18: end for
End

```

Node authentication is one of the major parameters of network security. The registration phase of the present scheme is able to ensure this part of security.

3.2 Query Response Phase

In the event of query processing, the base station broadcasts query. Upon receiving the query, a CH starts registration phase. Once the registration phase is over (section 3.1), the CH broadcasts the query for its member nodes. Depending on the query, specific member nodes send response and cluster head nodes forward the responses to base station. For example, if the query is related to temperature, then member nodes responsible for sensing temperature provide response to the query.

Query packet is of 8 bits. The MSB bit of the query packet is always 1, to distinguish between registration packet and query packet. So, 128 different queries can be generated with the remaining 7 bits of the query packet.

3.2.1 Encryption at Query-Response Phase

On receiving query response member nodes respond to the query in the form of a query response (QR) packet that contains 5-tuple (node-id, m , m' , $T_{MN_present}$, $T_{MN_previous}$) where m is an 8-bit response message for the corresponding query and m' is left shifted (B times) message. Further, $T_{MN_present}$ and $T_{MN_previous}$ are

present time stamps when the member node is sending the response packet for the query and previous time stamp of last communication (packet exchange) between member node and CH respectively. To start with, $T_{MN_present}$ of the 3-tuple data stored by the member node at registration phase replaces $T_{MN_previous}$.

3.2.2 Data Stored by Cluster Head

After receiving query response packet from member node, the cluster head updates the 6-tuple data stored at registration phase. During registration phase the stored 6-tuple data was (node-id, S, S', B, $T_{MN_present}$, $T_{MN_previous}$). During query response phase, S and S' attributes are replaced by m and m' respectively. $T_{MN_previous}$ which contained null value during registration phase is now replaced by $T_{MN_present}$ which was stored during registration phase and $T_{MN_present}$ is replaced by the time when the query response packet is being sent by the member node.

3.2.3 Decryption at Query-Response Phase

Cluster head on receiving the QR packet, finds the node-id of the member node from the packet. Then to decrypt the encrypted message m', the CH right shifts m', B times where B is found from the stored tuple corresponding to the node-id of the member node. The m' is converted to m" after decryption. Once the CH decrypts the message part of the QR packet, it checks for data integrity, data freshness and replay attack.

3.2.4 Check for Data Integrity

If it is found that the decrypted message m" is same as the original message m, data integrity is preserved.

3.2.5 Check for Replay Attack and Data Freshness

In replay attack the malicious node repeats the already sent packets and results in energy exhaustion of nodes and eventually collapse of the network. With the help of two attributes $T_{MN_present}$ and $T_{MN_previous}$ stored at cluster head corresponding to every registered/ authenticated member nodes, replay attack is detected. The cluster head compares the $T_{MN_present}$ and $T_{MN_previous}$ of an entry corresponding to a member node and if $T_{MN_present}$ stored by it and the $T_{MN_previous}$ in QR is equal, then it can be ensured that no replay attack has taken place and data freshness is preserved.

If a malicious node attempts to send query response packet posing as an authenticated member node, cluster head rejects the packet. Due to this unauthorized attempt, an authorized node's TDMA time slot will be consumed and as a result of which a packet to be sent by the authorized node is lost. Therefore, there may be two distinct effects of malicious attempt- one is packet loss and the other is packet rejection.

3.3 Diagrammatic Representation of the Scheme

The entire scheme is illustrated with the help of an activity diagram shown in Figure 1. The diagram shows how a cluster with n number of nodes and one cluster head works. To be more specific, it shows communication among different components of a wireless sensor network and control flow among various computational modules within the components of the network. Broadcast communication from cluster head is shown by lightning symbol whereas unicast communication between cluster head and member nodes and control flows are shown by solid lines with arrowhead. All communication symbols are labeled to make understand the content of the packet along with the timestamp at which it is being sent. For example, one broadcast communication labeled by RB_t means registration packet has been broadcast at time t. Timestamp helps to know the steps of operations visibly.

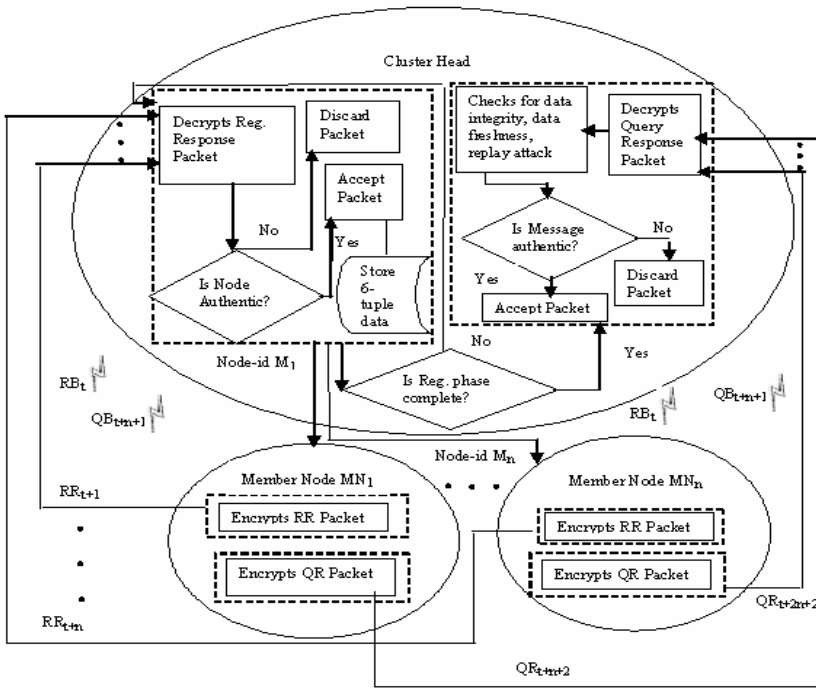


Fig. 1. Activity diagram of the proposed scheme (SQPS)

The notations used in Figure 1 has been described below-

- RB_t – Registration Packet broadcast at time t by cluster head.
- RR_{t+1} – Registration response packet transmitted at time t+1 by member node M_1 .
- RR_{t+n} – Registration response packet transmitted at time t+n by member node M_n .
- Node-id M_1 – Node-id of member node M_1 transmitted at time t+1 by cluster head.

Node – id M_n – Node-id of member node M_n transmitted at time $t+n$ by cluster head.

QB_{t+n+1} – Query packet broadcast at time $t+n+1$ by cluster head.

QR_{t+n+2} – Query response packet transmitted at time $t+n+2$ by member node M_1 .

QR_{t+2n+2} – Query response packet transmitted at time $t+2n+2$ by member node M_n .

4 Performance Evaluation

The effectiveness of the proposed security scheme reported in the earlier section is evaluated through simulation.

4.1 Simulation Environment

Simulation is performed using MATLAB (version 7.1). We consider 500 nodes in the network and number of malicious nodes is varied from 25 to 100.

Performance of the scheme is evaluated based on the following two metrics:

Authentication Rate – Number of authenticated nodes / total number of nodes in the network.

Data freshness (%) – (Number of received packets containing current data / total number of packets sent) x 100.

The relevant parameters and their associated values are listed below in Table 1-

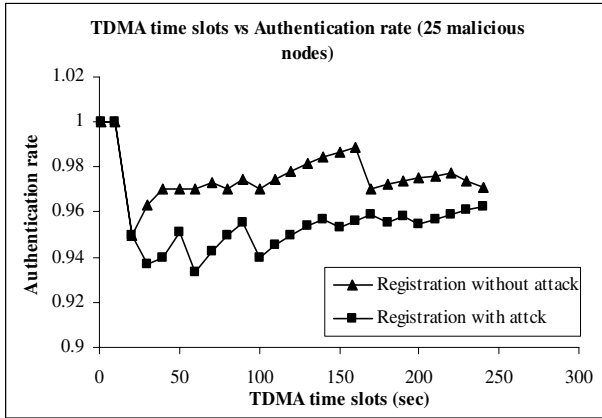
Table 1. Parameters and their corresponding values used in simulation

Parameters	Value
Initial Energy ($E_{initial}$)	2 J
Network area (user input)	25m × 25m to 100 m × 100 m
Communication range of sensor (R_c)	160 m
Sensing range of sensor (R_s)	80 m

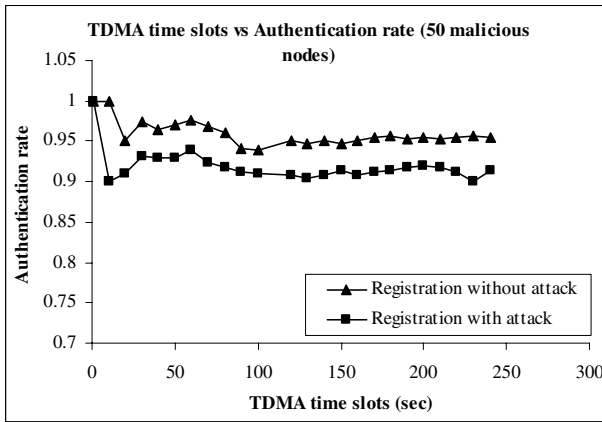
In presence of replay attack, data freshness is affected. In other words, if replay attack can be considered as a cause, data freshness is an effect. Therefore, the metric, data freshness is measured to cover analysis on both the security parameters replay attack and data freshness. Authentication rate is computed to cover another security parameter i.e. authentication. No separate experiment is carried out to measure data integrity as violation of data integrity due to presence of malicious node does not arise.

4.2 Simulation Results

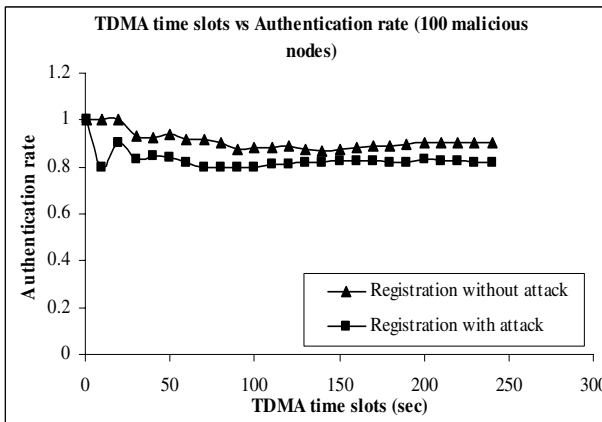
Authentication rate is measured and plotted with time in Figure 2 for varying number of malicious nodes. Two sets of experiments are conducted to compute authentication rate.



(a) with 25 malicious nodes



(b) with 50 malicious nodes



(c) with 100 malicious nodes

Fig. 2. Authentication rate over a period of time

In one set of experiment, we consider that malicious nodes are present and attempt to participate for sending data but stop participation once it is refused to do so. Results are plotted in Figure 2 for varying number of malicious nodes and designated as ‘without attack’. In the other set of experiment, malicious nodes attempt to send data repeatedly resulting in replay attack. Results of the same are plotted and designated as ‘with attack’.

We observe that in all the cases of Figure 2 ((a), (b), (c)) results ‘with attack’ show a fall of authentication rate compared to the results of ‘without attack’. This signifies presence of replay attack from malicious nodes. Further, if we compare results of all the plots ((a), (b), (c)), it is observed that authentication rate decreases with increase in number of malicious nodes.

Figure 3 shows percentage of data freshness over a period of time in presence of malicious nodes. Results in presence of 25 malicious nodes show that average data freshness is near about 96% whereas its values are 88% and 81% for 50 and 100 malicious nodes respectively. The results indicate that data freshness is inversely related to the number of malicious nodes present in the network.

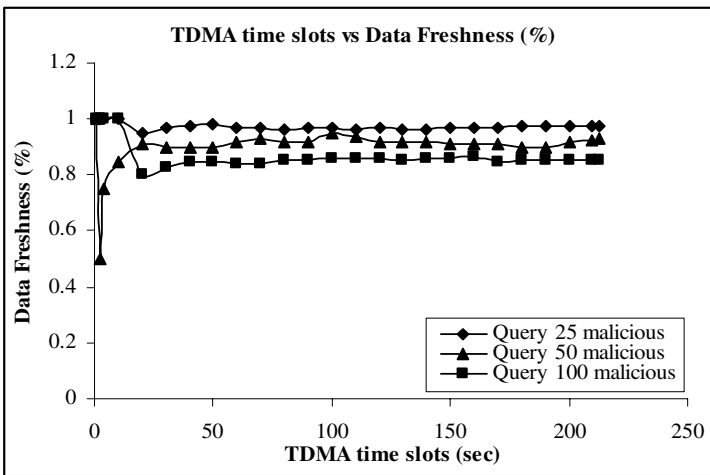


Fig. 3. Data freshness (%) over a period of time

4.3 Comparative Study

Due to unavailability of an existing suitable work so that our scheme can be compared based on all the security parameters considered here, we have chosen a work [11] on secured routing scheme to compare packet loss and data packet rejection ratio as comparison metrics as defined below.

Packet loss (%) – Total number of received packets by a CH / total number of packets sent by the member nodes.

Packet rejection ratio (%) – Total number of received packets by a CH containing tampered data / total number of packets sent by the member nodes.

Packet loss and packet rejection ratio (%) are computed and plotted for varying number of malicious nodes in Figure 4 and Figure 5 respectively.

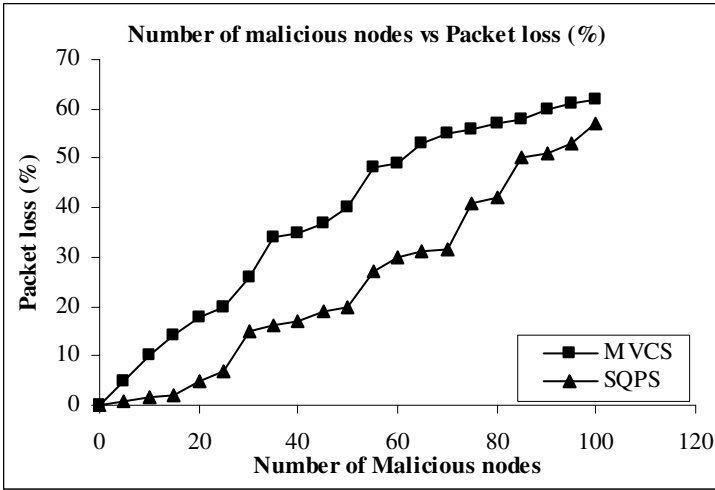


Fig. 4. Packet loss by varying number of malicious nodes

The figures also show the results for secured routing scheme named as MVCS (mitigating attack against virtual coordinate system) [11]. We observe that packet loss (Figure 4) increases with the increase of malicious nodes for both the schemes. However, in our scheme (SQPS), packet loss is about 16% less than MVCS for up to 25 numbers of malicious nodes. For 25 to 85 numbers of malicious nodes packet loss in SQPS is about 20% less than MVCS and this onwards packet loss is about 16% less in SQPS. Summarily, it can be said that packet loss in SQPS is less than MVCS for all sets of values of malicious nodes.

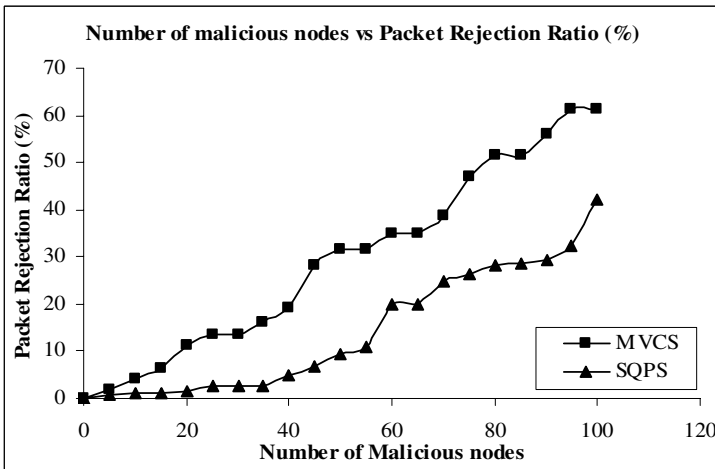


Fig. 5. Packet Rejection (%) by varying number of malicious nodes

Further we observe that packet rejection ratio (Figure 5) increases with the increase of malicious nodes for both the schemes. But in our scheme (SQPS), packet rejection is about 10% less than MVCS for up to 25 malicious nodes. For 25 to 100 malicious nodes packet rejection ratio in SQPS is 20% less than MVCS and this trend continues. So it can be inferred that packet rejection in SQPS is less than MVCS for all sets of values of malicious nodes. As all the member nodes' identities are verified through registration phase, there is very little chance that a malicious node is able to steal the identity of a legitimate node and passes through registration phase. That is why SQPS packet rejection ratio is lower than MVCS.

5 Conclusion

In this paper, we have proposed a scheme to defend replay attacks on nodes of WSN as well as preserve the essential basic security properties such as authentication, data integrity and data freshness of such a network. The scheme is designed in such a manner that no malicious node can take part in actual query processing thereby ensuring authentication.

Moreover, as there is no participating malicious node, violation of data integrity due to attack has been eliminated. However, a malicious node can attempt to participate stealing some time slots due to which there may be some packet loss. The merit of the scheme lies on the fact that simple symmetric key cryptography has been used to maintain security making the solution very lightweight. We have substantiated our claims by simulating the scheme in presence of attacks. Finally the scheme is compared with one of the existing routing schemes considering packet loss and data packet rejection ratio as comparison metrics. Results show our scheme outperforms the existing one.

As a future extension, the scheme may be made more realistic considering cluster head nodes are also vulnerable to attack. Further enhancement may be done to make it applicable for continuous data-flow application domain as well.

References

1. Akyildiz, I.F., Su, W., Sankarasubramaniam, Y., Cayirci, E.: A survey on sensor networks. *IEEE Communications Magazine* 40(8), 102–114 (2002)
2. DasBit, S., Ragupathy, R.: Routing in MANET and Sensor Network- A 3D position based approach. *Journal of Foundation of Computing and Decision Sciences* 33(3), 211–239 (2008)
3. Halder, S., Ghosal, A., Sur, S., Dan, A., DasBit, S.: A Lifetime Enhancing Node Deployment Strategy in WSN. In: Lee, Y.-h., et al. (eds.) *FGIT 2009*. LNCS, vol. 5899, pp. 296–308. Springer, Heidelberg (2009)
4. Ghosal, A., Halder, S., DasBit, S.: A Scheme to tolerate Jamming in multiple nodes in Wireless Sensor Networks. In: *Proceedings of Wireless VITAE*, pp. 948–951. IEEE Press, Los Alamitos (2009)
5. Ghosal, A., Halder, S., Chatterjee, S., Sen, J., DasBit, S.: Estimating delay in a data forwarding scheme for defending jamming attack in wireless sensor network. In: *Proceedings of 3rd International Conference NGMAST*, pp. 351–356. IEEE CS Press, Los Alamitos (2009)

6. Wood, A.D., Stankovic, J.A.: Denial of Service in Sensor Networks. *IEEE Computer* 35(10), 54–62 (2002)
7. Deng, J., Han, R., Mishra, S.: Limiting DoS attacks during multihop data delivery in wireless sensor networks. *Journal of Security and Networks* 1(3/4), 167–178 (2006)
8. Liu, D., Ning, P.: Efficient distribution of key chain commitments for broadcast authentication in distributed sensor networks. In: *Proceedings of 10th annual Network and Distributed System Security Symposium*, pp. 263–276 (2003)
9. Perrig, A., Szewczyk, R., Wen, V., Culler, D., Tygar, J.: SPINS: Security Protocols for Sensor Networks. *Wireless Networks Journal (WINET)* 8(5), 521–534 (2002)
10. Song, H., Zhu, S., Cao, G.: Attack-Resilient Time Synchronization for Wireless Sensor Network. In: *Proceedings of IEEE International Conference on Mobile Adhoc and Sensor Systems Conference*, vol. (7-7), pp. 772–779 (2005)
11. Dong, J., Ackermann, K.E., Bavar, B., Nita-Rotaru, C.: Mitigating Attacks against Virtual Coordinate Based Routing in Wireless Sensor Networks. In: *Proceedings of 1st ACM conference on Wireless Network Security*, pp. 89–99 (2008)
12. Soroush, H., Salajegheh, M., Dimitriou, T.: Providing Transparent Security Services to Sensor Networks. In: *Proceedings of IEEE International Conference on Communications*, pp. 3431–3436 (2007)
13. Heinzelman, W.R., Chandrakasan, A., Balakrishnan, H.: Energy Efficient Communication protocol for Wireless Microsensor Networks. In: *Proceedings of the 33rd Hawaii International Conference on System Sciences*, vol. 2, pp. 8020–8029 (2000)
14. Zia, T., Zomaya, A.: Security Issues in Wireless Sensor Networks. In: *Proceedings of International Conference on Systems and Network Communications*, p. 40 (2006)

A Review of Redundancy Elimination Protocols for Wireless Sensor Networks

Babak Pazand, Amitava Datta, and Rachel Cardell-Oliver

School of Computer Science & Software Engineering
The University of Western Australia
35 Stirling Highway, Crawley, W.A. 6009, Australia
{babak,datta,rachel}@csse.uwa.edu.au

Abstract. Sensing the same region of a physical environment by many sensor nodes results in wastage of energy. This behavior conflicts with the most important requirement of a wireless sensor network which is power efficiency. Thus, it is essential to have a mechanism to eliminate high redundancy while preserving an acceptable level of network coverage. In this paper, we review the proposed protocols in the literature that aim to eliminate the sensing redundancy in wireless sensor networks.

Keywords: Wireless Sensor Networks, Sensing Redundancy, Node Scheduling.

1 Introduction

Coverage is one of the most important metrics for evaluating the quality of service (QoS) in wireless sensor networks. A region in a sensor field is considered as a covered region if every point in this area is in the sensing radius of one or more sensors. However, it is important to consider how efficiently the sensor field is covered by the deployed sensors.

Two different aspects of the coverage problem have been discussed in the literature. The first aspect considers the efficiency of coverage of a deployment region in terms of the percentage of the total area covered. The second aspect of the coverage problem is to explore how to reduce redundancy in the coverage. The same area of a sensor field may be covered by many sensor nodes in a dense deployment. The aim is to ensure that no more than a few nodes cover a particular area at any time while the sensor field is fully covered, increasing sensing efficiency. In this paper, we review the proposed protocols from this sensing efficiency point of view.

We consider three main fields for redundancy elimination protocols. These fields are node scheduling schemes [1, 2, 3, 4, 5, 6, 7, 8, 9], approaches for solving the coverage problem (redundancy point of view) [10, 11, 12] and solutions for efficient routing [13, 14]. Although routing protocols proposed in [13, 14] are targeted for wireless ad-hoc networks. They have the potential to be deployed in wireless sensor networks.

Redundancy elimination protocols can also be classified into two major groups: location aware or location free schemes. The former utilizes the precise or relative

location information of sensor nodes. The latter does not consider any geographical position information of nodes.

In the following, these protocols are reviewed and a table is presented in the discussion section to summarize the attributes of each protocol.

2 Location Aware Protocols

Location aware schemes either utilize exact location information or obtain the relative position of nodes via some recursive trilateration or multilateration techniques [15]. For the first case, each node should be equipped with GPS (Global Positioning System) component or pre-programmed position information. The second case requires that some anchor nodes with GPS units are deployed in the sensor field.

Because of the knowledge of node's coordinates, these schemes can eliminate redundancy while preserving 100% network coverage. It means that there are no blind spots in the sensor field which is the main advantage of location aware schemes. Moreover, the warm-up stage of these techniques which organizes sensor nodes, usually results in lower overheads than location free schemes. Also, implementation and maintenance of these solutions are less complicated.

However, GPS units consume substantial amount of energy and so they decrease the total energy efficiency of redundancy elimination schemes. Additionally, these units increase the cost of wireless sensor nodes significantly and they do not work in all environments. In particular, they assume unobstructed line-of-sight communication with three or more geostationary satellites which is not available in many situations.

By exploiting range-based localization techniques [16, 17, 18], the relative location information of nodes are obtained. Nevertheless, we still need a large number of anchor nodes equipped with GPS units to reduce the accumulation of localization error due to iterative multilateration. Although range-free localization methods [19, 34, 20, 21] that utilize connectivity information do not require GPS units, they fail to achieve practical results with high level of accuracy [22].

2.1 GAF: Geographical Adaptive Fidelity

High redundancy in wireless networks results in several paths between source and sink. GAF [14] discovers a set of nodes that are equivalent in terms of routing. It then minimizes the number of paths by turning off the redundant intermediate nodes. A GAF simulation study confirms that in the presence of representative active nodes, a stable level of *routing fidelity* is preserved [14].

The first phase in GAF is to organize nodes in virtual grids. The monitoring area is divided into a number of grids. Each node is associated with one particular grid based on its location information. Nodes that are equivalent from the routing point of view are linked to the same grid. Two grids are adjacent when each node in one grid can communicate with every node in another grid. Grid size is determined based on the communication range of nodes.

After assigning a grid to each node, the main phase of GAF begins which determines and maintains the active node of each grid. In fact, each node has three different states. These are *sleeping*, *discovery* and *active*. Initially, all nodes are in the discovery mode. When the discovery timer of a node expires, it broadcasts a discovery message that

contains the grid id and node id of this node (with this approach equivalent nodes find each other). If a discovery node does not receive any response from the active node of its grid, it becomes the active node of this grid. Otherwise, it goes into the sleep state. The timer of sleeping nodes is set to a random number from the range of $[enlt/2, enlt]$. Here, $enlt$ is the expected lifetime of the active node. Upon the expiration of this timer, a sleeping node changes its state to discovery to detect the possible lack of an active node in its grid. GAF also provides fair energy consumption among sensor nodes. When a node becomes active, it sets a timer to a pre-defined active duration. When this timer fires, the active node returns to the discovery state. Now, discovery nodes that recently have become awake, have the chance to replace the active node of their grid.

GAF has some advantages. First, it sits over any underlying ad hoc routing protocol. Second, it dynamically changes the set of active nodes which provides load balancing. Third, it supports mobility of nodes. However, there are some drawbacks as well. Although it has been shown in [14] that the percentage of packet loss and route latency is as good as AODV and DSR routing protocols, it has been investigated under the relatively small number of traffic nodes. GAF does not propose any solution to preserve the capacity of network after eliminating the redundancy. Therefore, the performance of GAF can decrease significantly under high data traffic.

2.2 TTS: A Two-Tiered Scheduling Algorithm

The node scheduling presented in [3] poses a novel *two-tiered* node scheduling architecture. One layer is the *coverage-tier* which preserves full network coverage and the layer above is *connectivity-tier* that maintains a backbone to forward data traffic to the base station. This approach is targeted to put more nodes in sleep mode by considering different sleeping behaviors [3]. A lightweight localization method presented in [23] is utilized for determining the location of sensor nodes. A greedy algorithm is used to discover a set of sensor nodes for full coverage of the sensor field. The resulting cover set forms the coverage-tier. Subsequently, a dominating connectivity set is calculated so that each node in the cover set can directly communicate with at least one node of the connectivity set.

Node scheduling begins by considering three different patterns of activity. First, sensor nodes that neither exist in the coverage-tier nor in the connectivity-tier are put into sleep mode. Nodes in the coverage-tier wake up periodically to send measured data to the sink node. Sensor nodes in the connectivity-tier are always awake to forward data to the base station or receive a query from it.

There is an update interval that chooses new coverage and connectivity tiers dynamically to balance energy consumption of sensor nodes. At every update interval the algorithm is targeted to select sensor nodes so that the total remaining energy of every tier is maximized. A cost function is introduced that calculates the consumed power per coverage area of each sensor node. Therefore, when a sensor node consumes lower energy to monitor a larger area, it incurs lower cost. Thus, nodes with lower cost are selected for the next coverage and connectivity tier.

The main advantage of TTS is its tiered architecture. It is argued that having two different tiers for communication and coverage enables us to devise a wakeup and sleep pattern for covering nodes instead of keeping the nodes always awake. However, this architecture fails for the wireless sensor networks that need to monitor the region of interest constantly and in real time. Also, there is another drawback for TTS which

is its assumption of having the transmission range at least twice the sensing range. This assumption is necessary to determine the connectivity tier successfully [3].

2.3 K-Coverage Based on Perimeter-Coverage of Sensor Nodes

Authors in [10] propose a solution for the coverage problem. It determines whether the target area is adequately *k-covered* or not. Every point in the sensor field should be covered by at least *k* sensors. They further propose some applications such as node scheduling and finding insufficiently covered areas based on their approach.

The main idea of the solution is to determine whether the perimeter of each sensor is *k-covered* or not. First, they define that a sensor is *k-perimeter-covered* if every point on the perimeter of this sensor is covered by at least *k* sensors [10]. Then a low-cost algorithm is presented that discovers whether a particular sensor node is *k-perimeter-covered* or not. Also, another algorithm is proposed which works with sensor nodes with different sensing ranges. Finally, the perimeter-coverage of sensor nodes is related to the coverage property of the monitoring area. It has been proved that when all the sensor nodes are *k-perimeter-covered*, then the entire target area is *k-covered*.

In [10], several applications have been presented based on the proposed approach for the coverage problem. Determining insufficiently covered areas in the sensor field is one of the applications. Suppose *k* is the required value of coverage, first the base station broadcasts this value to all sensors. Then each sensor communicates with its neighbours to discover whether it has the segments in its perimeter which are less than *k-covered* or not. The base station collects the information about the insufficiently covered segments to find out the areas that are less than *k-covered*.

A basic node scheduling protocol is also proposed in [10]. In this approach, each sensor node S_i needs to communicate with its neighbours and ask them to reassess their perimeter coverage without node. If all the neighbours respond that still they are *k-perimeter-covered*, then node S_i can turn its radio off.

The proposed algorithm for coverage problem in this research work is quite efficient considering its cost which is $O(nd \log d)$, where *n* is the number of sensors and *d* indicates the highest number of nodes that might intersect a node. An interesting point about this scheme is its ability to work with irregular sensing ranges based on polygon approximation. However, its major drawback is its dependency on knowing the exact coordinates of sensor nodes.

2.4 Integrated Coverage and Connectivity

The basic aim in [5] is to maintain both sensing coverage and network connectivity as the result of its coverage maintenance protocol. A geometric analysis is presented to find out the relationship between coverage and connectivity [5]. It has been discovered that full network connectivity is preserved when sensing range is equal or less than half of the transmission range. As a result, a Coverage Configuration Protocol (CCP) is devised to eliminate redundancy based on a pre-defined coverage degree. If the coverage degree requested by the application is *k*, then the eligibility rule in CCP is as follows. If all the intersection points inside the sensing range of a particular sensor node are at least *k-covered*, then the sensor node is not entitled to be an active node. In order to calculate the intersection points inside the sensing circle, each node

should maintain a list of its active neighbours located within a distance of twice its sensing range. Three states have been considered for sensor nodes. These are *sleep*, *active* and *listen*. Initially all the nodes are in active mode. Then each node checks the eligibility rule to determine its next state. If the node becomes ineligible it goes to sleep mode. Then it wakes up periodically and enters into the listen state to re-evaluate its eligibility. Another solution is also considered in [5] to address the case when a node's sensing range is higher than half of the transmission range. SPAN [13] which is a redundancy eliminator protocol (reviewed later) is integrated with CCP. The former preserves the full connectivity while the latter maintains the requested coverage degree for the sensor field. The result of this integration is a new eligibility algorithm. With this new approach an inactive node executes the eligibility rules of both SPAN and CCP. If the result of one of them is positive then the inactive node is eligible to become active. Also, an active node will turn its radio off in case that it does not comply with the eligibility rule of either SPAN or CCP.

CCP provides very good flexibility in both preserving connectivity and addressing different levels of coverage. However, there are some shortcomings in this approach. The computation complexity of CCP is $O(n^3)$, where n is the number of active nodes within twice the sensing range of a node. This complexity is high for wireless sensor networks and it is not as efficient as other approaches. Moreover, this technique is infeasible for dense deployment of sensor nodes. Location error can also affect the performance of CCP. More importantly, the CCP eligibility rule assumes considering a perfect circular sensing model which is not true in the real world.

2.5 SET K-COVER

Slijepcevic *et al.* [12] focus on the set k -cover problem as a general form of the coverage problem. A heuristic is proposed to solve this problem by determining mutually exclusive sets of sensor nodes. These sets are called covers that monitor the entire area of the sensor field. The heuristic has been targeted to maximize the number of these covers. The proposed solution is based on dividing the deployment area into a number of *fields*. Each field is an area of the region of interest in which every point is covered by the same set of sensor nodes. Suppose A denotes the set of calculated fields. Then the set k -cover is formulated as follows. There are K disjoint covers, C_1, C_2, \dots, C_k , in which every member of A is covered by at least one element of every cover C_i [12]. The first part of the solution organizes points of the deployment area into a number of fields. A field is a unique area of the sensor field in which points in this area are covered by the same set of sensors. The second part of the solution utilizes a greedy algorithm to detect k covers. First, a field with lowest number of associated sensors is selected. An objective function is then utilized to calculate a probability value for every node in this field. This function measures the likelihood of redundant coverage for every sensor node if it is selected. Therefore, a node is selected which has the lowest likelihood. The selected field and its sensors are then ignored for the next repetition of the heuristic. This process is repeated until the current cover set attains full coverage. Then the members of calculated cover are deleted from the fields and the above process is repeated for the next cover set. The $K+1^{th}$ pass of the above heuristic which is the exit point, occurs when in the beginning of the process there is a field that has not been covered by any of the previously

calculated covers. Calculated covers are responsible to monitor the sensor field. At every time slice one cover is active while sensors belonging to the other covers are in sleep mode. The activity period is considered the same for all the covers. However, the authors have not implemented any scheduling scheme for covers.

The main advantage of this solution is its approach to organize the sensor network into mutually exclusive sets of nodes. However, the assumption of circular sensing area is not applicable in the real world. Therefore, the proposed scheme fails to deliver the promised full network coverage in real implementations.

2.6 Worst-Case & Best-Case Coverage

Research in [11] specifically manages the coverage problem in wireless sensor networks. Two general aspects of the coverage problem are considered. The *worst-case* coverage problem identifies the regions in the deployment area that have the lowest level of coverage compared to other areas. The *best-case* coverage problem identifies the regions that have the highest level of coverage by sensor nodes. For each case an algorithm is presented based on geometric computation and graph theory.

Worst-case coverage is identified by a path with pre-determined starting and end points in the sensor field. This path is called the *Maximal Breach Path*. The aim of the proposed algorithm is to find this path so that for every point on the path the distance of the point to the closest sensor node is maximized [11]. First, the Voronoi diagram [24] of the sensor field with its sensor nodes is generated. Then from this diagram a graph is created so that its nodes are the vertices of the diagram and its edges are the line segments of the diagram. A weight is also assigned to each edge that is the lowest distance from the closest sensor in the field. Finally, some graph search algorithms [11] are utilized to discover the path.

The best-case coverage is identified by a path whose starting point and end point are initially specified. This path is called the *Maximal Support Path*. The path is discovered in a way that the distance of every point on this path to the closest sensor is minimized. For this case another geometric structure is utilized that is Delaunay triangulation [25]. First, the Delaunay triangulation is computed based on the deployment area and its sensor nodes. Then a graph is formed such that its vertices are the sensor nodes and its edges are the lines of the triangles. The weight that is assigned to each edge is the length of its corresponding edge in the triangle. Finally, a search algorithm with specific search criteria [11] is applied to the graph to discover the maximal support path.

The most obvious advantages of the proposed scheme are the applications that can be derived from it. For example, the network administrator is able to discover the paths in sensor field which suffer from poor coverage. Then, coverage can be easily improved by deployment of new nodes in those areas. Additionally, by discovering maximal support paths, it is possible to schedule the redundant nodes of these paths to sleep. Therefore, the discovery of breach and support paths leads to a balanced network coverage. There are also some drawbacks in the proposed scheme. It only considers homogeneous sensor nodes. Moreover, the accuracy of the discovered paths is strongly related to the precision of location information of sensor nodes. Therefore, the localization techniques that rely on a few GPS devices can affect the quality of maximal breach and support paths.

2.7 Node Scheduling by Sponsoring Nodes

In [4] a node scheduling scheme is proposed in which every node is autonomously responsible for checking an off-duty eligibility rule. This rule considers the neighbours of a node to determine whether they can sponsor the node by covering its sensing area. The proposed solution comprises of two rounds. Round one is a self-scheduling process that is followed by a typical sensing round. These rounds are repeated periodically.

In the beginning of the self-scheduling round, each node is responsible for discovering its neighbours within its sensing range. This discovery includes location information for these nodes. Then, the node determines its eligibility to be an inactive node by performing some geometric calculations. These calculations include the measurement of the central angles of the sectors that are formed based on intersection of this node by each of its neighbours.

Every node is responsible for calculating all the central angles that exist between the node and its neighbours. If the summation of all these angles covers 360° , the sensing area of the node is fully sponsored by its neighbours and it is eligible to turn off its radio. In order to avoid any conflict that may occur if nodes sponsor each other concurrently, a standard back-off technique is utilized. Also, the connectivity between on-duty nodes is provided by LEACH [26], a cluster-based protocol.

It has been proved in [4] that the requirement of location information can be replaced by direction information in order to determine sponsor nodes. Additionally, the off-duty eligibility rule has been further extended to respond to the case that nodes have different sensing ranges. Nevertheless, there are some major shortcomings. It has been assumed that the sensing round should be long enough in order to compensate for the overhead of the self-scheduling round. This assumption can lead to large sensing and communication gaps in case of node failures. Moreover, in case of dense deployments, there is no approach to tackle the packet collisions that might occur during neighbourhood discovery. Loss of accuracy of neighbourhood information leads to higher numbers of on-duty nodes. The last disadvantage is the dependency of the off-duty eligibility rule on the existence of at least three neighbours for every sensor node.

2.8 Coverage Preserving Redundancy Elimination

Detecting and eliminating redundant sensor nodes while preserving complete network coverage is the main target of the solution presented in [1, 2]. Moreover, it presents an approach to discover the *coverage-boundary* of a sensor field. The proposed methods are based on Voronoi diagrams generated locally at each sensor node.

A redundant sensor is defined as a node whose sensing disk is completely under coverage of other nodes. This approach utilizes two different Voronoi diagrams. One is the Voronoi diagram of the network. The second one, constructed for every sensor node, is the Voronoi diagram of the Voronoi neighbours of that node, called the 2-Voronoi diagram. The vertices of this diagram are called 2-Voronoi Vertices (2-VV). In addition, the intersection point between the coverage circumcircle of the sensor node and an edge of its 2-Voronoi diagram is calculated. It is called a 2-Voronoi Intersection Point (2-VIP). Then, every sensor node is responsible for discovering whether it is redundant or not by utilizing this information. It has been proved in [1, 2] that a sensor node is redundant if its 2-VVs and 2-VIPs are under the coverage of its Voronoi neighbours. Every sensor node is responsible for checking this rule to determine whether it is a

redundant node or not. Redundant nodes that are Voronoi neighbours of each other should not decide to turn their radio off simultaneously. Since that could cause blind spots in the sensor field. In order to prevent it, an approach based on discovering a maximal independent set [27] of redundant nodes is implemented.

The second contribution of the research work presented in [1, 2] is dedicated to the detection of the coverage-boundary of the deployment area. It is formulated that the coverage-boundary includes the sensing area of those sensor nodes whose coverage areas are not fully covered by sensing disks of other nodes. It has been proved that a node is on the boundary area if its Voronoi cell is not covered by its sensing disk [1,2].

The main advantage of the proposed scheme is its ability to address the heterogeneous sensor networks. It is achieved by a modified version of Voronoi diagrams. Additionally, the proposed scheme utilizes a distributed approach compared to the centralized Voronoi calculations in [11]. It is capable of recomputing the Voronoi information locally in order to address node failures and deployment of new nodes. As a result, it is much more robust in this regard compared to [11]. Nevertheless, this protocol imposes a large amount of overhead for two reasons. The first issue is the regular maintenance of Voronoi cells. The second factor is the initial overhead imposed when every sensor node is supposed to find its Voronoi neighbours and compute its Voronoi cell. Also, similar to [11], the performance of the scheme is sensitive to the accuracy of the location information of each node. This is because of the dependency of these approaches to geometrical calculations.

3 Location Free Protocols

In location free schemes for node scheduling, nodes discover the network topology and react to its changes by sending control messages. Due to the lack of absolute position information, the main challenge for these techniques is to maintain a high level of network coverage after eliminating the redundant nodes. One of the most important metrics to evaluate the efficiency of these approaches is the number of active nodes. Obviously, the quality of coverage can be increased by a simple but wasteful decision of eliminating fewer nodes. However, the issue that determines the superiority of a proposed scheme is how to have as few as possible active nodes while preserving a high level of network coverage.

Location free approaches are more cost-effective and more energy efficient than location aware protocols because they do not require GPS devices. Therefore, they are better suited for wireless sensor networks since the major constraint is energy. Location independent schemes fail to ensure 100% network coverage and may introduce high initial packet and energy overhead in the network. However, they extend the lifetime of network considerably without additional cost of special hardware components.

3.1 SPAN: An Energy-Efficient Coordination Algorithm for Topology Maintenance

Chen *et al.* [13] propose a protocol to save energy in ad hoc wireless networks by keeping only the necessary nodes active to forward routing traffic. In SPAN each node periodically decides whether it is eligible to become a coordinator node and join the forwarding backbone or it should turn its radio off.

There is a basic rule which determines the eligibility of a node to be a coordinator node. The rule is that if a node finds a pair of its neighbours that cannot reach each other directly or through other coordinator nodes, then the node should be considered as a coordinator node. However, complying with the eligibility rule does not mean that the node will surely become an active node because there is the possibility that multiple nodes in a particular region of the network become eligible at the same time. This issue is addressed by using random backoff delay for every node that wants to announce its coordination role. The delay is based on two parameters. The first factor specifies the number of neighbours that will be connected if this node becomes a coordinator. The more neighbours get connected, the higher is the priority of this eligible node. The second factor is based on the amount of remaining energy of the node, higher priority is given to the node with higher amount of remaining energy.

Each node should announce its willingness to become a coordinator by sending a HELLO message after the delay time calculated based on the above factors. Therefore, at the end of the delay time, if a node has not received any HELLO message, it will broadcast a HELLO packet. Otherwise, the node should reassess its eligibility based on new conditions of the network. If still the node is eligible it repeats the above process. Otherwise it will sleep until the next election time.

One advantage of SPAN is its approach to check coordinator nodes regularly to keep their number as low as possible. A coordinator will be withdrawn if it discovers that every pair of its neighbours is connected directly or via other coordinators. Another useful feature of SPAN is its ability to support load-balancing in the network. A coordinator will be withdrawn after a period of time, if there is a connection between every pair of its neighbours directly or via some other neighbours.

In SPAN when a node is asleep the packets that were destined to it will be buffered at its coordinator. This is achieved based on power-saving feature of 802.11 MAC layer [28]. The main drawback of SPAN is its dependency on a power-saving feature of the MAC layer protocol. SPAN will not work with MAC layers that do not provide this option. Moreover, its performance will be affected according to the challenges that are related to buffering.

3.2 LUC: Location-Unaware Coverage

Younis *et al.* [9] propose a solution for node scheduling that relies on four different conditions. Two tests (RTest-D1 and RTest-D2) are geometrical. They determine that a node is redundant if its sensing area is covered geometrically by other active nodes. Two other tests (RTest-H1 and RTest-H2) specify the redundancy of a node if some conditions based on node density are met [9].

RTest-D1: A node v is redundant if these conditions are satisfied. 1) There are nodes v_j , $1 \leq j \leq 3$ that lie in the sensing range of node v . 2) v_j nodes are active. 3) v_j nodes are pairwise neighbours inside the sensing range. 4) Node v is placed inside the triangle formed by v_j nodes. 5) Boundary of node v is covered by the boundary of v_j nodes.

RTest-D2: A node v is redundant if these conditions are satisfied. 1) There are nodes v_j , $1 \leq j \leq 3$ that lie in the range $0.618 * R_s$. R_s is the sensing range of node v . (for more details please refer to [9]). 2) v_j nodes are active. 3) v_j nodes are pairwise non-neighbours within the sensing range of node v .

RTest-H1: A node v is redundant if these conditions are satisfied. 1) There are active neighbours v_j , $j \geq 4$, within distance R_s . 2) Every neighbour of node v is a neighbour of one or more of v_j nodes.

RTest-H2: First these conditions are considered. 1) Node v has at least y active or undecided neighbours inside its sensing range. 2) Its other neighbours have one or more neighbours among those y neighbours. Node v is considered redundant if it has the least remaining energy among those y active neighbours. In [9] the authors assumed that y is 6.

Each node has three states: *active*, *sleep* and *undecided*. A Location-Unaware Coverage (LUC) algorithm utilizes rules D1, D2, H1 and H2 to determine the state of each node in network. The first phase is started by a neighbourhood discovery when each node learns its one-hop and two-hop neighbours as well as its distance to its one-hop neighbours. Then the operational phase begins that determines the active and asleep nodes. Whenever a node notices that one of its neighbours has become an active node, it performs a redundancy check procedure to decide whether it is redundant or not. This procedure performs the tests in a special order. The order is the same as we listed above. If the result of any test determines that the node is redundant, it will quit the procedure and will change its state to sleep. Otherwise, the procedure is triggered by the next active neighbour. However, if a specific timer expires and the node still remains in an undecided state, the last round of the redundancy checking procedure is performed. Nodes become active if all the tests have a negative result.

There are some shortcomings of this scheme. First and foremost, it has been shown via the simulations in [9] that the total success ratio of the first two tests is almost 60%. Since the next two tests are exploited for dense networks, the proposed approach can not detect all the redundant nodes for low density networks. Moreover, the algorithm is limited by considering the assumption that the transmission range should be at least two times the sensing range. Last but not least, nodes should be equipped with special hardware to measure the distance between neighbours.

3.3 Expectancy of Coverage

The proposed solution in [8] comprises of two parts. First, the authors present some theoretical analysis to estimate the coverage level of the monitoring area based on some parameters of the wireless sensor network. Second, they devise a location-free node scheduling scheme.

In order to formulate the coverage expectancy, three factors are considered. These are sensing range, number of deployed nodes and area of the sensor field. First the expected coverage of each sensor node is predicted and then the expected coverage ratio of the entire deployment area is discovered as follows: (more details in [8])

$$E[C_n] = \left[1 - \left(1 - \frac{\frac{1}{2}r^4 - \frac{4}{3}lr^3 - \frac{4}{3}mr^3 + \pi r^2 ml}{m^2 l^2} \right)^n \right] lm \quad (1)$$

Here, $E[C_n]$ is the expected coverage level, r is the sensing range, n is the number of deployed nodes and the deployment area is represented by an l by m rectangle. Also,

this expression can be utilized to determine the minimum number of sensor nodes needed to maintain a pre-defined coverage ratio requested by the application.

This analysis for coverage estimation is further utilized to present a node scheduling scheme. The proposed solution is an uncoordinated scheme in which each node stochastically decides to become active or asleep. This decision is repeated independently. The duration of time that nodes remain in active or sleep state is based on an exponential distribution with parameters λ_a and λ_s , respectively [8]. If we consider that a particular node has changed its state m times. Then the probability that the node would be active is:

$$P_a = \frac{m\lambda_a}{m\lambda_a + m\lambda_s} \quad (2)$$

Then it can be estimated that the number of active nodes at any given time is nP_a . Now if we put the estimated number of active nodes in equation 1, the expected coverage of the network can be calculated. This calculation can help us to adjust λ_a and λ_s to attain the coverage ratio requested by an application.

Although, there are some basic analyses [35, 36] in the literature for estimating the coverage, this research takes into account the *border effect* to present a more accurate approximation. Border effect refers to nodes that are located close to the border of the monitoring area. Moreover, the proposed node scheduling scheme is an uncoordinated approach with very low overhead. Nevertheless, the scheduling process is completely stochastic which might not possess a stable behaviour for some issues such as preserving the network capacity. In addition, the proposed solution for coverage expectancy neither addresses sensor nodes with different sensing ranges nor real-world irregularity of sensing area.

3.4 PEAS: Probing Environment and Adaptive Sleeping

PEAS [6, 7] is a distributed redundancy eliminator protocol for wireless sensor networks. It has been devised to prolong the network lifetime by detecting a group of working nodes while redundant nodes are put into sleep mode.

PEAS comprises of two main components, *probing environment* and *adaptive sleeping*. The former is responsible for maintaining the necessary set of working nodes. The latter is responsible for adjusting the sleeping time of inactive nodes. Three different states have been considered for sensors; sleeping, probing and working. Initially all the nodes are in sleeping mode. When the sleep timer of a node expires, it wakes up and broadcasts a PROBE message based on a pre-defined probe range and waits to receive a REPLY message from each of the neighbouring working nodes. If the probing node does not receive any REPLY message, it becomes a new member of the group of working nodes. Otherwise it goes back to the sleep mode. A probe rate λ is considered for probing nodes.

The most remarkable aspect of PEAS is its feature to adjust the frequency of probing actions by an adaptive sleeping component. The goal is that every working node gets probed based on a desired probe rate, λ_d . In order to achieve this goal, every working node is responsible for measuring the aggregate probe rate, $\bar{\lambda}$, that it has received from its neighbouring probe nodes. Then, probing nodes use the aggregate probe rate to calculate a new probe rate as follows.

$$\lambda^{new} = \lambda \frac{\lambda_d}{\lambda} \quad (3)$$

Subsequently, a new sleeping time is generated according to the equation 4. Where, t_s is the previous sleeping time.

$$f(t_s) = \lambda^{new} e^{-\lambda^{new} t_s} \quad (4)$$

It has been proven in [6, 7] that with this kind of adjustment of probe rate as well as exponentially distributed sleeping times, the probe rate that every working node receives is equal to the desired probe rate.

PEAS is a location-free solution for the coverage problem that does not rely on distance information. Moreover, by simply varying the probe range it can provide different coverage levels based on application requirements. Additionally, its proposed approach of adaptive sleeping can address node failures quite robustly while keeping the overhead of this action low. Due to the high efficiency of PEAS in terms of handling node failures, we have utilized it as the fault-tolerant module of our node scheduling scheme. However, it needs to be mentioned that the main disadvantage of PEAS is its quality in terms of *temporal coverage*. Temporal coverage indicates the coverage ratio of the sensor field throughout the network lifetime.

3.5 Triangular: Determining Active Sensor Nodes for Complete Coverage by a Triangular Self-test Mechanism

Wi *et al.* [29] propose a location-free approach for selecting a set of active sensor nodes. The main contribution of this research is that in the absence of node coordinates, it can achieve 100% network coverage. However, it assumes that nodes are equipped with omni-directional antennas to measure the distance between nodes by received signal strength. The proposed solution comprises of two phases. In phase one the initial set of active nodes is determined. The second phase is dedicated to detecting blind spots in the sensor field.

First, all the nodes are inactive and each of them runs a back-off timer. Upon the expiration of this timer, a node wakes up and checks whether there is a working node within its sensing range or not. If the discovery is not successful, the node becomes a new working node. Otherwise, the node remains awake for a period of time to build a *Neighbour Table* based on some control messages. This table includes node id and distances of all the one-hop and two-hop neighbours. It is utilized by the second phase for detection of sensing holes. When this phase is finished an initial set of active nodes has been discovered.

In the second phase, a *triangular self-test* is performed by each inactive node to eliminate any possible blind spots. For every inactive node a triangle is formed based on the three closest active neighbours. Then relative coordinates are assigned to the endpoints of the triangle followed by the calculation of the relative coordinates of the circumcentre of the triangle. Based on calculating the distance from circumcentre to the endpoints of the triangle and comparing it with the sensing range, the existence of a blind spot is determined. If the result of this self checking approach is positive the inactive node becomes active to fill the sensing holes.

There are two main advantages of this solution. Firstly, a relative coordinate is assigned to each node based on virtual coordinates of its neighbours. This position information is not used for other location calculations. Therefore, the scheme is not affected by propagation of location errors. Secondly, the simulation results presented in [29] indicate that the proposed scheme is quite scalable so that the packet overhead does not change when node density increases. However, there are some drawbacks as well. Although the simulation results show that the set of active nodes covers the deployment area completely, the scheme achieves that by eliminating the margins of the sensor field. Also, there are some real-world factors that affect the efficiency of distance measurement based on received signal strength. This issue decreases the performance of the triangular self-test mechanism. Moreover, the proposed solution assumes that the transmission range should be at least three times the sensing range which limits the applicability of the solution.

4 Discussion

The first and foremost characteristic of a redundancy elimination protocol is its approach regarding location information. In sections 2 and 3 we discussed the advantages and disadvantages of location-free and location-aware solutions. However, an important point that needs to be mentioned is that some location-free approaches [29, 9] consider distance information. There are three main drawbacks in this regard. First, we need special hardware to measure distance. This increases the cost of the wireless sensor network. Second, environmental conditions affect the accuracy of distance estimation which might decrease the quality of coverage. The third drawback is related to the techniques that are used for measuring the distance between nodes. The RF signal strength approaches might have low performance when the transmission range is short. Fading, multipath and interferences are the factors that affect the signal strength when the radio range is short [30]. Moreover, time of flight approaches [31, 32, 33] might deliver inaccurate results mostly because of clock skew and drift between nodes.

Protocols [10, 12, 4, 5, 29, 9] determine redundant nodes based on some calculations on either the area of the sensing circle or the perimeter of these circles. These approaches might not deliver what they have promised to achieve. This is due to the assumption of perfect circular sensing areas. In the real-world, sensing is not only irregular but it might vary over time. Although other protocols consider circular sensing areas, they use it only for calculation of the coverage ratio and not as a building block of their schemes. Therefore, the only shortcoming of these approaches is some degrees of error in measuring their coverage ratio metric.

All the solutions that we have reviewed in this chapter utilize the perfect circular radio propagation model. This model does not reflect the real-life behaviour of the wireless channel. Therefore, the performance of these protocols when they are implemented can be different from what we have seen through simulation results. Although GAF [14] has investigated the effect of the shadowing radio model, this protocol is originally proposed for unrealistic radio models.

Elimination of redundant nodes might result in an unconnected set of representative nodes. In order to tackle this issue, some approaches [3, 29, 9] consider a special assumption that the transmission range should be at least twice the sensing range.

Therefore, they can not address particular applications where the communication range might be even lower than sensing range.

Finally, Table 1 compares the main characteristics of the protocols that we have reviewed in this paper.

Table 1. Comparison of redundancy eliminatro protocols

Redundancy Eliminator Protocol	Location Free	Location Aware	Distance Information	Distributed	Centralized	Coordinated	Uncoordinated	Use of Unrealistic Radio Model
GAF [14]	-	✓	-	✓	-	✓	-	¹
TTS [3]	-	✓	-	-	✓	-	✓	✓
K-Coverage based on Perimeter-Coverage of Sensor Nodes[10]	-	✓	-	✓	-	✓	-	✓
Integrated Coverage and Connectivity [5]	-	✓	-	✓	-	✓	-	✓
SET K-COVER [12]	-	✓	-	-	✓		✓	✓
Worst-Case & Best-Case Coverage [11]	-	✓	-	-	✓	N/A	-	✓
Node Scheduling by Sponsoring Nodes [4]	-	✓	-	✓	-	✓	-	✓
Coverage Preserving Redundancy Elimination [1, 2]	-	✓	-	✓	-	✓	-	✓
SPAN [13]	✓	-	-	✓	-	✓	-	✓
LUC [9]	✓	-	✓	✓	-	✓	-	✓
Expectancy of Coverage [8]	✓	-	-	✓	-		✓	✓
PEAS [6,7]	✓	-	-	✓	-	✓	-	✓
Triangular [29]	✓	-	✓	✓	-	✓	-	✓

¹ GAF is originally proposed for unrealistic radio models. However it has also been investigated under shadowing model.

References

1. Carbutar, B., Grama, A., Vlitek, J., Carbutar, O.: Redundancy and Coverage Detection in Sensor Networks. *ACM Transactions on Sensor Networks* 2(1), 94–128 (2006)
2. Carbutar, B., Grama, A., Vitek, J., Carbutar, O.: Coverage Preserving Redundancy Elimination in Sensor Networks. In: *Proceedings of the First IEEE Communications Society Conference on Sensor and Ad Hoc Communications and Networks (SECON 2004)* (October 2004)
3. Tezcan, N., Wang, W.: TTS: A Two-Tiered Scheduling Algorithm for Wireless Sensor Networks. In: *Proceedings of IEEE International Conference on Communications (ICC 2006)* (June 2006)
4. Tian, D., Georganas, N.D.: A Coverage-Preserving Node Scheduling Scheme for Large Wireless Sensor Networks. In: *Proceedings of the 1st ACM International Workshop on Wireless Sensor Networks and Applications (WSNA 2002)* (September 2002)
5. Wang, X., Xing, G., Zhang, Y., Lu, C., Pless, R., Gill, C.: Integrated Coverage and Connectivity Configuration in Wireless Sensor Networks. In: *Proceedings of the 1st International Conference on Embedded Networked Sensor Systems (SenSys 2003)* (November 2003)
6. Ye, F., Zhang, H., Lu, S., Zhang, L., Hou, J.: A Randomized Energy-Conservation Protocol for Resilient Sensor Networks. *Wireless Networks* 12(5), 637–652 (2006)
7. Ye, F., Zhong, G., Cheng, J., Lu, S., Zhang, L.: PEAS: A Robust Energy Conserving Protocol for Long-lived Sensor Networks. In: *Proceedings of the 23rd International Conference on Distributed Computing Systems, ICDCS 2003* (May 2003)
8. Yen, L.-H., Yu, C.W., Cheng, Y.-M.: Expected K-Coverage in Wireless Sensor Networks. *Ad Hoc Networks* 4(5), 636–650 (2005)
9. Younis, O., Krunz, M., Ramasubramanian, S.: Coverage Without Location Information. In: *Proceedings of the IEEE International Conference on Network Protocols (ICNP 2007)* (October 2007)
10. Huang, C.-F., Tseng, Y.-C.: The Coverage Problem in a Wireless Sensor Network. In: *Proceedings of the 2nd ACM International Conference on Wireless Sensor Networks and Applications (WSNA 2003)* (September 2003)
11. Meguerdichian, S., Koushanfar, F., Potkonjak, M., Srivastava, M.: Coverage Problems in Wireless Ad-Hoc Sensor Networks. In: *Proceedings of the twentieth Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM 2001)* (April 2001)
12. Slijepcevic, S., Potkonjak, M.: Power Efficient Organization of Wireless Sensor Networks. In: *Proceedings of the IEEE International Conference on Communications (ICC 2001)* (June 2001)
13. Chen, B., Jamieson, K., Balakrishnan, H., Morris, R.: Span: An Energy-Efficient Coordination Algorithm for Topology Maintenance in Ad Hoc Wireless Networks. *Wireless Networks* 8(5), 481–494 (2002)
14. Xu, Y., Heidemann, J., Estrin, D.: Geography-Informed Energy Conservation for Ad Hoc Routing. In: *Proceedings of the 7th Annual International Conference on Mobile Computing and Networking (MobiCom 2001)* (July 2001)
15. Bulusu, N., Estrin, D., Girod, L., Heidemann, J.: Scalable Coordination for Wireless Sensor Networks: Self-Configuring Localization Systems. In: *Proceedings of the 6th International Symposium on Communication Theory and Applications (ISCTA 2001)* (July 2001)
16. Priyantha, N.B., Chakraborty, A., Balakrishnan, H.: The Cricket Location-Support System. In: *Proceedings of the 6th Annual International Conference on Mobile Computing and Networking (MobiCom 2000)* (August 2000)
17. Savvides, A., Han, C.-C., Srivastava, M.B.: Dynamic Fine-Grained Localization in Ad-Hoc Networks of Sensors. In: *Proceedings of the 7th Annual International Conference on Mobile Computing and Networking (MobiCom 2001)* (July 2001)

18. Whitehouse, C.D.: The Design of Calamari: An Ad-hoc Localization System for Sensor Networks. Master Thesis, University of California at Berkeley (2002)
19. He, T., Huang, C., Blum, B.M., Stankovic, J.A., Abdelzaher, T.: Range-Free Localization Schemes for Large Scale Sensor Networks. In: Proceedings of the 9th Annual International Conference on Mobile Computing and Networking (MobiCom 2003) (September 2003)
20. Moore, D., Leonard, J., Rus, D., Teller, S.: Robust Distributed Network Localization with Noisy Range Measurements. In: Proceedings of the 2nd International Conference on Embedded Networked Sensor Systems (SenSys 2004) (November 2004)
21. Niculescu, D., Nath, B.: Ad Hoc Positioning System (APS). In: Proceedings of Global Telecommunications Conference (GLOBECOM 2001) (November 2001)
22. Zhong, Z., He, T.: MSP: Multi-Sequence Positioning of Wireless Sensor Nodes. In: Proceedings of the 5th International Conference on Embedded Networked Sensor Systems (SenSys 2007) (November 2007)
23. Cheng, X., Thaler, A., Xue, G., Chen, D.: TPS: A Time-Based Positioning Scheme for Outdoor Wireless Sensor Networks. In: Proceedings of the Twenty-third Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM 2004) (March 2004)
24. Aurenhammer, F.: Voronoi Diagrams - A Survey of a Fundamental Geometric Data Structure. *ACM Computing Surveys* 23(3), 345–405 (1991)
25. Preparata, F.P., Shamos, M.I.: *Computational Geometry: An Introduction*. Springer-Verlag New York, Inc., Heidelberg (1985)
26. Heinzelman, W.R., Chandrakasan, A., Balakrishnan, H.: Energy-Efficient Communication Protocol for Wireless Microsensor Networks. In: Proceedings of the 33rd Hawaii International Conference on System Sciences (HICSS 2000) (January 2000)
27. Luby, M.: A Simple Parallel Algorithm for the Maximal Independent Set Problem. In: Proceedings of the 7th Annual ACM Symposium on Theory of Computing (STOC 1985) (May 1985)
28. Wireless LAN Medium Access Control and Physical Layer Specifications. IEEE 802.11 Standard (IEEE Computer Society LAN MAN Standards Committee) (August 1999)
29. Wi, T.-T., Ssu, K.-F.: Determining Active Sensor Nodes for Complete Coverage without Location Information. *International Journal of Ad Hoc and Ubiquitous Computing* 1(1-2), 38–46 (2005)
30. Bulusu, N., Heidemann, J., Estrin, D.: GPS-less Low Cost Outdoor Localization for Very Small Devices. *IEEE Personal Communications* 7(5), 28–34 (2000)
31. Girod, L., Lukac, M., Trifa, V., Estrin, D.: The Design and Implementation of a Self-Calibrating Distributed Acoustic Sensing Platform. In: Proceedings of the 4th International Conference on Embedded Networked Sensor Systems (SenSys 2006) (November 2006)
32. Whitehouse, K., Culler, D.: Calibration as Parameter Estimation in Sensor Networks. In: Proceedings of the 1st ACM International Workshop on Wireless Sensor Networks and Applications (WSNA 2002) (September 2002)
33. Youssef, M., Youssef, A., Rieger, C., Shankar, U., Agrawala, A.: PinPoint: An Asynchronous Time-Based Location Determination System. In: Proceedings of the 4th International Conference on Mobile Systems, Applications and Services (June 2006)
34. Lazos, L., Poovendran, R.: SeRLoc: Secure Range-Independent Localization for Wireless Sensor Networks. In: Proceedings of the 3rd ACM Workshop on Wireless Security (WiSe 2004) (October 2004)
35. Hall, P.: *Introduction to the Theory of Coverage Processes*. John Wiley and Sons, Chichester (1988)
36. Philips, T.K., Panwar, S.S., Tantawi, A.N.: Connectivity Properties of a Packet Radio Network Model. *IEEE Transactions on Information Theory* 35(5), 1044–1047 (1989)

A Light-Weighted Misused Key Detection in Wireless Sensor Networks

Young-Ju Han¹, Min-Woo Park¹, Jong-Myoung Kim², and Tai-Myoung Chung¹

¹ Internet Management Technology Laboratory,
Department of Electrical and Computer Engineering, Sungkyunkwan University,
300 Cheoncheon-dong, Jangan-gu, Suwon-si, Gyeonggi-do, 440-746, Republic of Korea
Tel.: +82-31-290-7222; Fax: +82-31-299-6673
{yjhan,mwpark,tmchung}@imt1.skku.ac.kr

² Incidents Analysis Team, Korea Internet Security Agency, Republic of Korea
jmkim@kisa.or.kr

Abstract. In wireless sensor networks, sensor nodes are generally distributed in hostile environments so the security services such as confidentiality, authentication and integrity are very important. The basis of these security services is the key management. The majority of key managements use the random key pre-distribution mechanism based on the probability model. Thus, an adversary can easily get the keys between any two non-compromised nodes by compromising some number of nodes. If the adversary gets the keys, he can modify or insert messages into the communication channels. These kinds of attacks are very critical because the adversary may mislead the operation of the networks or the application. Recently, Liu and Dong proposed a method to detect the keys which is misused by an adversary. They suggested the additional protection for non-compromised sensor nodes even if an attacker has learned the shared key between them. To detect a misused key, the traffic overhead is inevitable but their method has some inefficient aspects by performing a unidirectional misused key detection and generating unnecessary traffic. In this paper, we enhance the misused key detecting mechanism by detecting a misused key bidirectionally and removing the unnecessary traffic while the security functionality of the Liu and Dong's scheme is preserved. Our simulation shows that the energy consumption of our scheme is 36% more efficient than that of the Liu and Dong's scheme on average.

1 Introduction

Recent advances in wireless communication, micro system and low-power technologies have enabled the wireless sensor networks(WSNs). Generally, WSNs consist of many, low-cost, low-power and small sensor nodes and the base station(BS). Once the sensor nodes are deployed at a target area, they collect some information from the area and report it to the BS. The BS takes the role of gateway to the traditional networks(:Internet), so we can use the collected information for specific applications. For example, some applications of WSNs

are habitat monitoring, object tracking, fire detection, traffic monitoring and so on [1,2,7].

The security services in WSNs are very important because the sensor nodes are often deployed in hostile environments. An adversary can easily overhear the radio channels, insert or modify some messages, or capture physically some sensor nodes. To guarantee security service such as the integrity, authentication and confidentiality in the WSNs, cryptographic methods are essential and the key management is the most fundamental one. However, the traditional key management methods such as KDC(Key Distribution Center) and public key mechanism are the impractical solutions due to the limited resources of sensor nodes [14].

Many key management techniques in WSNs have been introduced [4,8,9,10,13]. Eschenauer and Gligor proposed the basic random key distribution scheme based on the probability model to overcome the storage limitation of the sensor node for storing keys [8]. The main idea is that k keys are randomly selected by each node out of a large pool of P keys. However, it has the critical drawback that the attacker can easily threat the network survivability by capturing some nodes. Chan suggested the q -composite random key distribution scheme [6]. This scheme makes the resilience against node capture more strong but makes the key sharing probability between any pair of nodes decreased. Liu and Ning proposed more enhanced key management scheme using t -degree bivariate polynomials and two-dimensional grid [3]. It perfectly provides the network resilience until $t + 1$ nodes are captured for any specific keys and reduces the overhead on establishing pairwise keys. An attacker can get the keys between any two non-compromised nodes because the all of the key management schemes mentioned above are commonly based on the probability model. If An adversary gets some keys using physical node capture attack, the attacker might infer some keys between non-compromised nodes by using the collected keys from compromised nodes, and then he can attempt some kinds of attacks such as eavesdropping, insertion and modification of messages. The insertion and modification attacks are more harmful than the eavesdropping because they can lead the abnormal operations of the applications in WSNs. Figure 1 describes this situation. Node 3, 6 and

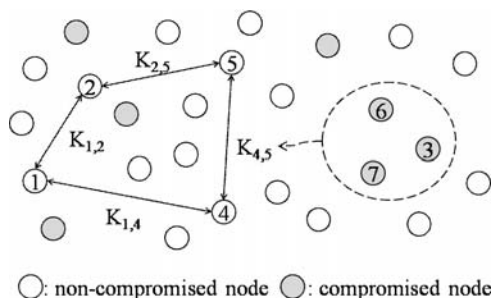


Fig. 1. An example of misused keys

7 were compromised by the attacker and the attacker obtained the shared key, $K_{4,5}$, between node 4 and 5. He can insert or modify some messages and finally threaten the overall functionality of an application.

Recently, Liu and Dong has proposed a misused key detection scheme in WSNs to prevent these kinds of attacks [5]. According to them, the misused key can be defined under the following three conditions.

- It should be the key shared between non-compromised nodes.
- The key is exposed to an adversary.
- The key must be used for insertion and modification attacks.

To detect a misused key efficiently in distributed manner, they has introduced committing values and detecting nodes. Whenever a sender node transfers a message to a receiver node, the sender node has to attach committing values to the message and the receiver node performs the detecting procedure by generating a report message using the committing values under the predefined probability. Finally, the detecting nodes and the BS check the report message to determine whether the key was misused.

However, Liu and Dong's scheme has four limitations. First, the unused committing values at the receiver node are big overhead. Second, the reporting procedure can detect only the unidirectional key misuse from a sender to a receiver per one report message. Third, it is impractical that the detecting node lists must be exchanged between sensor nodes because the pairwise key establishment depends on the routing protocol deployed for the purpose of efficient operation of specific applications in WSNs. Additionally, the detecting procedure in Liu and Dong's scheme occurs unbalanced traffic overhead among sensor nodes.

In this paper, we introduce a light-weighted misused key detection mechanism based on Liu and Dong's scheme. Our scheme has better two properties, which are the reduction of the detecting overhead and the high detecting probability, than Liu and Dong's scheme by performing the bidirectional misused key detection, removing of unnecessary message and using the reasonable probability method for the balanced misused detection among the sensor nodes.

The rest of this paper is organized as follows. The next section briefly reviews the Liu and Dong's scheme and discusses about its limitations. Section 3 introduces our scheme: a light-weighted misused key detection mechanism. In section 4, we will analyze our mechanism from the viewpoints of security and energy efficiency compared with the Liu and Dong's method. Finally, section 5 summarizes our paper and discusses about future research.

2 The Liu and Dong's Work

An adversary can compromise easily a few sensor nodes and may obtain many shared keys between non-compromised nodes. By using these exposed keys, the adversary can insert forged or modified messages into the networks. The application of WSNs might be misled by the adversary since the modified or forged

messages make the application choice wrong decisions. To prevent these kinds of attacks, Liu and Dong suggests the additional protection for non-compromised sensor nodes even if an attacker has learned the shared keys between them.

2.1 The Assumptions

Before we summarize their method, we briefly review their assumptions as follows.

- The adversary can eavesdrop, modify, replay, forge or block any network traffic.
- It is computationally impossible to break the cryptographic primitives.
- The adversary can compromise a few nodes and obtain some shared keys between non-compromised nodes.

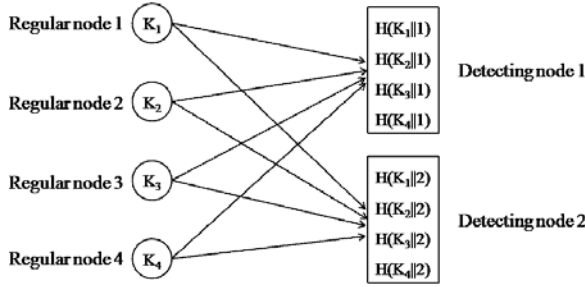
2.2 The Ideas

The main idea of this scheme is based on the committing value(CV) and the secret key between a sensor node and the BS. The CV is obtained from the cumulative hashed value about outgoing messages(for a sender) or incoming messages(for a receiver). If node u and v has established a pairwise key $K_{u,v}$, each node maintains a variable $C_{u,v}$ to track the history of the communication from u to v ($C_{v,u}$ is also maintained for the inverse direction). The initial value of the $C_{u,v}$ is 0. If the node v receives the un-forged, un-modified or un-inserted messages, the $C_{u,v}$ of u would be the same as that of v . If the two nodes want to check whether their pairwise key was misused, they simply calculate the CVs using the secret key shared with the BS and send it to the BS. Finally, the BS checks whether the key was misused and updates the key if the key was misused. In addition to this basic idea, they introduced the detecting node, which are distributed in WSNs and detects a misused key on behalf of the BS to avoid the communication overheads resulted from the centralized detection by the BS. The detecting nodes have the keying materials so they can check the misused keys without the help of the BS.

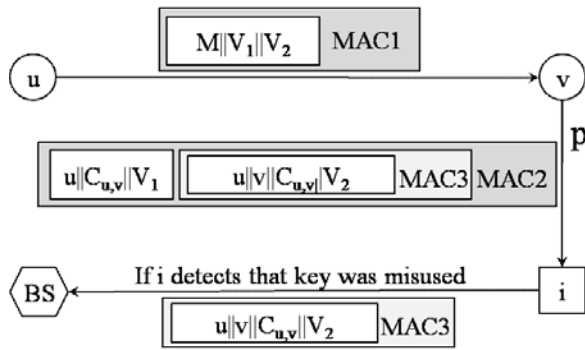
2.3 The Operation

The operation of Liu and Dong's scheme is divided into three parts;initialization phase, message committing phase and detecting phase.

Initialization Phase. Let m be the total number of detecting nodes and n be the total number of sensor nodes. Before deployment, each sensor node u pre-loads with a unique secret key K_u shared with only the BS. Every detecting node i stores all the hashed values $H(K_u, i)$ for every sensor node, $i(i \in n)$. Figure 2(a) describes this step and shows an example of pre-loaded keys for the sensors and the detecting nodes where $m = 2$ and $n = 4$.



(a) Initialization



(b) Committing, sampling and detection

Fig. 2. The operation of Liu and Dong’s scheme

Message Committing Phase. After deployment, each sensor node u makes the detecting node ID list (u) of the detecting nodes around it. If u and v establish a pairwise key $K_{u,v}$, they exchange their detecting node ID lists (u and v) securely using $K_{u,v}$. For every new message M from u to v , u updates the $C_{u,v} = H(C_{u,v}, M)$ then computes two committing values, $V_1 = H(C_{u,v}, H(K_u, i))$ and $V_2 = H(C_{u,v}, K_u)$, where the i is the ID of the detecting node i in the (v). To synchronize the detecting node ID, u and v select the detecting node i at the position $1 + (s \bmod S)$ in (v), where the s is the sequence number of M and the S is the size of (v). Finally, u transmits $M||V_1||V_2$ with $MAC1$ which is generated with the key, $K_{u,v}$.

Detecting Phase. When v receives the message M , v updates $C_{u,v} = H(C_{u,v}, M)$. Depending on the probability p , v makes a report message that contains a two pieces of information. The first one is $u||C_{u,v}||V_1$ for the detecting node i and the other one is $u||v||C_{u,v}||V_2$ for the BS. Figure 2(b) shows these messages. The $MAC2$ is generated by using $H(K_v, i)$ and the $MAC3$ is generated by using K_v .

When the detecting node i receives the report message, i verifies whether V_1 equals $H(C_{u,v}, H(K_u, i))$. If the verification fails, i simply forwards $u||v||C_{u,v}||V_2||MAC3$ to the BS. Finally, the BS checks if V_2 equals $H(C_{u,v}, K_u)$. If not, the BS judges the pairwise key between u and v to be misused and then distributes the new pairwise key to u and v using the unique secret key K_u and K_v .

2.4 The Limitations

The redundancy of CVs. Sensor node u transmits the V_1 and V_2 to node v for every message M but v uses the CVs only when only the probability p is satisfied by the receiver-side v . The unused CVs are big overheads since the packet size of the sensor networks generally is very small. For example, sensor nodes may report the sensed data such as humidity, temperature and illumination. This information can be expressed within a few bytes but the size of CV may be 8 bytes which is the general size of hashed value in the WSNs. Unlike in the traditional network, 8 bytes is considered as a big overhead in WSNs. Besides, the transmission of unused CVs occurs the unnecessary communication overhead.

This overhead comes from the location where the probability p is applied on. The receiver-side v sends the report message to the detecting node under the probability p . Indeed, the sender-side u does not have to send the CVs which will not be used in the receiver-side v .

Unidirectional detection. The detection of a misused key is only limited to one direction ($u \rightarrow v$ or $v \rightarrow u$), since the detecting node can check only one variable ($C_{u,v}$ or $C_{v,u}$) through one report message. To check the misused key of the opposite direction, u has to create a new report message containing $C_{v,u}$ because $C_{u,v}$ is the only unidirectional cumulative hashed value from u to v .

Exchanging of the detecting node lists. The selection of detecting node depends on the sequence number of the message and the size of the detecting node list. The sender node must previously know the detecting node list of receiver node in secure manner. However, this may be infeasible since the establishment of pairwise key between the sender node and the receiver node is not predictable. The reason is that the establishment of a pairwise key not only depends on the routing or the application of the network but also frequently changes as the environment varies.

Unbalanced detection among the sensor nodes. The detecting procedure in Liu and Dong's scheme occurs the unbalanced traffic overhead among sensor nodes. To limit the traffic overhead of the misused key detection, they use the probability p for generating a report message. However, a sensor node u can generate a new report message although u has reported just before. In this case, sending the second report message occurs unnecessary traffic overhead. Rather doing like this, it is efficient that other sensor nodes which did not send the report message recently have to generate a new report message. For doing so,

our scheme will make each sensor node perform the detecting procedure one per $1/p$ message transmitting rounds on average for randomized rotation of the detection procedure among the sensor nodes.

3 The Proposed Scheme

In this section, we present our scheme: a light-weighted misused key detection mechanism. The main focuses of our scheme are the reduction of overhead caused by the report message for detecting misused keys and the maximization of the misused key detection probability while the functionality of the Liu and Dong's scheme is preserved.

For convenience, we define some notations like bellow list.

- u and v are non-compromised nodes.
- $K_{u,v}$ is the pairwise key between u and v .
- f_c is the fraction of detecting nodes that are compromised.
- $C_{u,v}$ is the cumulatively hashed value of the communication from u to v and is concerned with $K_{u,v}$.
- $C_{v,u}$ is the cumulatively hashed value of the communication from v to u and is concerned with $K_{u,v}$.
- (u) and (v) are the detecting node ID lists of u and v in order.

3.1 Our Ideas

Location of the node requesting the detection procedure. To remove the overhead resulted from the unused CVs in Liu and Dong's scheme, we do not send the CVs for every new messages by applying the probability p at the sender-side u . This guarantees that the unused CVs do not transmitted and the traffic overheads are reduced. If T denotes the traffic overhead caused by the CVs, we can reduce the overhead as $T \times p$.

Bidirectional misused key detection. We will check the two variables($C_{u,v}$ and $C_{v,u}$) with one report message. This idea can reduce the network traffic for detecting misused keys by half.

No exchange of the detecting node lists According to the section [2.4](#), the exchange of the detecting node list between two nodes is impractical. In our scheme, the selection of the detecting node is performed by the node which would create a report message. As a result, the detecting node lists do not need to be exchanged previously.

Balanced rotation of detecting procedure among the sensor nodes. In the Liu and Dons's scheme, the misused key detection is unbalanced among sensor nodes when considering the traffic overhead. To overcome this limitation, we apply the probability model introduced in LEACH [\[12\]](#). LEACH uses the

probability model for the randomized cluster leader election of clustering scheme in WSNs. Whenever u sends a message M to v , u generates a random number δ between 0 and 1. If the $\delta < P_{tr}(u, r)$, u performs the detecting procedure which will be described in section 3.2. Equation (1) defines the $P_{tr}(u, r)$, where p is the optimal reporting probability for detecting a misused key of a sensor node, r is the current message transmission round and G is the set of sensor nodes that did not perform the reporting procedure during last $r_{max} = 1/p$ message transmitting rounds. According to Eq. (1), every sensor node must perform the reporting procedure only once during r_{max} message transmission rounds. In this way, the sensor nodes which did not perform the reporting procedure recently have a more chance for detecting the misused key than the nodes that have already performed during last r_{max} rounds. As a result, our scheme ensures that each sensor node performs the detecting procedure only one per $1/p$ message transmitting rounds on average and it achieves the balanced and fair detection of a misused key among the sensor nodes.

$$P_{tr}(u, r) = \begin{cases} \frac{p}{1 - p \times \left(r \bmod \frac{1}{p}\right)} & , \text{if } n \in G \\ 0 & , \text{otherwise} \end{cases} \tag{1}$$

3.2 The Operation

The operation of our scheme consists of four phases;initializing, tracking, reporting and detecting phases.

Initializing Phase. Initialization phase is exactly same with the Liu and Dong’s scheme. Each regular sensor node pre-loads the unique secret key shared with only the BS. The detecting node i stores $H(K_u, i)$ for every sensor node u as described in Figure 2(a). After deployment, each regular sensor node u discovers the detecting nodes around itself and makes the detecting node List (u).

Tracking Phase. In this phase, each sensor node only tracks the history of messages. Whenever u sends a new message M to v , u updates the $C_{u,v} = H(C_{u,v}, M)$ where the initial value of $C_{u,v}$ is zero. v updates the $C_{u,v}$ with $H(C_{u,v}, M)$ if it receives the message M from u . The $C_{v,u}$ also is tracked by u and v for the communication from v to u .

Reporting Phase. Whenever sensor node u sends a message M to v , u appends the detecting node ID, i , to M with the probability $P_{tr}(u, r)$, where i is the ID of detecting node which is randomly selected from (u). As the message ① shown in Fig. 3, u sends $M||i$ with $MAC1$ to v . The $MAC1$ is generated by using the pairwise key $K_{u,v}$ between u and v . Whenever this message is sent, u should set a timer T_1 for detecting the blocking attacks.

Once the sensor node v receives the ID i from u , v updates $C_{u,v} = H(C_{u,v}, M)$. If there is a message M' to send, v updates $C_{v,u} = H(C_{v,u}, M')$. Otherwise, v

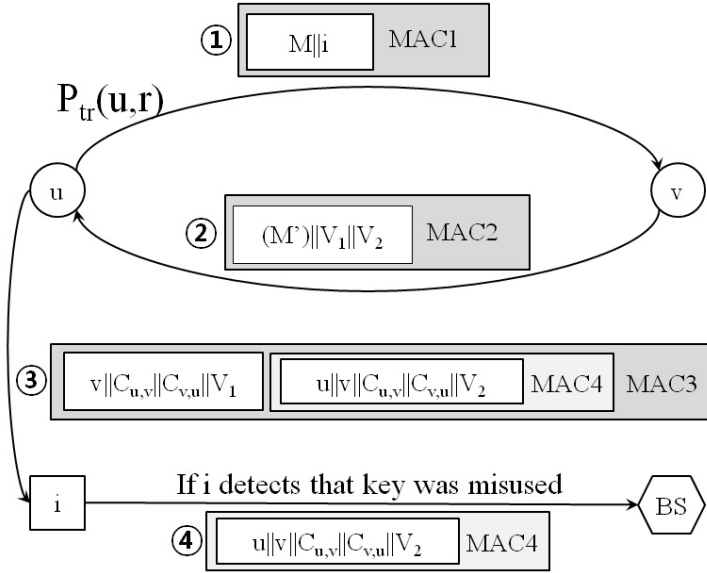


Fig. 3. The operation of the proposed scheme

updates $C_{v,u} = H(C_{v,u}, M)$ and then immediately sends back the committing message, $(M')||V_1||V_2$ with $MAC2$ as the message ② in Fig. 3 (M') denotes the piggybacked message heading for u if there is a message M' to send. If not, only the $V_1||V_2$ is transferred. The V_1 is the hashed value with $H(C_{u,v}, C_{v,u}, H(K_v, i))$ and V_2 is $H(C_{u,v}, C_{v,u}, K_v)$. The $MAC2$ is generated by using the pairwise key $K_{u,v}$ between u and v .

The sensor node u waits the message ② in Fig. 3 until the T_1 expires. If T_1 expires, u reports it to the BS that the message was blocked. Otherwise, u generates a report message as ③ in Fig. 3, where the $MAC3$ and $MAC4$ are generated by using $H(K_u, i)$ and K_u in order. The first piece of information, $v||C_{u,v}||C_{v,u}||V_1$, is for the detecting node i and the second one, $u||v||C_{u,v}||C_{v,u}||V_2$, is for the BS.

Detecting Phase. When the detecting node i receives the report message ③ in Fig. 3, it verifies whether the V_1 equals $H(C_{u,v}, C_{v,u}, H(K_v, i))$. If it is not same, it means that the key $K_{u,v}$ was misused so the detecting node i simply forwards the second piece of information, $u||v||C_{u,v}||C_{v,u}||V_2||MAC4$ as the message ④ in Fig. 3. Otherwise, the detecting node i sets a timer T_2 for detection of the blocking the report message.

The BS checks whether the V_2 equals $H(C_{u,v}, C_{v,u}, K_v)$. If it is not equal, the BS generates a new pairwise key $K'_{u,v}$ and transmits it to the both u and v securely using K_u and K_v . The BS concludes a correct verification with a high probability if the detection procedure at the detecting node was successfully performed.

4 Evaluation

To enhance the Liu and Dong’s scheme, we have introduced our ideas in the previous section. This section compares our scheme with Liu and Dong’s scheme in the aspects of security and energy efficiency.

4.1 Security

False Key Revocation. Any non-compromised key must not be revoked if the key is shared between non-compromised regular sensor nodes. There are three cases that we have to consider.

1. If an adversary sends a forged report message to a detecting node i to revoke a non-compromised key $K_{u,v}$, i will never report a false alarm to the BS. The adversary cannot generate the $MAC2$ to authenticate his forged report message to i since he cannot know $H(K_u, i)$.
2. If a detecting node i is compromised and the adversary tries to revoke a non-compromised key $K_{u,v}$ through i , the adversary will never revoke the key. He cannot create the $MAC4$ because the attacker does not know the shared key between u and the BS, K_u . As a result, the forged message by the compromised detecting node will be always rejected at the BS.
3. If a report message is authenticated by a detecting node i or the BS and the detecting node i is not compromised, the detection of a misused key always is guaranteed because the V_1 and V_2 are protected by the shared key, K_v , with only the BS.

Liu and Dong’s scheme has proved that their scheme satisfies these three cases. Our scheme also satisfies the security requirements as described above.

The Probability of Finding Detecting Node. In the Liu and Dong’s method, they have introduced a detecting node for the purpose of distributing the centralized detecting overheads and also analyzed about the number of additional detecting nodes. According to their analysis, Eq. (2) shows the probability of finding at least one detecting node in any given regular node’s neighborhood, where m is the number of detecting nodes, n is the number of sensor nodes and b is the average number of neighboring regular sensor nodes for every regular sensor node. When $b = 50$, only 10% additional detecting sensor nodes make a regular sensor node find at least one detecting node with a probability 0.99. Our scheme has adopted the detecting node method so this property is also available in our scheme.

$$P_{cover}(m, n) = 1 - \left(1 - \frac{b}{n}\right)^m \approx 1 - \left(\frac{1}{e}\right)^{\frac{bm}{n}} \tag{2}$$

Detection Probability. In the Liu and Dong’s scheme, the probability of a misuse detection is $(1 - (1 - (1 - f_c \cdot p))^{x+1}) \times P_{cover}(m \times (1 - f_c), n)$, where x is the number of messages that v has received since the first misusing of the $K_{u,v}$. The probability of our misused key detection can be inferred like follows.

1. we can define the probability of not performing the detecting procedure until r messages have been transmitted within $r_{max} = 1/p$ rounds as Eq. (3).

$$P_{notr}(u, r) = \prod_{i=1}^r (1 - P_{tr}(u, i)) \tag{3}$$

2. If we assume that u did not performed the detecting procedure until r messages have been transmitted within r_{max} rounds, the probability that u does not perform the reporting procedure until the additional x messages have been transmitted from u to v can be calculated from Eq. (4).

$$P_{notdet|notreport}(u, x, r) = \begin{cases} \prod_{i=r+1}^{r+x} \{1 - P_{notr}(u, i-1) \times P_{tr}(u, i) \times (1 - f_c)\} & , \text{if } x + r \leq r_{max} \\ f_c^{\lfloor (x+r)/r_{max} \rfloor} \times \prod_{i=1}^{(x+r) \bmod r_{max}} \{1 - P_{notr}(u, i-1) \times P_{tr}(u, i) \times (1 - f_c)\} & , \text{otherwise} \end{cases} \tag{4}$$

3. If we assume that u has performed the detecting procedure when r messages have been transmitted within r_{max} rounds, the probability that u does not perform the reporting procedure until the additional x messages have been transmitted from from u to v can be calculated from Eq. (5).

$$P_{notdet|reported}(u, x, r) = \begin{cases} 1 & , \text{if } x + r \leq r_{max} \\ f_c^{\lfloor (x+r)/r_{max} \rfloor - 1} \times \prod_{i=1}^{(x+r) \bmod r_{max}} \{1 - P_{notr}(u, i-1) \times P_{tr}(u, i) \times (1 - f_c)\} & , \text{otherwise} \end{cases} \tag{5}$$

4. Since our scheme performs the bidirectional misused key detection, node v also can perform the detecting procedure for $C_{u,v}$ and $C_{v,u}$. Equation (6) thus denotes the probability that both u and v does not perform the reporting procedure until the additional x and y messages have been transmitted.

$$P_{notdet}(u, v, x, y, r_1, r_2) = \{ P_{notdet|notreport}(u, x, r_1) \times P_{notdet|notreport}(v, y, r_2) + P_{notdet|notreport}(u, x, r_1) \times P_{notdet|reported}(v, y, r_2) + P_{notdet|reported}(u, x, r_1) \times P_{notdet|notreport}(v, y, r_2) + P_{notdet|report}(u, x, r_1) \times P_{notdet|reported}(v, y, r_2) \} / 4 \tag{6}$$

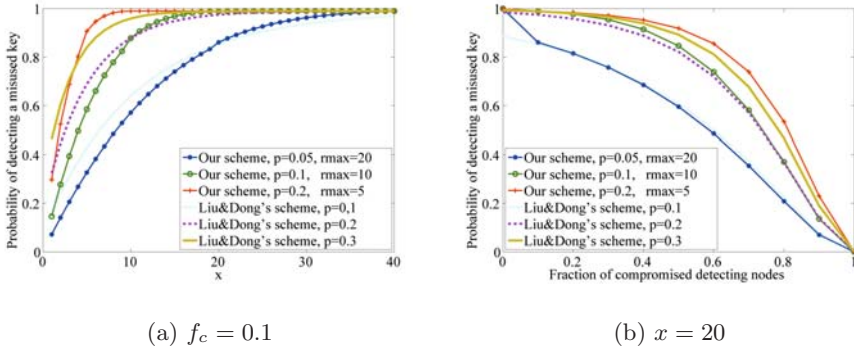


Fig. 4. The probability of detecting a misused key

5. Considering the probability of finding at least one detecting node in any given regular node’s neighborhood, Eq. (7) is the probability of a misused key detection.

$$P_{det}(u, v, x, y, r_1, r_2) = (1 - P_{notdet}(u, v, x, y, r_1, r_2)) \times P_{cover}((1 - f_c) \times m, n) \tag{7}$$

6. Finally, we can conclude the average probability of a misused key detection by removing the variables(r_1 and r_2) like following Eq. (8).

$$P_{detavg}(u, v, x, y) = \sum_{i=1}^{r_{max}} \sum_{j=1}^{r_{max}} \left(\frac{P_{det}(u, v, x, y, i, j)}{r_{max}^2} \right) \tag{8}$$

Figure 4 depicts the probability of detecting a misused key under assuming that the average number of neighbors is $b = 50$, the total number of detecting nodes is $m = \frac{n}{10}$ and $y = x$. In Fig. 4(a), the detecting probability of our scheme is higher than that of the Liu and Dong’s scheme under the same probability $p (= \frac{1}{r_{max}})$. Moreover, the detecting probability of our scheme with $p = 0.1$ is almost same with that of the Liu and Dong’s scheme with $p = 0.2$. Figure 4(b) shows that the probability of detecting a misused key is slowly decreased as the fraction of the compromised detecting nodes is increased. Our performance is more slowly decreased than the Liu and Dong’s scheme under the same probability p . This is because our scheme performs the bidirectional misused key detection and uses the probability model of LEACH.

The p is very important factor since it controls a tradeoff between the traffic overhead and the detecting probability. The higher p is, the higher the detecting probability is but the more the traffic overhead is increased. Our scheme increases the probability of a misused key detection while consuming the lower traffic overhead compared to the Liu and Dong’s scheme.

No Exchange of Detecting Node List. In the Liu and Dong’s method, the detecting node lists (u) and (v) have to be exchanged immediately after

Table 1. The parameters of the simulation

Parameter	Value
E_{elec}	$50nJ$
E_{fs}	$10pJ/bit/m^2$
E_{mp}	$0.0013pJ/bit/m^4$
Size of CV	8 bytes
Size of $C_{u,v}$	8 bytes
Size of message	80 bytes
Size of node ID	4 bytes
Size of MAC	8 bytes
Average of d	6 m

deployment. However, this is very impractical because establishing a pairwise key may depend on the routing protocol deployed for the purpose of an efficient operation of specific applications in WSNs. Especially, in the case of establishing a path key through another sensor nodes, exchanging the detecting node lists is very dangerous because a modified or forged detecting node list may frustrate the whole misused key detection procedure.

In our scheme, we do not have to exchange the detecting node lists. Each sensor node u just keeps its own detecting node list (u) in local repository. If u wants to do the detecting procedure, u just randomly selects one ID i from the (u) and informs it to v . At this time, although an adversary modifies the ID, this modification can be detected at the detecting node i through the detecting procedure since the adversary cannot know $H(K_v, i)$ and K_v . Even though the adversary blocks this message, u can detect this blocking attack using the timer T_1 introduced in the section [3.2](#).

4.2 The Energy Efficiency

In our scheme, we do not send the CVs per message for energy efficiency unlikely Liu and Dong's scheme. In the most general WSNs, the size of message is very small because it only represents the temperature, humidity, illumination and so on. So, it can be expressed within a few bytes but the CVs are generally hashed value which may be about 8 bytes. Sending the CVs per message could cause a big overhead because the size of CV is very large compared to the size of message.

To compare the energy efficiency, we have simply simulated the energy consumption of two nodes u and v using Matlab. For the realistic, the first order radio model in [\[11\]\[12\]](#) is used as a communication model between sensor nodes. Equation [\(9\)](#) and [\(10\)](#) show the amount of energy consumption in transmitting and receiving a packet with l bits over d distance according to the first order radio model. E_{elec} is the amount of energy consumption per bit to run the transmitter or receiver circuitry. E_{fs} and E_{mp} are the amount of energy per bit dissipated in the RF amplifier. The simulation environments are described in Table [1](#). For

simplistic, we have assumed that u and v always have a message to send and transmit to each other by turns.

$$E_{trans} = \begin{cases} l \times (E_{elec} + E_{fs} \times d^2) , if d \leq \sqrt{\frac{E_{fs}}{E_{mp}}} \\ l \times (E_{elec} + E_{mp} \times d^4) , if d > \sqrt{\frac{E_{fs}}{E_{mp}}} \end{cases} \quad (9)$$

$$E_{recive} = l \times E_{elec} \quad (10)$$

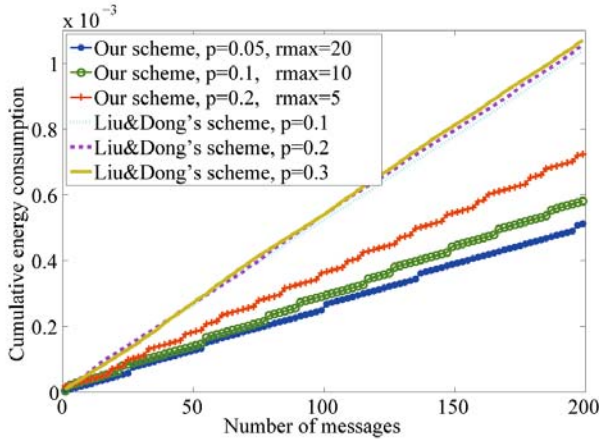


Fig. 5. Energy consumption of a sensor node

Figure 5 shows that the cumulative energy consumption as the number of exchanged message is increased. The energy efficiency of our scheme compared to the Liu and Dong’s scheme is shown in Table 2. As a result, our scheme is more energy efficient than the Liu and Dong’s scheme.

Table 2. Average Energy Efficiency

Probability p	Energy Saved
$p = 0.05$	29%
$p = 0.1$	44%
$p = 0.2$	35%
Average	36%

5 Conclusion and Future Works

In this paper, we have presented an efficient method for detecting a misused key in the WSNs. We have applied our four ideas into the Liu and Dong’s scheme.

Our scheme has better two properties, which are the reduction of the detecting overhead and the high detecting probability, than Liu and Dong's scheme by performing the bidirectional misused key detection, and removing of unnecessary message and using the reasonable probability method for the balanced misused detection among the sensor nodes. We can finally achieve an efficient mechanism with a high detecting probability and low energy consumption. On the average, our scheme is 36% more energy efficient than the Liu and Dong's scheme and the detecting probability of our scheme is higher than that of Liu and Dong's under the same reporting probability. The analysis shows that our scheme is good enough for a misused key detection in the WSNs.

A further direction of this study will be applying our scheme into the existing sensor networks with some specific routing protocols, key management mechanisms, applications and so on. By applying, we will analyze the effectiveness of our scheme.

References

1. Cerpa, A., Elson, J., Estrin, D., Girod, L., Hamilton, M., Zhao, J.: Habitat Monitoring: Application Driver for Wireless Communications Technology. In: Proceedings of the ACM SIGCOMM Workshop on Data Communications in Latin America and the Caribbean, San Jose, Costa Rica (2001)
2. Shen, C., Srisathapornphat, C., Jaikaeo, C.: Sensor Information Networking Architecture and applications. In: IEEE Personal Communications, August 2001, pp. 52–59 (2001)
3. Liu, D., Ning, P.: Establishing pairwise keys in distributed sensor networks. In: Proceedings of the 10th ACM conference on Computer and communications security, pp. 52–61. ACM Press, New York (2003)
4. Liu, D., Ning, P., Du, W.: Group-Based Key Pre-Distribution in Wireless Sensor Networks. In: Proc. 2005 ACM Wksp. Wireless Security (WiSe 2005), September 2005, pp. 11–20 (2005)
5. Lui, D., Dong, Q.: Detecting Misused Keys in Wireless Sensor Networks. In: Performance, Computing, and Communications Conference (IPCCC 2007), April 2007, pp. 272–280 (2007)
6. Chan, H., Perrig, A., Song, D.: Random Key Predistribution Schemes for Sensor Networks. In: Proc. IEEE Sec. and Privacy Symp., pp. 197–213 (2003)
7. Akyildiz, I.F., Su, W., Sankasubramaniam, Y., Cayirci, E.: Wireless Sensor Networks: A Survey. *Computer Networks*, 393–422 (2002)
8. Eschenauer, L., Gligor, V.: A Key Management Scheme for Distributed Sensor Networks. In: Proc. 9th ACM Conf. Comp. and Commun. Sec., November 2002, pp. 41–47 (2002)
9. Eltoweissy, M., Moharrum, M., Mukkamala, R.: Dynamic Key Management in Sensor Networks. *IEEE Communications Magazine* 44, 122–130 (2006)
10. Camtepe, S.A., Yener, B.: Key distribution Mechanisms for Wireless Sensor Networks: A Survey, Technical Report TR-05-07, Rensselaer Polytechnic Institute (March 2005)
11. Mhatre, V., Rosenberg, C.: Design Guidelines for Wireless Sensor Networks: Communication, Clustering and Aggregation. *Ad. Hoc. Networks*, 45–63 (2004)

12. Heinzelman, W.B., Chandrakasan, A.P., Balakrishnan, H.: An application-Specific Protocol Architecture for Wireless Microsensor Networks. *IEEE Transaction on Wireless Communications* 1, 660–670 (2002)
13. Du, W., Deng, J., Han, Y.S., Chen, S., Varshney, P.K.: A Key Management Scheme for Wireless Sensor Networks Using Deployment Knowledge. In: *IEEE INFOCOM 2004*, vol. 1, pp. 7–11 (March 2004)
14. Xiao, Y., Rayi, V.K., Sun, B., Du, X., Hu, F., Galloway, M.: A survey of key management schemes in wireless sensor networks. *Computer Communications Special issue on security on wireless ad hoc and sensor networks* 30(11-12), 2314–2341 (2007)

Identifying Mudslide Area and Obtaining Forewarned Time Using AMI Associated Sensor Network

Cheng-Jen Tang and Miao Ru Dai

The Graduate Institute in Communication Engineering
Tatung University, Taipei, Taiwan
ctang@ttu.edu.tw, d9610002@ms2.ttu.edu.tw

Abstract. An accountable disaster prediction and the appropriate forewarned time are the key issues to reduce the possible damages. Around the globe, landslides and mudslides are serious geological hazards affecting people, and cause significant damages every year. The stability of a slope changed from a stable to an unstable condition that spawns a landslide or mudslide. In most of mudslide-damaged residences, the electricity equipments, especially electricity poles, are usually tilted or moved. Since the location and status of each electricity pole are usually recorded in AMI (Advanced Metering Infrastructure) MDMS (Meter Data Management System), AMI communication network is a substantial candidate for constructing the mudslide detection network. To identify the possible mudslide areas from the numerous gathered data, this paper proposes a data analysis method that indicates the severity and a mechanism for detecting the movement. According the detecting result and the gathered data, this study calculates the remaining forewarned time when an anomaly happens.

Keywords: Advanced Metering Infrastructure, Mudslide Detection, Sensor Network, Meter Data Management System, Data Analysis, Disaster Prevention System.

1 Introduction

Typhoon has a large low-pressure center and numerous thunderstorms that produce strong winds and heavy rains. The strong winds and heavy rains usually cause landslides and mudslides. In Taiwan, there are sometimes over a dozen of typhoons visited in summer. In fact, the name Typhoon means “the deadly storm from Taiwan.” On August 8th, 2009, Typhoon Morakot hit southern Taiwan. It is the deadliest typhoon that sweeps Taiwan in recorded history. The storm brought a record-setting rainfall, nearly 3000mm (almost 10 feet) rainfall accumulated in 72 hours. The rainfall spawned mudslide made a devastating damage to several villages and buried hundreds of lives. There was not any warning mechanism worked when the mudslide embodied as an advent devil. Threatened by such nature disaster, how good a disaster warning system performs decides life and death.

A sensor system is one of the key tools to predict a disaster. Several previous studies proposed some detection techniques. One of the techniques uses the operating theorem of a mechanical mouse as basic. Another uses the image processing to get the

differences between two images. A key issue of a disaster detection system is the communication network. With the rapidly emerging Advanced Metering Infrastructure (AMI), AMI communication network is a substantial candidate for constructing the mudslide detection network. AMI Meter Data Management System (MDMS) records since the location and status of each electricity pole. With such data, a mudslide warning system performs the detection and prediction processes. This study also quantifies the degree of possible damage as the severity indicator. The mudslide warning system triggers different alarms according to the given severity indicator. To obtain a forewarned time, AMI MDMS needs to record the distances information for each electricity pole and its related residence areas. This study then uses the distance information to calculate the remaining forewarned time when an anomaly happens.

2 Background

Previous studies developed many detection and identification techniques for the landslides and mudslides, including Image enhancement [1, 2, 3], Image differencing [4, 5, 6], Vegetation index differencing [7], Image classification [3, 8] and Image registration [9]. Image enhancement emphasizes the landslides or mudslides in the satellites images. This technique requires human experience and knowledge of the observed areas for visual interpretation. Image differencing is a straightforward method and easy to interpret the results. However, this method does not provide a detail changed matrix. The selection of thresholds affects the correctness of the result. Vegetation index differencing emphasizes the differences in spectral response. This method reduces the impacts of topographic effects and illumination. However, this technique also enhances the random noise or coherent noise. Image classification minimizes the impact of atmospheric, sensor and environmental differences between multi-temporal images. However, this method requires a sufficient set of the training samples for classification. Image registration detects landslide/mudslide movements with sub-pixel accuracy. The main concern is its requirements of the high computational cost in terms of CPU time. A simple method of change detection and identification uses a local similarity measure based on the mutual information and the image thresholding [11]. This method uses the mutual information to measure the similarity. It then detects the landslides and mudslides from different images. However, this method is not suitable for detecting a small area landslide. This method fails, if the size of the used imagery is too large in compared with the landslide areas. Table 1 shows the ratio of the approximate size of the landslide over the size of the imagery.

Some studies use the wireless sensor networks (WSN) as the communication infrastructure of their disaster detection system, such as: Slip Surface Localization in Wireless Sensor Networks [10]. This method uses a sensor network consisting of a lot sensor columns. This detection system deploys the sensor columns on hills to find the early signals preceding a mudslide. This sensor network consists of a collection of sensor columns placed inside the vertical holes drilled during the network deployment phase. Installers arrange these sensor columns on a semi-regular grid over the monitored area. The size of the grid depends on the characteristics of the site where the

Table 1. Comparison on the ratio of the approximation of the landslide area with size of imagery used. (Reprinted from [11])

Case	Image size	Box size	Ratio	Land-slides
Niigata	1033 × 1044	450 × 720	1:33.33	✓
Aceh	270 × 270	170 × 130	1:32.99	✓
Muzaffarabad	565 × 768	40 × 65	1:166.89	×
Subimages 1	100 × 100	40 × 65	1:3.85	✓
Subimages 2	200 × 200	40 × 65	1:15.38	✓
Subimages 3	300 × 300	40 × 65	1:34.62	✓

sensors are deployed. Each sensor column has two components: the sensing component that is buried underground and contains all the instruments, and the computing component that stays above ground and contains the processor and radio module.

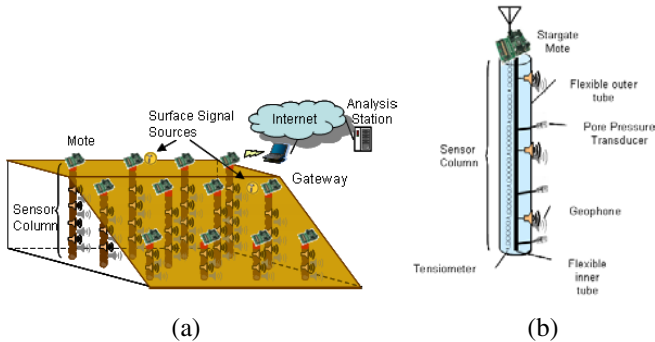


Fig. 1. (a) Landslide Warning System (b) Sensor Column (Reprinted from [10])

The major problem of this technique is the deployment cost. To obtain a “good” prediction, the density of the sensors must be high enough. In other words, if there is a broad area requiring monitoring, the investment must be sky-high. Ironically, the mudslide usually happens in an area that draws little economical interests. Therefore, to protect people living in these areas, the detection devices must be inexpensive. Furthermore, the mounting of these devices must be associated with some public utilities system to reach where people are living.

3 System Architecture

A disaster warning system conducts the environment sensing, data analysis, and communication. A database stores the periodically collected sensor data. The proposed system analyses the gathered sensor data from the database that is a part of AMI system. This disaster warning system uses the tilting degree of an electric pole to define the severity of the location that the pole is set. The system stores every

analyzed data in the same database. Using a specialized suffix-tree algorithm, the warning system sends the forewarning signals according the information. Both data gathering and message broadcasting use AMI-associated communication network.

3.1 AMI-Associated Sensor Network

Recently, many advanced areas such as Europe, America, and Japan etc. commit themselves to develop the Smart Grid, SG. SG is an integration of electricity usage monitoring, generation and distribution automation, meter data management and efficiency improving. SG achieves the energy saving and carbon reduction. SG has functions of energy transmission, energy management and two-way communication, which involve energy generation, transmission, distribution and the customer site. Advanced Metering Infrastructure (AMI) is regarded as the fundamental and key technology of enabling SG.

AMI has a communication network which consists of many advanced metering and sensing devices, including Smart meters. Smart meters provide interval usage (at least hourly) and collect at least daily. According to the analysis from Gartner, AMI system consists of following characteristics: data acquisition, data transform, data cleansing, data processing, information storage/persistence.

Table 2. AMI Process step and involved technologies [SOURCE: GARTNER]

Process Step	Involved Technologies
1. data acquisition	Meter device
2. data transform	Broadband over Power Lines(BPL), Wireless, RF Satellite
3. data cleansing	validation Editing Estimation (VEE) Tools, Meter Data Management System(MDMS)
4. data processing	MDMS
5. information storage/persistence	MDMS
6. information delivery/presentation	Portals, Web Services, Electronic Data Interchange (EDI), MDMS

AMI architecture is shown in Figure 2. The blue lines represent the power lines that connect the transforming stations, utility devices, and customer sites. Each customer site has a meter, a meter interface unit (MIU), and a communication module. The information of power consumption is sent to the data collector through the communication module. The data collectors also use the same communication channel to send back messages. Each smart meter is connected to an energy management device. The smart meter communicates with the control center to manage the energy consumption at its location. The communication lines at the customer sites are presented in green.

For storing the collected data, there is a meter data management system (MDMS) that communicates with the customer information portal, the customer information management system, the dispatch automation system (DAS), and the data collector.

The status of each device within an AMI installation is sent to the data collector, and then forward to the MDMS.

This study proposes an inexpensive detection device that is mounted on each electricity pole. The movement or tilting of an electricity pole triggers the mounted detection device to send a message back to the data collector through the AMI communication network. This displacement message then forwards to the MDMS that activates the mudslide analysis module for further calculations.

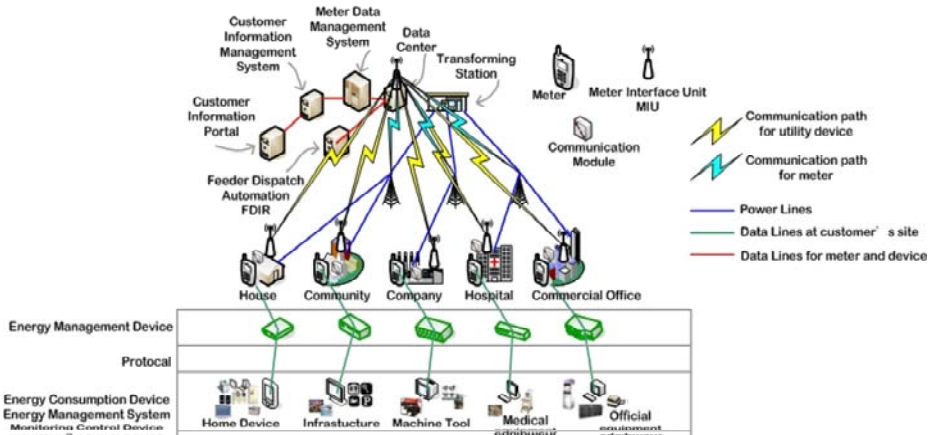


Fig. 2. Advanced Metering Infrastructure

3.2 Detection Method

In most of mudslide-damaged residences, the electricity equipments, especially electricity poles, are usually tilted or moved. In order to obtain the displacement, a movement detector attaches to a pole. A simple method is to install a camera in a pole. Each camera photographs the neighbor pole. Comparing different imagery through image processing techniques identifies whether a movement has happened.

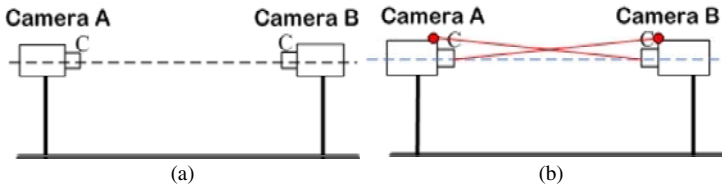


Fig. 3. (a) Camera photographs neighbor one; (b) Camera adding the laser light

The weather affects the performance of this simple detection method. For example, the heavy fog hampers the image identification. The unclear images increase the detecting error. Adding a laser light projector in the upside of camera enhances this situation. Laser light is usually spatially coherent, which means that the light either is emitted in a narrow or low-divergence beam. Coherent light typically means the source produces light waves that are in step.

Another detection method, instead of using an expensive sensor device, the proposed detection technique just uses a portion of a mechanical mouse as the detection device. Everybody knows that when a mouse is moved, the ball inside rolls. This motion turns two axles, and each axle spins a slotted wheel. A light-emitting diode (LED) sends a path of light on one side of each wheel through the slots to a receiving photo-transistor on the other side. The pattern of light to dark is translated to an electrical signal, which reports the position and the moving speed of the mouse. In a word, the operation of a mechanical mouse has five steps as the following: i) moving the mouse turns the ball. ii) X and Y rollers grasp the ball and transfer movement. iii) Optical encoding disks include light holes. iv) Infrared LEDs shine through the disks. v) Sensors gather light pulses to convert to X and Y velocities.

To make this device, one needs to attach a small piece of iron or other material that makes one side of the ball always facing down, and then mounts this “tilt-aware mouse” on an electricity pole. If the pole is moved or tilted, the ball rolls and the displacement messages are sent.

According to how the rollers turn, the tilting direction is obtained. If the RollerA turned to the left, the ball is moving forward, as shown in the Figure 4.

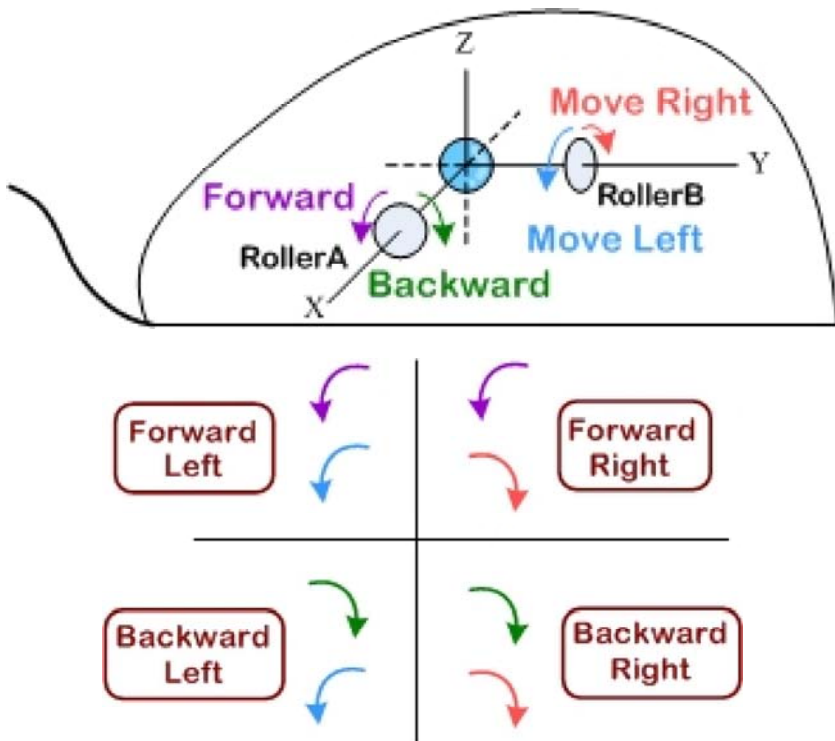


Fig. 4. Detection of the Tilted Direction

To obtain the tilted angle, the ball separates into x-z plane and y-z plane. From the x-z plane, the radian represents the angle of turning left/right, and the vector is \overline{OD} . From the y-z plane, the radian represents the angle of turning forward or backward, and the vector is \overline{OC} . The two vectors \overline{OD} and \overline{OC} construct a plane ODCP with the vector sum \overline{OP} . The angle between \overline{OP} and z-axis is the tilted angle of the electricity pole.

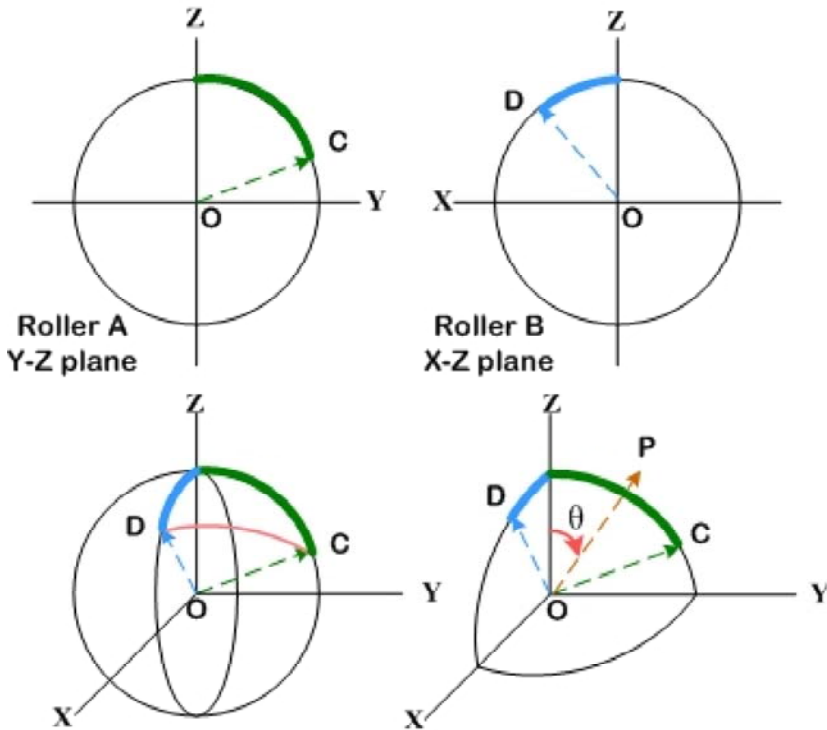


Fig. 5. Detection of the Tilted Angle

4 Data Analysis

The displacement information of the tilted electricity pole is sent to an AMI data collector periodically. The data collector then forwards the information to the MDMS. MDMS then inserts a displacement record in the Pole_Displacement table. A displacement record has the following fields: pole id, time, tilted direction, tilted angle, difference of tilted angle, checking bit, and the tilted severity of an electricity pole. $\{t_1, t_2, t_3, \dots, t_n\}$ denotes the set of recording time of each displacement. a_n denotes the tilted angle. The differences between the continued tilted angles are obtained by $a_n - a_{n-1} = d_n$ and $a_{n-1} - a_{n-2} = d_{n-1}$. d_k denotes a difference between two continuous tilted

angles of a pole. When d_k is not zero, the checking bit c_k is set and the tilted severity is also calculated according to d_k . If $|d_k| \leq 2^m$, the tilted severity s_k is m . If c_{k-1} is set, c_k is also set. If c_k is set and d_k is zero, s_k is set to s_{k-1} .

The pairs of checking bit and tilted severity of all electricity poles constructs the alphabets in the proposed mudslide analysis system. The neighboring poles are assigned with adjacent numbers. The pair $\langle c_i, s_i \rangle$ of time t_i of all numbered electricity poles constructs the S_i according to the pole numbers. The mudslide analysis system then builds a suffix tree ST_i of two adjacent S_i and S_{i+1} . The longest repeating pattern in the ST_i of highest tilted severity indicates the most dangerous area that needs to put most attentions.

Suffix tree is one of the data structure which uses to process a stream data. For example, there is a string “mississippi”.

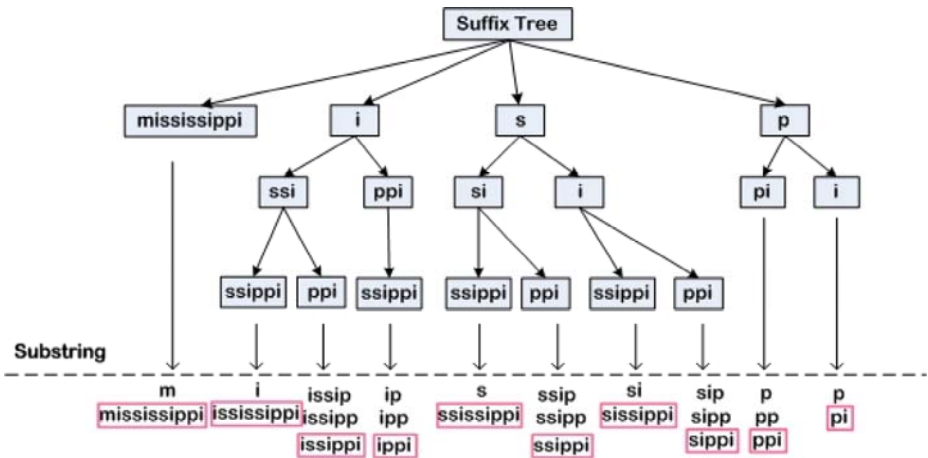


Fig. 6. The suffix tree of “mississippi”

The main reason of using Suffix Tree as the analysis model is for its well-known computation complexity of $O(m)$.

Several substrings have the same suffix such as “i”, “m”, “p”, and “s”. The substrings use ‘is’ as the suffix that also use ‘issi’ as the suffix. In this situation, these substrings share a node in the suffix tree. In the proposed mudslide analysis, this study assumes that the repeating patterns with ci been set indicate where mudslide will probably happen.

A table Pole_Displacement in the MDMS stores the status of all poles. The status of a pole includes its id, time, tilt direction, tilt angle, the difference of tile angle, checking bit, and the degree.

Field “pole_id” is a unique value that indicates an electricity pole in the system. Field “time” uses the UNIX timestamp format, and denotes the time to store data. Figure 4 shows the detection of tilted direction. Through the detection of two rollers A and B, a detector finds the tilted direction. The roller A stays at X-axle for detecting

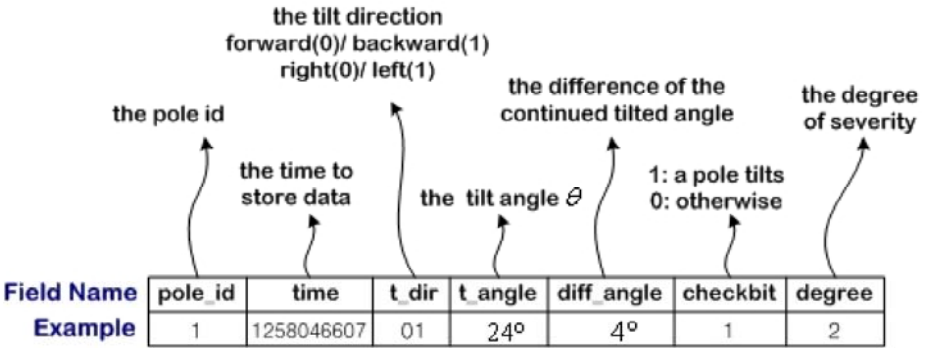


Fig. 7. The fields of table Pole_Displacement

a back-forth movement. The roller B stays at Y-axle for detecting the left-right movement. Field “t_dir” is a 2-bit value that is for recording the back-forth, and left-right movement. Field “t_angle” stores the tilted angle. Field “diff_angle” records the differences of “t_angle”. The field “checkbit” denotes whether a pole is tilted or not. A tilted pole marks “1”, otherwise “0”. Field “degree” denotes the severity degree of a pole. $|diff_angle| \leq 2^m$, and the value m is the degree of severity. An example is shown in Figure 7.

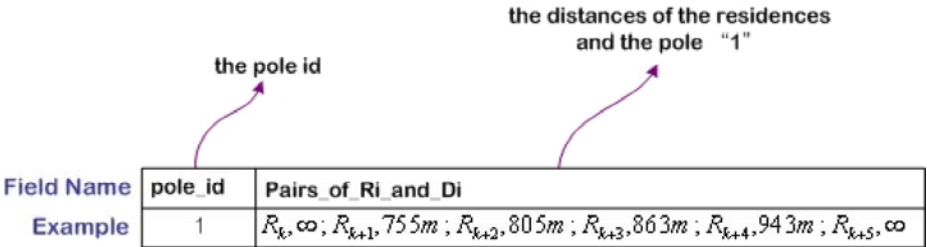


Fig. 8. The fields of table Distance_Ri_Di

The MDMS stores the distances between each pole and the related residences to obtain a forewarning time. Table Distance_Ri_Di stores the distances. There are two fields in this table. Field “pole_id” is the unique id of a pole. The other one is “Pairs_of_Ri_and_Di”. A pair (R_i, D_i) denotes the distance D_i between a pole E_k and a residence R_i . An example is shown in Figure 8.

4.1 Forewarning Time Calculation

Obtaining enough forewarning time is the main purpose of a disaster prevention system. The database of the proposed system stores the possible titling direction of each electricity pole. With these tilting directions, the database stores the affecting scope of each pole. The system gives each residential area a unique identification.

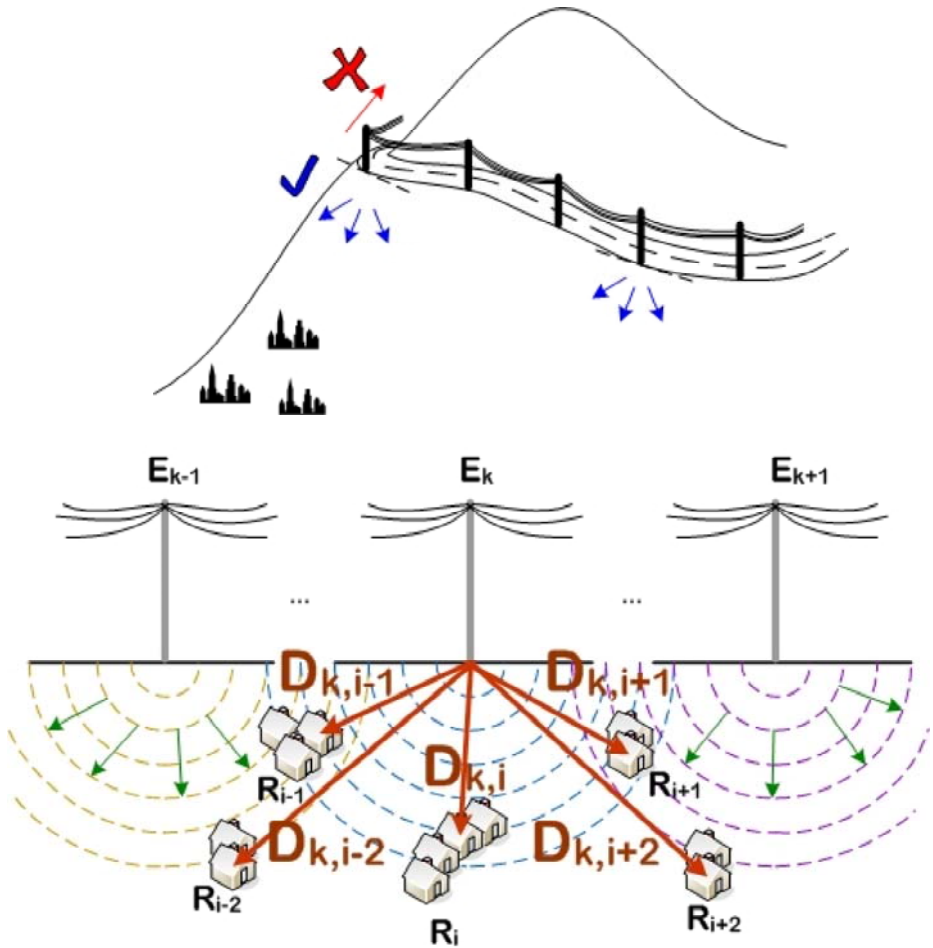


Fig. 9. The tilted directions of poles

Each affecting scope of an electricity pole than contains a set of residence ID. This system associates each electricity pole and every residence ID of this pole with an offset that indicates the distance between the residence area and the pole.

$\{E_1, E_2, \dots, E_k\}$ denotes all poles. $\{R_1, R_2, \dots, R_i\}$ denotes the residence IDs of an electricity pole k . D_{ik} denotes the distance between R_{ik} and E_{ik} . Assume the probability of landslide is $P(x)$.

$$P(x) = \begin{cases} 1 & , \text{landslide happened} \\ < 1 & , \text{otherwise} \end{cases}$$

At any given time after E_k has tilted, the first recorded time is t_1 , and t_0 is subtracting Δt from t_1 . Δt is the period time for retrieving data. The tilted displacement for pole E_k

is m that is $h_k \sin \Delta\theta$. Therefore, the tilted velocity for pole E_k is V_k . t_{0_k} is the time that pole hasn't tilted for the pole E_k .

$$V_k = \frac{c_k \times h_k \sin \Delta\theta}{t_{now} - t_{0_k}} \tag{1}$$

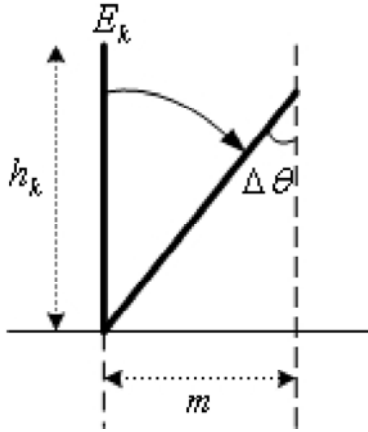


Fig. 10. The tilted displacement of the pole E_k

The forewarned time of the pole E_k alerting to the residence R_i denotes $R_{i,AlertTime}$.

$$R_{i,AlertTime} = \frac{D_i}{c_k \times h_k \sin \Delta\theta / t_{now} - t_{0_k}} \tag{2}$$

5 Conclusion

The mudslides and landslides have consumed many lives. AMI communication network is able to reach every residence that makes a substantial candidate for constructing a mudslide detection network. The location and status of each electricity pole are recorded in AMI MDMS. With an extra movement detector attached to each pole, an AMI-associated sensor network is able to identify the possible mudslide areas from the numerous gathered data. To obtain a “good” prediction, the density of the sensors must be high enough. In other words, if there is a broad area requiring monitoring, the investment must be sky-high. Ironically, the mudslide usually happens in an area that draws little economical interests. Therefore, to protect people living in these areas, the detection devices must be inexpensive. Furthermore, the mounting of these devices must be associated with the public utilities system to reach where people are living. Therefore, this paper proposes an inexpensive detection device associated with AMI network, and a mudslide analysis method that indicates the severity and urgency of a mudslide.

This proposed device and method is still a concept rather than a realization. The future work is to install such system in an AMI pilot installation to prove this concept.

References

1. Whitworth, M.C.Z., Giles, D.P., Murphy, W.: Airbone remote sensing for landslide hazard assessment: a case study on the Jurassic escarpment slopes of Worcestershire, UK. *The Quarterly Journal of Engineering Geology and Hydrogeology* 38(2), 197–213 (2005)
2. Ostir, K., Veljanovski, T., Podobnikar, T., Stancic, Z.: Application of satellite remote sensing in natural hazard management: the Mount Mangart landslide case study. *International Journal of Remote Sensing* 24(20), 3983–4002 (2003)
3. Nichol, J., Wong, M.S.: Satellite remote sensing for detailed landslide inventories using change detection and image fusion. *International Journal of Remote Sensing* 9, 1913–1926 (2005)
4. Cheng, K.S., Wei, C., Chang, S.C.: Locating landslides using multi-temporal satellite images. *Advances in Space Research* 33, 296–301 (2004)
5. Hervas, J., Barredo, J.I., Rosin, P.L., Pasuto, A., Mantovani, F., Silvano, S.: Monitoring landslides from optical remotely sensed imagery: the case history of Tessina landslide, Italy. *Geomorphology* 1346, 1–13 (2003)
6. Rosin, P.L., Hervas, J.: Remote sensing image thresholding methods for determining landslide activity. *International Journal of Remote Sensing* 26, 1075–1092 (2005)
7. Lin, W.T., Chou, W.C., Lin, C.Y., Huang, P.H., Tsai, J.S.: Vegetation recovery monitoring and assessment at landslides caused by earthquake in Central Taiwan. *Forest Ecology and Management* 210, 55–66 (2005)
8. Nichol, J., Wong, M.S.: Detection and interpretation of landslides using satellite images. *Land Degradation and Development* 16, 243–255 (2005)
9. Yamaguchi, Y., Tanaka, S., Odajima, T., Kamai, T., Tsuchida, S.: Detection of a landslide movement as geometric misregistration in image matching of SPOT HRV data of two different dates. *Int. J. Remote Sensing*, preview article:1 12 (2002)
10. Terzis, A., Anandarajah, A., Moore, K., Wang, I.-J.: Slip Surface Localization in Wireless Sensor Networks for Landslide Prediction. In: *IPSN 2006*, Nashville, Tennessee, USA, April 19–21 (2006)
11. Khairunniza-Bejo, S., Petrou, M., Ganas, A.: Landslide Detection Using a Local Similarity Measure. In: *Proceedings of the 7th Nordic Signal Processing Symposium, NORSIG 2006* (2006)
12. Mouse (computing) Wikipedia. Wikipedia (September 24, 2009), http://en.wikipedia.org/wiki/Mechanical_mouse#Mechanical_mice

Utilization of Ontology in Health for Archetypes Constraint Enforcement

Anny Kartika Sari^{1,2}, Wenny Rahayu¹, and Dennis Wollersheim¹

¹ Department of Computer Science and Computer Engineering,
La Trobe University, Victoria 3086, Australia
aksari@students.latrobe.edu.au,
{w.rahayu,d.wollersheim}@latrobe.edu.au

² Department of Computer Science, Gadjah Mada University,
Yogyakarta 55281, Indonesia

Abstract. Most existing works in ontology deployment within the health industry are mainly focusing on the standardization and interoperability goals. In this paper, we propose the utilization of an ontology to apply a new constraint in health archetypes, i.e. the slot filling constraint. An archetype is a model that represents functional health concept such as *admission record*. It can reuse other existing archetypes through a slot. The name of a slot represents a more specific health concept such as *head*. The slot filling constraint restricts the selection of archetypes to fill in that specific slot so that only relevant archetypes are chosen from the available ones. Ontology is used to enforce this constraint. An approach on how to apply the constraint is presented based on the semantic similarity/relevance concept. The evaluation shows that the approach is a better alternative to the current slot filling process which depends on manual decision by the archetype author.

Keywords: ontology, health, archetype, slot filling constraint, semantic relevance.

1 Introduction

Ontologies have been used widely in many areas today. It is primarily employed in the applications where formal specifications of concepts are needed. One of the fields that utilize ontology is health. There are even some ontologies related to health which have been defined and become standards such as SNOMED CT (Systematized Nomenclature of Medicine-Clinical Terms)¹ and LOINC (Logical Observation Identifiers Names and Codes)². However, the main use of ontologies in health so far is still limited to the effort of achieving the uniform semanticity of health terms, in which the clinical ontologies become standards. Another common aim of ontology usage in health is to achieve interoperability between the different frameworks of electronic health records (EHR), which are presented in e.g. [1], [2], and [3].

¹ See: <http://www.ihtsdo.org/>

² See: <http://loinc.org/>

This paper will present another kind of ontology employment in health, specifically in the area related to EHR. Our focus is the EHR architecture proposed by openEHR³, in which archetype lays as the basis of the model. The archetype concept was first mentioned as part of the Australian project GEHR (Good Electronic Health Record) [4]. The development of the archetype concept is mainly motivated by the inability of the existing medical information system to keep up with the dynamic changing of health knowledge. In the existing framework, which is commonly referred as “single-level” methodologies, both informational and knowledge concepts are built into one level of object and data models [5]. This methodology will become obsolete soon after there is knowledge updating, and then should be replaced. This produces high maintenance cost. The archetype concept is proposed to cope with the shortcoming of the existing medical information model. The key feature of an archetype-based system is that it is made up of a two-level architecture: a RIM (Reference Information Model), consisting of a set of reference types, and a domain or knowledge model, consisting of a set of archetypes [4]. The reference model is the repository of the patient clinical data. As clinical knowledge is stored completely separately in archetypes, resulting EHR systems are more flexible as changes in the clinical knowledge can be embraced by modifications in archetypes, without compromising the integrity of information in the repository [5]. An archetype itself formally expresses a distinct health concept such as *blood pressure*. The concept is expressed in the form of constraints on data whose instances conform to some RIM [6].

In this work, we are proposing a new archetype constraint, named as the *archetype slot filling* constraint, which is defined as the complement of the existing constraints specified in the archetype specifications. The proposal of the new constraint is motivated by an example of the data inconsistency on the absence of this constraint. With the enforcement of the constraint, we believe that new archetypes can be defined more accurately and semantic inconsistency can be avoided. At the end, this will improve the quality of the whole clinical systems in which archetypes are utilized.

As the constraint is not specified in the archetype specifications, a method will also be proposed to make the constraint feasible to be applied without influencing the internal structure of archetypes. Ontology, in this case clinical ontology, will be employed for the application of the approach. The method also uses the concept of semantic similarity/relevance of an ontology. Of course, the concept needs to be adjusted to fit in the clinical ontology to be used.

Actually, the concept of ontology has been used in an archetype as one of its sections, but the usage is limited to terminologies (clinical ontology) linkings and bindings. The binding is aimed at the uniformity in the meaning of terms contained in the archetype, even though there is no special mechanism to check the uniformity. The ontology is simply viewed as collection of terms, not as a ‘real’ ontology. This is different to the approach in this work in which ontology will be utilized as the integration of all its concepts, relationships, and attributes.

The rest of the paper will be organized as follows. Section 2 will present the motivating example on the proposal of the new constraint in archetype concept. Related works is presented in Section 3. The approach to enforce the application of the

³ See: <http://www.openehr.org/>

constraint, as well as the related algorithms, is then explained in Section 4. The approach is evaluated in Section 5. Section 6 closes this paper with some conclusion and future works.

2 Motivating Example and the Slot Filling Constraint

According to [7], the term archetype is used to denote knowledge level models which define valid information structures. The domain concepts defined in the models are expressed using constraints. Constraints are used to govern the variability that may occur in the concepts definition. This idea leads to the definition of archetype as *a model defining some domain concept, expressed using constraints on instance structures of an underlying reference model* [7]. From this definition, it can be inferred that basically archetype contains constraints. Some important built-in constraints are related to the occurrence of an object node, the existence of an attribute, the cardinality and ordering of a container attribute, and the containment of an archetype inside another archetype. The syntax used to represent the constraints is cADL (constraint ADL), while dADL (data ADL) syntax is deployed to express instance data.

In the archetype definition using cADL, constraint in an archetype that is being developed can be specified by reusing the existing archetypes. This mechanism is known as an archetype ‘slot’ or ‘chaining point’ [8]. Basically, a slot is the means for an archetype to link to another archetype and reuse the entire constraints defined in the linked archetype. The definition of an archetype slot uses the keyword `allow_archetype` and utilizes two types of assertion statement: `include` and `exclude`. The assertion statements are the constraints in the slot filling. Archetypes contained in the `include` statement can be chosen as the slot fillers, while those contained in the `exclude` statement must not be chosen as the slot fillers. The archetype slots are available in two of the archetype classes, i.e. the Composition archetypes and Section archetypes. The classes of archetype that can fill in the slots are Section, Entry, Observation, Evaluation, Action, and Administration Entry archetypes.

The wrong choice of the slot fillers may raise semantic inconsistency. An example of the problem is described here. Consider the openEHR archetype *openEHR.EHR.SECTION-admission_record.v1*. Fig. 1. shows the view of that archetype as displayed by the Archetype Editor tool. The archetype contains Physical Examination sub section which has Regional sub section. One of the sub sections of Regional is Head. Head sub section includes Observation archetypes as its slot. In this slot, the inclusion and exclusion assertion are not defined. Thus, all archetypes of Observation class can be used to fill in the archetype slot. However, some archetypes like *openEHR.EHR.OBSERVATION-autopsy.v1*, *openEHR.EHR.OBSERVATION-menstruation.v1*, and *openEHR.EHR.OBSERVATION-urinalysis.v1* are not relevant to the sub section since they are not semantically related to the term ‘head’. Since there is no semantic relation between the name of the sub section (slot) with some of the slot fillers, semantic inconsistency happens in this situation. At the end, this can lead to the semantic inconsistency of in the data level which is captured using the archetype.

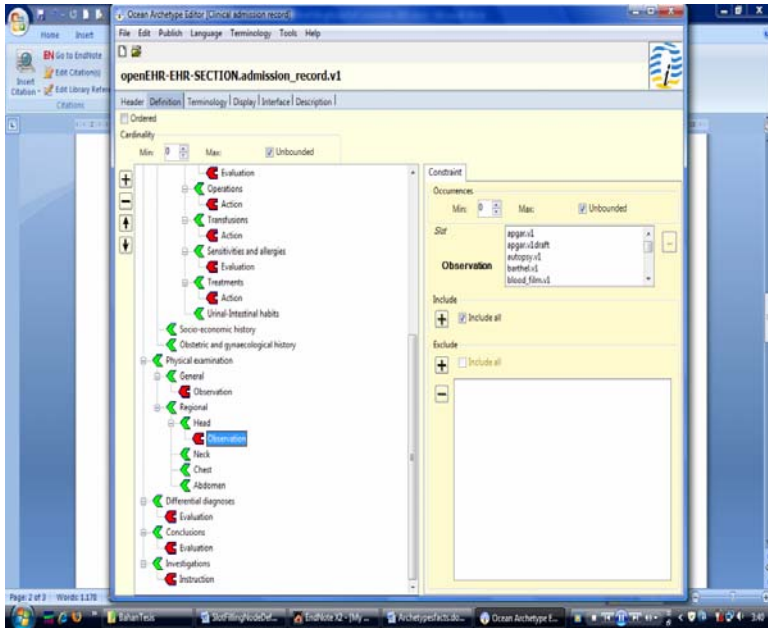


Fig. 1. The *openEHR.EHR.SECTION-admission_record.v1* and its slot in the Head sub-section

The semantic violation problem occurs because the user, i.e. the archetype author should fill in the archetype slots without the guidance on which archetypes can be chosen and which ones should be excluded due to their irrelevancy to the sub section name. In this example, user should inspect the list of the fifty Observation archetype names and decide by himself which names can fill in the slots. Since archetype concept is expected to be used widely in the future, and thus, more and more new archetypes will be created, it will be more complicated for the user to choose the relevant archetypes when there will be many more than fifty Observation archetypes.

This problem could be prevented if there is a constraint in filling in the slot for archetypes. This constraint should determine the dependency of the archetype slot filling to the sub section it is related to or to the archetype itself. We refer this constraint as *archetype slot filling constraint*. A method that utilizes ontology can be applied to enable the slot filling constraint. This method will provide the user with the relevant archetypes to be filled in a particular slot. The relevancy is measured based on a particular clinical ontology. However, since the archetype should neutral with respect to terminologies [6], the utilization of terminology/ontology should be conducted outside the archetype definition. The detail of the method is discussed in Section 4.

3 Related Works

There are not many works which utilize ontologies for the works related to archetypes. However, to the best of our knowledge, there is no existing paper which discusses the employment of ontology for the validation of archetype constraints.

Rather, most of the works use ontology as a mediator to achieve semantic interoperability between the different frameworks of EHR. For example, [1] shows that the interoperability between archetypes based web service messages could be accomplished by using OWL (Web Ontology Language) as the representation of both archetype definitions and instances to semantically annotate the web service messages. Similarly, in [2] one of clinical ontologies i.e. SNOMED CT is used as the standard into which archetypes and messages are bound. The aim is to allow different systems to interoperate. And finally, Fernandez-Breis, et al., in [3], presents an ontological approach to promote interoperability among CEN(European Committee for Standardization)/TC 251⁴ and openEHR compliant information systems by facilitating the construction of interoperable clinical archetypes. This is achieved by constructing a one to one mapping between the models, utilizing OWL as the ontological representation.

The concept of semantic relevance is also employed in this paper. It is used in the approach to apply the new constraint in the archetypes. Some works which are related to the relevance semantic issue have been done. In most of them, e.g. [9], [10], [11], [12], [13], and [14], the *semantic relevance* term is related to, or based on, or even used interchangeably with the terms *semantic similarity* and/or *semantic distance*. Hence, in this work we will also base the semantic relevance measurement on semantic similarity measurement.

There have been many papers related to semantic similarity. According to [15], based on the structures and input considered, there are three kinds of similarity measures: edge-based, information content-based, and feature-based approaches. Edge-based approach is used in a graph representation of an ontology. This approach is used in e.g. [16] and [17]. Information content-based method, which is applied in [18], [19], uses text corpus to evaluate the semantic similarity. The feature-based technique utilizes the properties of concepts for assessing the similarity and is used in [20]. Some works, such as [21], [22], and [23], use the combination of those methods or compare them.

To the best of our knowledge, none of the papers related to semantic similarity discuss its application to clinical ontology, which is the focus of this paper. Clinical ontology, in this case SNOMED CT, has some special characteristics. First, the meaning of *attribute* is not the same as *attribute* in common ontologies. In SNOMED, *attribute* is a kind of *relationship* which relates a concept to another. This different meaning should be taken into account when feature-based similarity techniques are used, as they commonly base the measurement on *attribute*. The next characteristic is that in SNOMED, the hierarchy is based only on the IS-A relationships. Other types of relationships are not considered. This will influence the calculation when the edge-based similarity approach, which is based on all types of relationships, is used. Thus, there should be adjustment when the approach is applied for SNOMED CT clinical ontology.

⁴ CEN/TC 251 is standard for health informatics. The complete published standards, see: <http://www.cen.eu/CENORM/Sectors/TechnicalCommitteesWorkshops/CENTechnicalCommittees/CENTechnicalCommittees.asp?param=6232&title=CEN/TC%20251>

4 The Approach for the Slot Filling Constraint Application

This section describes the approach to enforce the slot filling constraint application in the development of new archetypes. The following sub sections explain the approach in detail.

4.1 Description of the Method

As previously mentioned, we propose the ontology based approach to apply the slot filling constraint to archetype. We believe that ontology is suitable to be utilized to force this constraint for the following reason. The constraint is much related to the semantic data inconsistency. As ontology is a formal specification, it can maintain semantic consistency of the concepts it specifies. As the result, the problem of data inconsistency can be avoided.

In this work, SNOMED CT is chosen as the ontology for the approach. The choice is based on some reasons. Firstly, SNOMED CT is considered to be the most comprehensive, multilingual clinical healthcare ontology in the world [24]. It covers areas that are closely related to EHR, the main concern of archetypes. Secondly, some works such as [2] and [25] use SNOMED CT as the ontology for archetype related works. This shows that SNOMED CT fits well to the concepts defined in archetypes.

Since the slot filling constraint involves only particular classes of archetypes, i.e. Composition, Section, and Entry, this paper only considers those types of archetypes as well. For the purpose of clearness of the method used, we categorize those classes of archetypes into three class layers: the bottom layer, the middle layer, and the top layer. The bottom layer is the classes of archetypes which can only be the filler of slots. They are not able to have any slot to be filled in by another archetype. All the entry classes, i.e. the Instruction class, the Observation class, the Action class, the Evaluation class, and the Admin Entry class, are included in this layer. The top layer is the class of archetypes which can only have slots inside them. They themselves cannot fill in another archetype's slot. This layer includes the Composition class. The middle layer comprises the Section class. This layer of archetypes can have slots to be filled in by other archetypes as well as be the filler of other archetypes. In the rest of this paper, we will use the layers categorization to refer to the classes of archetypes. Table 1. summarizes this categorization.

Table 1. The layer categorization of archetype classes used in this work

Name of Layer	Included Classes of Archetypes
Bottom	Instruction, Observation, Action, Evaluation, Admin Entry
Middle	Section
Bottom	Composition

The ontology based approach to force the constraint application is explained intuitively in this following scenario. An archetype author needs to fill in the slots with the `include` and `exclude` assertion. The archetypes to be included or excluded in the slot are chosen from some groups of a specific type of archetypes which can be one of the bottom or middle layer class. The archetype type is chosen by the author. Each

archetype of the determined type is a *candidate* for the `include` assertion. Since there can be many candidates, it is essential to assure that the chosen archetypes are indeed relevant to the sub section or the Section/Composition archetype that contains the slot. For this purpose, ontology is employed. It is used to annotate each archetype with the concepts it specifies which are related to the archetype. The annotation process is performed directly after an archetype of bottom or middle layer is built. The archetypes of upper layer, i.e. the Composition archetypes, do not need to be annotated because they are never used as slot fillers.

When an archetype is included as an archetype candidate to fill in a slot, its annotation will be used in the filtering process. This process filters the archetypes in the candidate list which are relevant to the slot. The relevancy is measured by inspecting the relevance of each annotation of each archetype candidate with the concept in the ontology which is related to the keyword. The user chooses the keyword from the name of the sub section where the slot is contained, the names of all the ancestral sub sections, the name of the archetype, or the combination of them which then become a set of keywords. Each archetype candidate is then evaluated whether it is considered relevant or not to fill in the slot. If it is, then the candidate is included in the list of the relevant archetypes for the slot, otherwise it is excluded. The list is then used to force the slot filling constraint to the slot filling activity by the user.

From the above description of the method, we determine two technical steps in the approach: the preliminary process and the filtering process. In this paper, we focus on the filtering process, thus we will describe the preliminary process briefly, while the filtering process is described in detail in Section 5.2.

The preliminary process refers to the archetype annotation with respect to the SNOMED CT ontology. Each archetype is annotated with some terms which represent the concepts in SNOMED CT which the archetype relates to. The annotation is applied to any archetype of bottom and middle layer upon its development. Thus, it is not limited to the archetypes in the candidate list. The process involves two main steps: the ontology extraction and the archetypes annotation. The ontology extraction is performed to obtain only the relevant part of the whole ontology, which is called *sub-ontology* or *ontology view*. We reuse the ontology extraction proposed in [26] for this work. The annotation is used in this work to annotate a specific archetype document with the concepts specified in SNOMED CT which are related to its content. Some annotation tools, such as SMORE [27] and Ontomat [28] can be utilized for the purpose.

4.2 The Filtering Process

This is the main process of the method. In this process, the relevant archetypes for the slot to be filled in are determined using a specific evaluation method. First, the archetype candidate list is prepared. This is based on the specific type of archetype chosen by the user, i.e. the archetype author, for the archetype slot to be filled in. then, the evaluation of relevant archetypes is conducted. Two steps need to be conducted here, i.e. the determination of the set of keywords and the relevance calculation, which result in the list of relevant archetypes. For the user, choosing archetypes not included in the list will violate the *archetype slot filling constraint*.

The method for the relevance evaluation will be elaborated. The relevance of a concept included in the annotation of a specific archetype to the keyword is evaluated in a simple way. Algorithm 1 below describes the evaluation method.

Algorithm 1: Relevance evaluation

Variables :

K : list of keywords

k_h : the h -th keyword

P : list of archetypes candidate

R : list of relevant archetypes

h, i, j : counter

n : the number of archetypes in list P

P_i : the i -th archetype in list P

m : the number of annotations for archetype P_i

P_{ij} : the j -th concept included in the annotation of archetype P_i

T : ontology (SNOMED CT)

relv: boolean variable to indicate if the annotation P_{ij} is determined to be *relevant* to the keyword

Input : $K, P = [P_1 \dots P_q], T$

Output : R

Algorithm:

$R \leftarrow []$

For $h = 1$ to n do

$i \leftarrow 1$

Do while $i \leq q$

$rel \leftarrow false$

$j \leftarrow 1$

Repeat

If P_{ij} is relv to k_h based on T then $relv \leftarrow true$

else $j \leftarrow j + 1$

until ($relv = true$) or ($j > m$)

If ($relv = true$) then

$R \leftarrow R + [P_i]$

Delete P_i from P

$q \leftarrow q - 1$

endif

enddo

endfor

The relevance evaluation is calculated for each keyword. The keywords are chosen by the user based on the name of archetypes and the name of sub sections and all its descendants where the slot is placed. For each keyword, each archetype is evaluated its relevance to the keyword using its annotations. Each annotation is inspected whether it is considered *relevant* to the keyword or not. Once it finds that an annotation of the

archetype candidate is considered *relevant* to the keyword based on the ontology, it will add that archetype to list R, which includes the relevant archetypes for the slot. Thus, the evaluation of the rest of the annotations is not needed. Moreover, the archetype will be deleted from list P (list of candidate archetypes). It means that for the relevance evaluation of the next keyword, it will not be considered anymore.

For evaluating the relevance between the keyword k with concept P_{ij} in the archetype annotation, first one or more concepts which have the names similar to k are searched in SNOMED. Each concept found is then evaluated its relevance to P_{ij} based on SNOMED CT ontology. Since the basis is an ontology, certainly the word *relevance* here refers to the *semantic relevance*, which means the closeness of the relationship between two terms which in this work refers to the concepts contained in SNOMED CT ontology. As previously mentioned in Section 3, semantic relevance is much related to semantic similarity. Thus, some previous works on semantic similarity are reused here with some adjustment.

To measure the semantic relevance between two concepts in SNOMED CT, we first determine some features which influence the assessment. The factors are based on the edge-based and information content based approaches of similarity calculation. They are described below.

1. The distance between the concepts

In edge-based approach, the shortest path between two concepts in the semantic net is commonly used to evaluate their similarity. The shortest path of concept a and b is the path where the number of edges in the path is minimum among all the paths connecting a and b in the nets. Rada et.al. [16] propose a basic semantic distance measure by only using taxonomic links and computing the shortest path between the concepts. A variant of the work is [29] in which distance is substituted by similarity. Thus, we consider the equation is more suitable to our work and it is reused in this paper. By adopting the equation in [29], we define the function for semantic relevance as follows:

$$rel_e(a,b) \propto -\log \frac{dist(a,b)}{2 \times Max} . \quad (1)$$

where $rel_e(a,b)$ denotes the edge-based relevance between concept a and concept b , $dist(a, b)$ represents the number of edges of the shortest path between concept a and b , and Max is the maximum depth of the taxonomy.

2. The link type connecting the concepts

There are 4 types of relationships in SNOMED CT: defining characteristics (IS-A relationships and defining attributes), qualifying characteristics (non-defining, qualifying attributes), historical relationships (relate inactive concepts to active concepts), and additional relationships (other non-defining characteristics such as PART OF) [30]. The type of relationship that form the hierarchy of SNOMED CT is the IS-A relationship. The other types of relationships are not represented as links in the hierarchy. Two special cases occur here and we consider that the relevance between two concepts is absolute if at least one of them happens. The first case is when the two concepts are actually the same concept; the other is when one of the concepts or one of their descendants becomes an attribute of the

other concept. The second case is based on the consideration that if the annotation (which is a SNOMED CT concept) of an archetype becomes an attribute of the keyword (which is also a SNOMED CT concept) of a slot, then the archetype should be relevant for the slot. It is unlikely that the attribute-based relationships form a path. Thus, we only consider their existence, not their path length. On the other hand, the PART-OF relationships often form a path. Thus, the PART-OF path distance between concept a and b needs to be considered. To differentiate the IS-A and PART-OF relationships, we define the equations below which are based on (1):

$$\text{sim}_{\text{IS}}(a, b) = -\log \frac{\text{dist}_{\text{IS}}(a, b)}{2 \times \text{Max}_{\text{IS}}} / \log(2 \times \text{Max}_{\text{IS}}). \quad (2)$$

$$\text{sim}_{\text{PART}}(a, b) = -\log \frac{\text{dist}_{\text{PART}}(a, b)}{\text{Len}_{\text{PART}}(a, b)}. \quad (3)$$

where $\text{sim}_{\text{IS}}(a, b)$ denotes similarity between concepts a and b by using the IS-A hierarchy, $\text{dist}_{\text{IS}}(a, b)$ is the shortest path between concepts a and b based on the IS-A hierarchy, Max_{IS} is the maximum depth of the IS-A taxonomy, $\text{sim}_{\text{PART}}(a, b)$ represents the similarity between concepts a and b based on the PART-OF relationship, $\text{dist}_{\text{PART}}(a, b)$ is the shortest distance between concepts a and b on the PART-OF path, and $\text{Len}_{\text{IS}}(a, b)$ is the maximum length of PART-OF path which contains a and b .

Then, we need to combine the equations and the consideration of attribute relationship to get the relevance value between concepts a and b , as follows:

$$\text{rel}_{\text{E}}(a, b) = \begin{cases} 1, & \text{one of special cases happens} \\ \max \{ \text{sim}_{\text{IS}}(a, b), \text{sim}_{\text{PART}}(a, b) \}, & \text{otherwise} \end{cases} \quad (4)$$

3. The information content of the concepts

The value of the information content of a concept is shown by the probability of occurrence of this class in a large text corpus [31]. Since we work with the collection of archetypes, we consider that it is more suitable to use the collection of archetypes as text corpus. As presented in [32], the functions for this factor is defined as follows.

$$\text{rel}_{\text{I}}(a, b) = \frac{1}{\text{IC}(a) + \text{IC}(b) - 2 \times \text{IC}(\text{LCS}(a, b))}. \quad (5)$$

where $\text{rel}_{\text{I}}(a, b)$ is the information content-based relevance between concepts a and b , $\text{IC}(a) = -\log P(a)$ is the information content of concept a and $\text{LCS}(a, b)$ is the least common subsumer of concept a and b . $P(a)$ is the probability of encountering an instance of concept a in the corpora.

To get the final relevance equation, we combine (4) and (5). The following is the final equation.

$$\text{Rel}(a, b) = \alpha_1 \text{rel}_{\text{E}}(a, b) + \alpha_2 \text{rel}_{\text{I}}(a, b). \quad (6)$$

$Rel(a, b)$ is the relevance between concepts a and b , while α_1 and α_2 are adjustable factors, $\alpha_1 + \alpha_2 = 1$. α_1 and α_2 can be adjusted to the requirement. For example, when the information about the text corpus is missing, α_2 can be set to 0 to eliminate the information content-based relevance from the equation.

Finally, equation (6) is connected with Algorithm 1 to determine the value of variable `relv`. Equation (9) is used for the purpose.

$$Relv = \begin{cases} \text{true, } Rel(a, b) \geq \beta \\ \text{false, otherwise} \end{cases} \quad (7)$$

β is the value which is determined by the user or the archetype author.

5 Evaluation and Discussion

To evaluate the method, we should apply it to the process of filling in some slots. We choose two types of slots. The first type includes the slots which are contained in sub section(s). These slots can have more than one keyword to be chosen. On the other hand, the second type, which is not contained in any sub section, can only have one keyword, i.e. the archetype's name. We will try to find out whether the number of keywords influences the result. Four slots from each type are picked from the archetypes collection of the Archetype Editor version 1. Not many archetypes have sub sections; even, we only find the *SECTION.admission_record.v1* archetype which has more than one level of sub section. Thus, we choose four slots from it as the representation of the first slot type. The slot names are *Injury*, *Sensitivities and allergies*, *Childhood illness*, and *History*. For the second type, we choose 4 archetypes which does not has sub section: *SECTION.family_history*, *SECTION.encounter*, *SECTION.summary*, and *SECTION.antenatal_check* archetypes. For the filtering process, we set $\alpha_1 = 0.95$ and $\alpha_2 = 0.05$ as we consider that the number of archetype collection, as the source of text corpus, is still limited. β is set to be 0.55.

Table 2 shows the comparison on some statistical metrics between the slot filling result using the method and the current slot filling in which the method is not used. The precision, recall, accuracy, and F-measure are calculated based on the manual decision of the slot filling process. As the number of archetypes is still low at this moment, we believe that manual decision can be used as the valid basis for the calculation.

From the table, we can see that the precision value of the current slot fillings tend to be higher than the slot filling result using the method, except for the Encounter and Antenatal check slots. This can be understood as in the other slots, the percentage of archetypes to be included is much lower than the actual relevant archetypes. For instance, the current number of included archetypes in the Injuries slot is only 2, compared to 15 actual relevant archetypes for that particular slot. For the Encounter and Antenatal check slots, all Evaluation archetypes are included to fill in the slots. That is why the precision is lower than the result on using the method, as not all archetypes are actually relevant to those slots.

Table 2. The comparison on some statistical metrics between the result using the method and

Slot Names	Number of Keywords	Precision		Recall		Accuracy		F-Measure	
		Using the Method	Without the Method	Using the Method	Without the Method	Using the Method	Without the Method	Using the Method	Without the Method
Injury	4	0.77	1.00	1.00	0.13	0.77	0.50	0.83	0.24
Sensitivities and allergies	4	0.85	1.00	1.00	0.06	0.85	0.38	0.89	0.11
Childhood illness	4	0.82	1.00	0.88	0.06	0.81	0.42	0.85	0.12
History	2	1.00	1.00	0.65	0.04	0.69	0.15	0.79	0.08
Family history (archetype)	1	0.60	0.80	0.75	0.33	0.65	0.65	0.67	0.47
Encounter (archetype)	1	0.85	0.23	0.79	1.00	0.81	0.23	0.81	0.38
Summary (archetype)	1	0.75	0.83	0.60	0.45	0.77	0.73	0.67	0.59
Antenatal check (archetype)	1	0.71	0.50	0.83	1.00	0.77	0.50	0.77	0.67

On the other hand, the recall value of the result using the method is generally higher than the current slot filling. Again, the exception is for two slots, i.e. the Encounter and the Antenatal check slots, as those two slots included all Evaluation archetypes. This makes the recall value to be maximum, because the inclusion of all archetypes results in the inclusion of all relevant archetypes as well.

The accuracy and F-measure metric calculations normalize the opposite results of the recall and precision metrics. The method generally has higher accuracy than the current slot filling condition, which means that the method produces the closer results to the actual relevant archetypes for each specific slot. The result of the F-measure calculation is also higher. This shows that the method is more effective in retrieving relevant archetypes for a specific slot than the current slot filling. From the result of the two metrics, the slot filling process using the method seems to be a better alternative for the current slot filling condition.

The number of keywords does not seem influencing much on the result. For the slots with multiple keywords, the accuracy value ranges from 0.69 to 0.85, while that of the slots with single values ranges from 0.65 to 0.81. It shows that the number of keywords does not influence the accuracy. A slightly difference measure happens to the F-measure metric, in which the slots with multiple keywords have a little bit higher F-measure value. This means that the method is more effective to be applied to the slots which are contained in a sub section rather than directly in an archetype.

6 Conclusion

In this paper, we have presented the new constraint, i.e. the slot filling constraint as the complement of the existing archetypes' built-in constraints. The constraint is important to maintain semantic consistency between the archetype or sub section names containing the slot and the filling archetypes. An ontology-based approach has been proposed to validate the constraint on the development of new archetypes. This approach also shows a different way of utilization of ontology in health from most of the existing ontology employment, which goal is for interoperability purpose.

The evaluation result shows that the method is a better alternative to the current slot fillings condition. Based on the manual decision of the slots filling, the precision, the recall, the accuracy, and the F-measure of the slot filling results using the method are calculated. Compared to the current slot filling condition, the accuracy and F-measure of the slot filling using the method is generally better. The evaluation also shows that the F-measure for the slots contained in a sub section is better than those contained directly in an archetype. However, the accuracy of the method is not much different.

While there are two classes of archetypes, i.e. Section and Composition, can have slots, in the evaluation the method is only applied to some Section archetypes. Further application to all other archetypes, including the Composition archetypes, can be done as part of our future works. Another interesting issue is the application of the approach to another existing clinical ontology, as some archetypes have the binding to them as well.

Acknowledgement

This work is partially supported by the Ministry of National Education of the Republic of Indonesia through the scholarship granted to the first author.

References

1. Bicer, V., Kilic, O., Dogac, A., Laleci, G.B.: Archetype Based Semantic Interoperability of Web Service Message in the Health Care Domain. *International Journal on Semantic Web and Information Systems* 1, 1–23 (2005)
2. Qamar, R., Rector, A.: Semantic Issues in Integrating Data from Different Models to Achieve Data Interoperability. In: *MEDINFO 2007, Brisbane* (2007)
3. Fernández-Breis, J.T., Menárguez-Tortosa, M., Moner, D., Valencia-García, R., Maldonado, J.A., Vivancos-Vicente, P.J., Miranda-Mena, T.G., Martínez-Béjar, R.: An Ontological Infrastructure for the Semantic Integration of Clinical Archetypes. In: Hoffmann, A., Kang, B.-h., Richards, D., Tsumoto, S. (eds.) *PKAW 2006. LNCS (LNAI)*, vol. 4303, pp. 156–167. Springer, Heidelberg (2006)
4. Wollersheim, D., Sari, A., Rahayu, W.: Archetype-based Electronic Health Records: a Literature Review and Evaluation of their Applicability to Health Data Interoperability and Access. *Health Information Management Journal* 38, 7–17 (2009)
5. Leslie, H., Heard, S.: Archetypes 101. In: Westbrook, J., Callen, J. (eds.) *HIC 2006, Sydney* (2006)

6. Beale, T., Heard, S.: Archetypes Definitions and Principles, The openEHR Foundation (2007)
7. Beale, T.: Archetypes: Constraint-based Domain Models for Future-proof Information Systems. In: Baclawski, K., Kilov, H. (eds.) The International Conference on Object Oriented Programming, Systems, Languages and Applications 2002, Seattle, Washington, pp. 1–18 (2002)
8. Beale, T., Heard, S.: Archetype Definition Language ADL 1.4, The openEHR Foundation (2007)
9. Ruotsalo, T., Hyvönen, E.: A Method for Determining Ontology-Based Semantic Relevance. In: Wagner, R., Revell, N., Pernul, G. (eds.) DEXA 2007. LNCS, vol. 4653, pp. 680–688. Springer, Heidelberg (2007)
10. Lombardi, L., Sartori, G.: Concept Similarity: An Abstract Relevance Classes Approach. In: Fum, D., Missier, F.D., Stocco, A. (eds.) The 7th International Conference on Cognitive Modeling, Trieste, pp. 190–195 (2006)
11. Sartori, G., Lombardi, L.: Semantic Relevance and Semantic Disorders. *Journal of Cognitive Neuroscience* 16, 439–452 (2004)
12. Zhang, G., Yu, C., Cai, D., Song, Y., Sun, J.: Research on Concept-Sememe Tree and Semantic Relevance Computation. In: The 20th Pacific Asia Conference on Language, Information and Computation, vol. 20, pp. 398–402 (2006)
13. Rhee, S.K., Lee, J., Park, M.-W.: Ontology-based Semantic Relevance Measure. In: The First International Workshop on Semantic Web and Web 2.0 in Architectural, Product and Engineering Design, Korea (2007)
14. Ricklefs, M., Blomqvist, E.: Ontology-Based Relevance Assessment: An Evaluation of Different Semantic Similarity Measures. In: Meersman, R., Tari, Z. (eds.) OTM 2008, Part II. LNCS, vol. 5332, pp. 1235–1252. Springer, Heidelberg (2008)
15. Raftopoulou, P., Petrakis, E.: Semantic Similarity Measures: a Comparison Study. Technical University of Crete, Department of Electronic and Computer Engineering (2005)
16. Rada, R., Mili, H., Bicknell, E., Blettner, M.: Development and Application of a Metric on Semantic Nets. *IEEE Transactions on Systems, Man, and Cybernetics* 19, 17–30 (1989)
17. Lee, J.H., Kim, M.H., Lee, Y.J.: Information Retrieval Based on Conceptual Distance in IS-A Hierarchies. *Journal of Documentation* 49, 188–207 (1993)
18. Resnik, P.: Semantic Similarity in a Taxonomy: An Information-Based Measure and Its Application to Problems of Ambiguity in Natural Language. *Journal of Artificial Intelligence Research* 11, 95–130 (1999)
19. Ljubešić, N., Boras, D., Bakarić, N., Njavro, J.: Comparing Measures of Semantic Similarity. In: 30th International Conference on Information Technology Interfaces, Cavtat (2008)
20. Rodriguez, M.A., Egenhofer, M.J.: Comparing Geospatial Entity Classes: An Asymmetric and Context-Dependent Similarity Measure. *International Journal of Geographical Information Science* 18, 229–256 (2004)
21. Liu, M., Shen, W., Hao, Q., Yan, J.: A Weighted Ontology-based Semantic Similarity Algorithm for Web Service. *Expert Systems with Applications* 36, 12480–12490 (2009)
22. Guisheng, Y., Qiuyan, S.: Research on Ontology-Based Measuring Semantic Similarity. In: International Conference on Internet Computing in Science and Engineering, Harbin, pp. 250–253 (2008)
23. Xu, X.-h., Huang, J.-l., Wan, J., Jiang, C.-f.: A Method for Measuring Semantic Similarity of Concepts in the Same Ontology. In: 2008 International Multi-symposiums on Computer and Computational Sciences, Washington DC, pp. 207–213 (2008)
24. IHTSDO: SNOMED-CT, <http://www.ihtsdo.org/snomed-ct/>

25. Sundvall, E., Qamar, R., Nyström, M., Forss, M., Petersson, H., Åhlfeldt, H., Rector, A.: Integration of Tools for Binding Archetypes to SNOMED CT. In: SMCS 2006, Copenhagen, pp. 64–68 (2006)
26. Wouters, C.: A Formalization and Application of Ontology Extraction. PhD Thesis. Department of Computer Science and Computer Engineering, La Trobe University, Melbourne (2005)
27. Mindswap: SMORE - Create OWL Markup for HTML Web Pages, <http://www.mindswap.org/2005/SMORE/>
28. Ontomat Homepage - Annotation Portal, <http://annotation.semanticweb.org/ontomat/index.html>
29. Leacock, C., Chodorow, M.: Combining Local Context with WordNet Similarity for Word Sense Identification. In: WordNet: A Lexical Reference System and its Application, pp. 265–283. MIT Press, Cambridge (1998)
30. IHTSDO: SNOMED Clinical Terms User Guide, http://www.ihtsdo.org/fileadmin/user_upload/Docs_01/Technical_Docs/SNOMED_CT_User_Guide_20080731.pdf
31. Jiang, J.J., Conrath, D.W.: Semantic Similarity Based on Corpus Statistics and Lexical Taxonomy. In: International Conference Research on Computational Linguistics X, Taiwan (1997)
32. Mihalcea, R., Corley, C., Strapparava, C.: Corpus-based and Knowledge-based Measures of Text Semantic Similarity. In: AAAI 2006 Conference, pp. 775–780 (2006)

Time-Decaying Bloom Filters for Efficient Middle-Tier Data Management

Kai Cheng

Faculty of Information Science, Kyushu Sangyo University
3-1 Matsukadai, 2-chome, Higashi-ku, Fukuoka, Japan
chengk@is.kyusan-u.ac.jp

Abstract. Distributed enterprise applications are typically based on a multiple-tier client-server architecture where large volume of data is transferred between tiers frequently. When the amount and frequency of data to be transferred become large, network bandwidth will become a bottleneck and efficient middle-tier data management is critical. In this paper, we propose a semi-persistence model to capture the evolving nature of data in a middle tier data management system. We also propose to use Bloom Filters (BF) as an efficient data structure to maintain the time-sensitive frequency profile of the underlying data items. We first extend the standard Bloom Filters by replacing the bit-vector with an array of counters. We then optimize it by allocating lowest space necessary for each counter to store its value. The preliminary experiments show that the optimized BF achieves considerable improvement on space usage while providing the same results of frequency profile.

1 Introduction

Distributed enterprise applications are commonly adopted in almost every business, large or small. In such applications a multi-tier client-server architecture is typically used to provide adaptability, interoperability for enterprise computing [1]. Because data in this architecture are frequently transferred between tiers, network bandwidth will become a performance bottleneck, especially in large-scale distributed enterprise applications. To alleviate this problem, a middle tier is often used to reduce the data transferring, by storing frequently used data nearby to the data consumers. A middle tier dynamically maintains a set of data items in response to the changing access patterns. Thus data in the middle tier is not necessarily persistent.

According to persistence requirement, one can classify data in current data management researches into two categories. On the one hand, *persistent data* is the main topic of traditional database research. Database management systems are designed to guarantee all data can be stored persistently and consistently all the time. Persistent data will be there at any time unless it has been explicitly moved to somewhere else or removed completely. In other words, persistent data is data with ∞ lifespan. On the other hand, *stream data* has recently gained

most importance. In many applications from IP network management to telephone fraud detection, data arrives in high-speed streams, and queries over those streams need to be processed in an online fashion to enable real-time responses. Data streams pose a serious challenge for data management systems as the traditional DBMS paradigm of set-oriented processing of disk-resident tuples does not apply. Stream data implies data with no persistence requirement and it is with 0 lifespan.

In middle-tier data management systems, data with a finite life bound, or $(0, \infty)$ lifespan, we identify such data as *semi-persistent data*. Efficient management of semi-persistent data requires new framework and technique. Particularly, it is important maintain frequency of access and other profiles of each data item because semi-persistent data management system manages such profiles as first class objects. This means that profiles are explicitly materialized, and managed in the same way that data items are managed. This allows the system to answer trend-report query such as **finding the most frequently used data items** on 'mining unstructured data'.

In this paper, we first propose a semi-persistent data model for modeling data in the middle tier. We also propose to use Bloom Filters (BF) as an efficient data structure to maintain the time-sensitive frequency profile of the underlying data items. We then extend the standard Bloom Filters by replacing the bit-vector with an array of counters and optimize it by allocating lowest space necessary for each counter to store its value. The preliminary experiments show that the optimized BF achieves considerable improvement on space usage while providing the same results of frequency profile.

2 Related Work

As the multi-tier client/server architecture becomes increasingly common in the web, *mid-tier data management*, i.e. cache with database management capabilities has recently gained importance [2,3]. In such a multi-tier architecture, application servers implementing most process logic connect to a backend (central) DBMS, and the latter often becomes the bottleneck of performance. We can improve the backend DBMS by storing frequently used results in the attached cache. This technique can be regarded as "materialized view" of DBMSs. Such kind of application of cache will not contribute to the reduction of network traffic. Thus the objectives are different.

Another relevant work is *multi-level store* of persistent objects. M. Stonebraker [4] proposed the extension of disk-resident database to include multiple storage levels, where time-critical objects reside in main memory, other objects are disk resident, and the remainder occupy tertiary memory. It is possible that more than three levels will be present and that some of these levels will be on remote hardware. Distribution of objects in a storage hierarchy is based on semantic criteria. For example,

**main memory representation: EMP where age \geq 30
and age $<$ 60**

disk representation: age < 30
archive representation age >= 60

The distribution criteria can be changed dynamically by either an application or a database administrator. A specific application can temporarily redistribute instances by temporary redistribution criteria prior to execution. A database administrator can permanently change the distribution by defining permanent distribution criteria. A special program called *vacuum cleaner* is dedicated to enforcement of the distribution criteria as well as management of buffers for data from lower levels.

3 Semi-persistence Model for Middle-Tier Data Management

In this section, we describe the semi-persistence model. After giving a few motivation examples, we shall analyze the insufficiency of the cache model and show the features that a semi-persistence model should have.

To motivate the necessity of a semi-persistence model for the middle-tier data management, we suppose there is a corpus of data items in a middle tier downloaded from remote servers. We also suppose we have an extended SQL-like query language with the capability to query profiles maintained for each data item. For example,

```
SELECT title, author, abstract
FROM Document
WHERE title LIKE '%Ajax%'
ORDER BY ::profile(frequency) DESC
OFFSET 0 Limit 10
```

Here, we ask the system to retrieve the title, author and abstract of the top 10 most frequently accessed documents about 'Ajax'. The ORDER BY clause is extended to include query for profiles and $::profile(\dots)$ is a function provided by the system. The function $::profile(\dots)$ relies on the efficient management of statistics about the underlying data items. We will discuss this in Section 4.

3.1 A Semi-persistence Model

We propose semi-persistence model to model data managed in a middle-tier context. In A data item is associated with a vector of statistics (f_i, r_i, a_i, v_i) .

- f_i : frequency of reference
- r_i : recency of reference
- a_i : age of the data item
- v_i : time to live of the data item

The value f_i for frequency of reference is obtained by a algorithm for efficient monitoring frequency counts. The r_i is the timestamp when the most recent reference occurs. The age a_i represents time since last synchronization. If $a_i \leq v_i$

we say data is living, otherwise we say it is dead. Later we will discuss in details how to set the value of v_i 's adaptively.

The recency statistic can be obtained by timestamp when the data item was referenced. As users may not actually download a data item from the system, instead they may go directly to the origin site. Thus we have to clarify the meaning of data item reference.

4 The Space-Optimized Counting Bloom Filters (SBF)

As so far discussed, the semi-persistence model is dependent heavily on the scalable approach to management of various statistics. In this section, we propose a space-efficient scheme for this purpose. Our scheme is based on a randomized data structure, called Bloom Filters.

The characteristics of a semi-persistence model can be described as follows:

(1) *Relaxed persistence requirement.* Evolving data has relaxed persistence requirement, where data may have a life cycle and may have multiple versions.

(2) *Autonomous life cycle management.* Each data item has a life cycle. Any changes during its life cycle may not be initiated by the user operation. Instead it is managed by the system autonomously in terms of the lifespan of the data item. This is essential for management of the evolving data.

(3) *Retrievable metadata.* To capture the trends of data access, statistics about the data should be maintained efficiently. Furthermore, in order to support trend-report queries, the statistics should be managed as first-level citizen in a semi-persistent data management system.

(4) *Trend-report queries.* One of the most appealing features of the semi-persistence model is that it supports queries for the trends of data evolution as well as the reference trends by the user. Although data warehouse or OLAP systems also support some forms of trend-report queries, the difference lies in the fact that the data under those systems are non-volatile, static, read-only and once committed never over-written or deleted. Whereas data in a semi-persistence model takes the dynamics of the underlying data as useful trends on its own right.

4.1 Bloom Filters for Membership Query

A Bloom Filter is a space-efficient data structures that maintain a very compact inventory of the underlying data, supporting membership queries over a given data set [5]. The space requirements of Bloom filters fall significantly below the information theoretic lower bounds for error-free data structures. This efficiency is at the cost of a small false positive rate (items not in the set have a small constant probability of being listed as in the set), but have no false negatives (items in the set are always recognized as being in the set). Bloom filters are widely used in practice when storage is at a premium and an occasional false positive is tolerable.

The standard Bloom Filters use a bit-vector to hold information about the underlying data set. Initially, all bits in the bit-vector are turned off. Each item

in the set is hashed into several locations of the bit-vector using different hash functions. Bits at these locations are then turned on. In order to find whether an item is a member of the set, we first compute the locations to which the item is mapped according to the same hash functions. We then check if all bits at these locations are on. If so, the answer is yes; otherwise answer no. As a bit can be set by other items due to the hash collisions, false positives are possible.

Standard Bloom Filters have several limitations when extended to deploy in other applications. First, Bloom Filter does not support deletes as simply turning off the corresponding bits may introduce false negative errors (some bits of other items, although still in the set, may be turned off). Second, it is not suitable for dealing with multiset (set with duplicates), where multiplicities of items should be reported. To extend standard Bloom Filters several variants have been developed [6,7,8,9].

4.2 Time-Decaying Counters

In a previous work, we have proposed Time-decaying Bloom Filters (TBF), an improvement of Counting Bloom Filters[9] by employing the exponentially aging model [10]. This is done by introducing *time-decaying counter* or *decaying counter* for short, based on some form of *decay functions*. In [7], a number of decay functions are given, such as exponential decay, sliding window decay, polynomial decay, poly-exponential decay and choral and polygonal decay.

A *time-decaying counter* or simply *decaying counter* is a counter whose value decays periodically. The fashion of how a counter decays with time is determined by a special non-increasing, non-negative function $\phi(t)$, called *time-decaying function (tdf)*. A tdf function should satisfy the following conditions:

1. $\phi(0) = 1$
2. $\phi(t)$ is non-increasing
3. $0 \leq \phi(t) \leq 1$ for all $t \geq 0$

Let f_e be the frequency count of item e in S up to the current time, t_n . Let $f_e^{(i)}$ be the frequency count of item e up to the time t_i since t_{i-1} . That is, $f_e = \sum_i f_e^{(i)}$. Then, the decayed value of the counter for item e is

$$f_e^* = \sum_{i=1}^n \phi(t_i) \cdot f_e^{(i)}$$

Cohen, E. et al [11] give a variety of tdf functions. In this paper, we adopt *exponential tdf*, one of the most widely used tdf's in practice.

$$\phi(t) = \lambda^{t/T}, \quad 0 \leq \lambda < 1, \quad T > 0$$

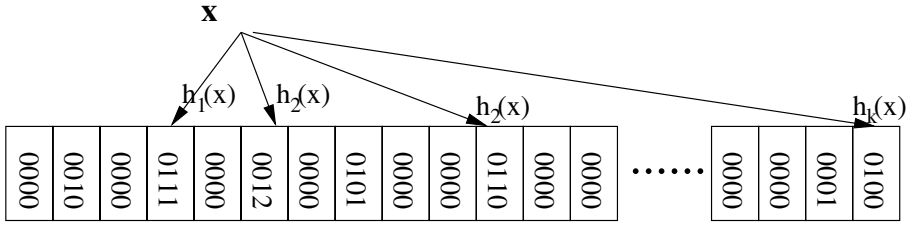


Fig. 1. Basic Time-Decaying Bloom Filters

T and λ are parameters that control the behavior of how fast the counter will decay. T is the period between two consecutive applications of ϕ . In other words, it controls the granularity of time-sensitivity. A time period of T time units is referred to as an *epoch*. The parameter $\lambda \in [0, 1]$, called *exponential decaying factor* or simply *decaying factor*, controls the speed of exponential decaying. In a time-decaying Bloom Filter, basically we use the exponential decay since it is most commonly used decay function in practice. An exponential decay is a function of the the form $g(x) = \lambda^x, \lambda \in [0, 1]$. A decaying counter c is maintained in such a way that its value c_t at time t is obtained as follows:

$$c_t = \lambda \cdot \Delta_t + (1 - \lambda) \cdot c_{t-1} \tag{1}$$

where Δ_t is the fluctuation of c between time $t - 1$ and t . Suppose c is updated at an time interval of T . The value of c_s at time $s = kT$ will be decayed to be $c_s \cdot (1 - \lambda)^{m-k}$ at time $t = mT (k \leq m)$ which is as per the definition of $g(x)$.

Time-decaying counter defined in (1) can not be easily maintained using single counter, because Δ_t and c_{t-1} should be maintained separately. To remedy this, we proposed the following scheme when $\lambda > 0$. Let $c'_t = c_t/\lambda$, we obtain the following equation.

$$c'_t = \Delta_t + (1 - \lambda) \cdot c'_{t-1}, \quad \lambda \in (0, 1] \tag{2}$$

Using this equation, we need only one counter for each value. We can increment the counter directly after the decay function was applied, that when a new $(1 - \lambda)c'_{t-1}$ was obtained. Whenever, say at time t we need the value of c , the result of c_t can be obtained by the following equation.

$$c_t = \lambda \cdot c'_t \tag{3}$$

Note that when the value of a counter can become very large and the exact counts are not necessary, some form of probabilistic counting can be adopted [2][3].

4.3 The Space-Optimized Counting Bloom Filters

To avoid allocating large counters to small values, we optimize the extended BF by using allocating minimal size of memory needed by each counter. To do so, we maintain a free counter pool. The data structures include:

1. A standard Bloom Filter, $BF(m, k)$, with a vector of m bits, k pairwise independent hash functions.
2. A basic $SBF(m, k)$ with k pairwise independent hash functions and m small counters
3. A free counter pool $Pool_j$ with $m/2$ free counters
4. A lookup table (Lookup)

First, a standard Bloom Filters is used to enable quick membership query. The basic SBF with small counters hold frequency counts for most "cold" data items with smaller values. For few "hot" items, however, extra counters can be allocated dynamically when necessary. When a small counter in the SBF gets overflowed, we allocate a new counter from the free counter pool (Pool) for the carried digits. A lookup table (Lookup) maintains the bi-directional links between the overflowed counters and the extra counters. A counter in the SBF can extend to a linked counter list to represent much larger values.

To insert a new item q into the optimized SBF filter, we simply insert the item into SBF, incrementing each of its counters $h_{1,i}(i = 1, \dots, k)$ by 1 unless the counter becomes overflow. If ALL these counters get overflowed, and if there is an extra counter for this SBF counter, increment that counter by 1; otherwise allocate an extra counters for the SBF counter. Repeat the process if a extra counter gets overflow too.

To query for the frequency count of an item, say, q , we first check if q is recorded in BF. If not, return 0; otherwise, we check the SBF and the lookup table to construct the whole counter for q from higher significant bits to lower ones. In Fig 2, q is recorded in the basic BF, and there are 2 extra counters with

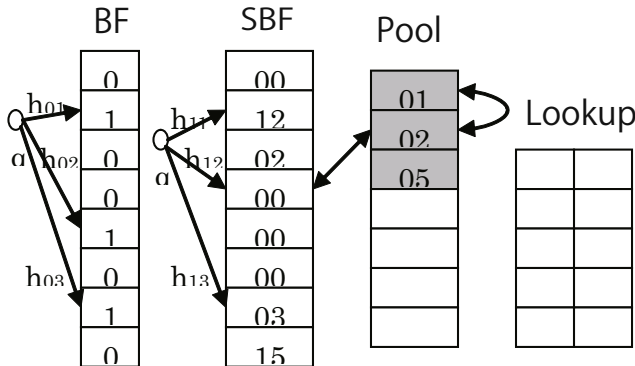


Fig. 2. Space-Optimized Counting Bloom Filters

values of 01 and 02 respectively. Assuming all are 4-bit counters, the frequency count for q by the SBF counter $C_{h_{12}}(q)$ is $256 \times 1 + 2 \times 16 + 0 = 288$.

The usage of web is well known to be quite biased, with a small fraction of popular sites getting very high hits, while the rest and also the majority are rarely used. As the values of counters vary significantly from very small values for most "cold" pages to thousands of hits for a few "hot" pages it is not suitable to allocate the bits to count each of these items.

To avoid allocating large counters for many "cold" pages, we optimize the basic SBF as well as TMF using a dynamic representation of large counters. The optimization is done in the following way:

1. Allocate a small counter to each cell, for example 4 bits or 8 bits.
2. Introduce a *free counter pool (Pool)* and dynamically allocate to large counters, each may be 16 or 32 bits. The *Pool* can also be organized as a hierarchy with counters of different length.
3. Allocate a flag bit for each small counter to indicate whether it contains the real frequency value or just an index to a large counter. Its initial value is on or 1 indicating real value contained. When a small counter becomes overflow, turn off the flag bit and allocate a free large counter from the *Pool* and record its index in the small counter.

The query process will begin with a flag bit check. If it is on, return the value in small counter as the frequency count. Otherwise use the value as an index to find the large counter and return value of that counter.

5 Preliminary Experiments

We evaluate the optimization of the proposed Space-Optimized Counting Bloom Filters. The data set we will use is proxy web access logs of 12 days from June 13 to 24, 2003 provided by IRLCache¹. The sanitized cache access logs contain 541.4 MB data in gzipped form. We preprocess the data set by excluding accesses without a success status code 200 and keep requests of htmltext mime type. The cleaned data set consists of as many as total 2,626,434 requests in total and 1,027,679 distinct requests. Fig. 3 shows that it follows Zipf-like distribution.

First we compare the basic TBF, or B-TBF with its optimization called O-TBF, with different values of λ . The result is shown in Figure 4. With the increase of λ , the number of large counters in use increase. Thus more memory space is consumed.

We then evaluate the the basic TBF, or B-TBF with its optimization called O-TBF, with different window sizes. The result is shown in Figure 5. With the increase of window size, the number of large counters in use increase and more memory space is required.

From Figure 4 and Figure 5, we see the optimization by allocating large counters dynamically on demand is effective.

¹ <http://www.ircache.net/>

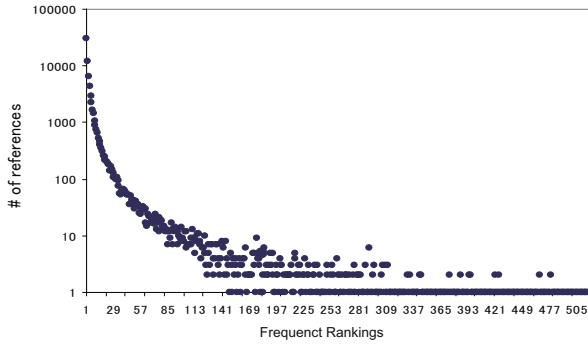


Fig. 3. The workload with a Zipf-like distribution

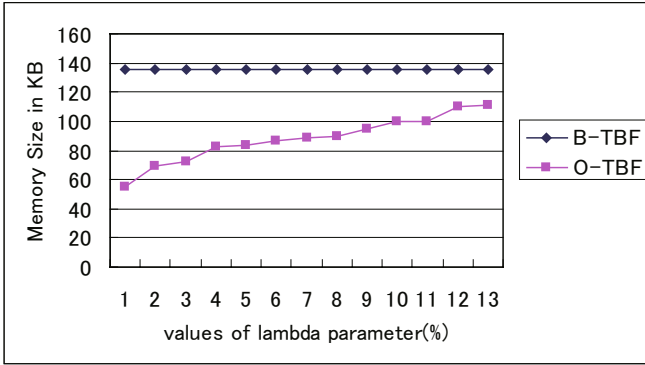


Fig. 4. Comparison of the basic TBF and its optimization

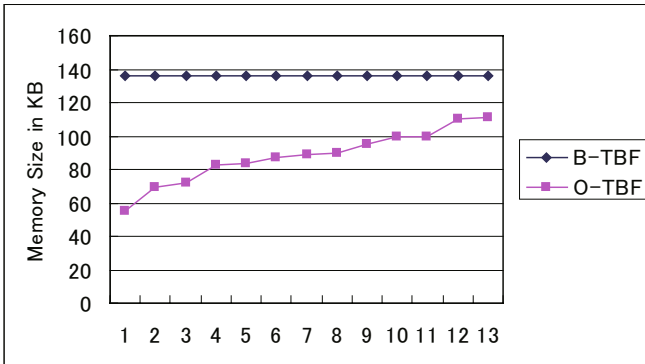


Fig. 5. Comparison of the basic TBF and its optimization

6 Concluding Remarks

To capture the evolving nature of data in middle-tier data management, we proposed a semi-persistence model, which enhances the cache approach by features like trend-report queries as well autonomous life cycle management. The proposed model relies on efficient management of statistics of the underlying data over time. We proposed Time-Decaying Bloom Filters as an efficient data structure to manage time-sensitive access patterns of the underlying data. As a future direction, a suitable scheme for adaptive management of the time to live value of v_i for each element is importance. Results from cache coherency maintenance can be used but more semantic factors should be considered.

References

1. Cooper, B.F., Sample, N., Franklin, M.J., Olshansky, J., Shadmon, M.: Middle-tier extensible data management. *World Wide Web* 4, 209–230 (2001)
2. Luo, Q., Krishnamurthy, S., Mohan, C., Piraahesh, H., Woo, H., Lindsay, B.G., Naughton, J.F.: Middle-tier database caching for e-busines. In: *Proceedings ACM SIGMOD Conference on Management of Data (SIGMOD 2002)*, Madison, Wisconsin, pp. 600–611 (2002)
3. Team, T.: Mid-tier caching: The timesten approach. In: *Proceedings ACM SIGMOD Conference on Management of Data (SIGMOD 2002)*, Madison, Wisconsin, pp. 588–593 (2002)
4. Stonebraker, M.: Managing persistent objects in a multi-level store. In: *Proceedings ACM SIGMOD Conference on Management of Data (SIGMOD 1991)*, Denver, CO, pp. 2–11 (1991)
5. Bloom, B.H.: Space/time trade-offs in hash coding with allowable errors. *Commun. ACM* 13, 422–426 (1970)
6. Bonomi, F., Mitzenmacher, M., Panigrahy, R., Singh, S., Varghese, G.: An improved construction for counting bloom filters. In: Azar, Y., Erlebach, T. (eds.) *ESA 2006*. LNCS, vol. 4168, pp. 684–695. Springer, Heidelberg (2006)
7. Cohen, S., Matias, Y.: Spectral bloom filters. In: *SIGMOD 2003*, pp. 241–252. ACM Press, New York (2003)
8. Deng, F., Rafiei, D.: Approximately detecting duplicates for streaming data using stable bloom filters. In: *SIGMOD 2006*, pp. 25–36. ACM Press, New York (2006)
9. Fan, L., Cao, P., Almeida, J., Broder, A.Z.: Summary cache: a scalable wide-area web cache sharing protocol. *IEEE/ACM Trans. Netw.* 8, 281–293 (2000)
10. Cheng, K., Xiang, L.: Efficient web profiling by time-decaying bloom filters. *DBSJ Letters* 4, 137–140 (2005)
11. Cohen, E., Strauss, M.: Maintaining time-decaying stream aggregates. In: *PODS 2003: Proceedings of the twenty-second ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pp. 223–233. ACM Press, New York (2003)
12. Morris, R.: Counting large numbers of events in small registers. *Commun. ACM* 21, 840–842 (1978)
13. Whang, K.Y., Vander-Zanden, B.T., Taylor, H.M.: A linear-time probabilistic counting algorithm for database applications. *ACM Trans. Database Syst.* 15, 208–229 (1990)

Soft Decision Making for Patients Suspected Influenza

Tutut Herawan^{1,2} and Mustafa Mat Deris¹

¹ Faculty of Information Technology and Multimedia
Universiti Tun Hussein Onn Malaysia, Johor, Malaysia

² Department of Mathematics Education
Universitas Ahmad Dahlan, Yogyakarta, Indonesia
tutut81@uad.ac.id, mmustafa@uthm.edu.my

Abstract. Computational models of the artificial intelligence such as soft set theory have several applications. Parameterization reduction under soft set theory can be considered as a technique for medical decision making. One possible application is the decision making of patients suspected influenza. In this paper, we present the applicability of soft set theory for decision making of patients suspected influenza. The proposed technique is based on maximal supported objects by parameters. At this stage of the research, results are presented and discussed from a qualitative point of view against recent soft decision making techniques through an artificial dataset.

Keywords: Soft set theory; Decision making; Influenza artificial dataset.

1 Introduction

The reliance on information obtained from databases is very critical and important. Almost in every part of our life, there are lots of instances where we have either direct or indirect dealing with databases. One aspect that database plays an important role is in the field of decision making. Inputs obtained from data are stored in terms of records and attributes in databases do contribute a lot in the process of decision making. To this, one practical problem is faced: for a particular property, whether all the attributes in the attributes set are always necessary to preserve this property [1]. Soft set theory [2], proposed by Molodtsov in 1999, is a new general method for dealing with uncertain data. We note that the soft set is designed to replace a Boolean-valued information system. The research described in this paper is a part of our short term effort in applying soft set theory in order to make a decision and further grouping of patients under certain symptoms of influenza. Since soft set theory and soft decision making techniques are not yet widely known, we start from a tutorial introduction. The notions are first introduced by means of simple examples and later formalized. The second part of this research elucidates a soft decision making technique in significant detail. We present a problem of decision making through an artificial data of patient suspected influenza. The data contains 50 patients with 7 symptoms. Using soft set and maximal symptoms co-occurrences in patients, we explore how soft decision making technique can be used to reduce the number of dispensable symptoms and further make a correct decision. This technique may potentially

contribute to lowering the complexity of medical decision making without loss of original information.

The rest of this paper is organized as follows. Section 2 describes the notion of information system. Section 3 describes the theory of soft set. Section 4 describes soft Parameterization reduction and decision making under soft set. Section 5 describes soft decision making using maximal supported objects by parameters. Section 6 describes an application of soft set theory for decision making and grouping patients suspected influenza. Finally, the conclusion of this work is described in section 7.

2 Information System

The syntax of information systems is very similar to relations in relational data bases. Entities in relational databases are also represented by tuples of attribute values. An *information system* is a quadruple $S = (U, A, V, f)$, where $U = \{u_1, u_2, \dots, u_{|U|}\}$ is a non-empty finite set of objects, $A = \{a_1, a_2, \dots, a_{|A|}\}$ is a non-empty finite set of attributes, $V = \bigcup_{a \in A} V_a$, V_a is the domain (value set) of attribute a , $f : U \times A \rightarrow V$ is an information function such that $f(u, a) \in V_a$, for every $(u, a) \in U \times A$, called information (knowledge) function. An information system is also called a knowledge representation systems or an attribute-valued system. In an information system $S = (U, A, V, f)$, if $V_a = \{0,1\}$, for every $a \in A$, then S is called a *Boolean-valued information system*. In many medical information systems, there is an outcome of classification that is known. This *a posteriori* knowledge is expressed by one (or more) distinguished attribute called decision attribute; the process is known as *supervised learning*. An information system of this kind is called a decision system. A *decision system* is an information system of the form $D = (U, A \cup \{d\}, V, f)$, where $d \notin A$ is the decision attribute. The elements of A are called condition attributes. A simple example of a decision table can be found in Table 1.

Example 1. A decision system of 6 patients with six symptoms (conditions) and a decision.

Table 1. An example of a decision system

U	Fever	Cough	Runny Nose	Lethargic	Tired	Headache	Flu
p_1	yes	yes	no	yes	no	yes	yes
p_2	no	no	yes	no	yes	no	no
p_3	yes	yes	no	yes	no	yes	yes
p_4	yes	no	no	no	yes	no	no
p_5	yes	yes	yes	no	no	yes	yes
p_6	no	no	no	yes	yes	no	no

3 Soft Set Theory

Throughout this section U refers to an initial universe, E is a set of parameters, $P(U)$ is the power set of U .

Definition 2. (See [2].) A pair (F, E) is called a soft set over U , where F is a mapping given by

$$F : E \rightarrow P(U).$$

In other words, a soft set over universe U is a parameterized family of subsets of universe U . For $e \in E$, $F(e)$ may be considered as the set of e -elements of the soft set (F, E) or as the set of e -approximate elements of the soft set. Clearly, a soft set is not a (crisp) set.

Definition 3. (See [3].) The class of all value sets of a soft set (F, E) is called value-class of the soft set and is denoted by $C_{(F,E)}$.

Clearly $C_{(F,E)} \subseteq P(U)$.

Example 4. Let we consider a soft set (F, E) which describes the “conditions of patients suspected influenza” that a hospital is considering to make a decision. The six influenza symptoms, i.e., fever, respiratory, nasal discharges, cough, headache and sore throat are adopted from [4] and one symptom added is lethargic. Suppose that there are six patients in the hospital under consideration,

$$U = \{h_1, h_2, h_3, h_4, h_5, h_6\},$$

and E is a set of decision parameters

$$E = \{e_1, e_2, e_3, e_4, e_5, e_6, e_7\},$$

where

- e_1 stands for the parameter “fever”,
- e_2 stands for the parameter “respiratory”,
- e_3 stands for the parameter “nasal discharges”,
- e_4 stands for the parameter “cough”,
- e_5 stands for the parameter “headache”
- e_6 stands for the parameter “sore throat”
- e_7 stands for the parameter “lethargic”.

Consider the mapping $F : E \rightarrow P(U)$ given by “patients(\cdot)”, where (\cdot) is to be filled in by one of parameters $e \in E$. Suppose that

$$\begin{aligned} F(e_1) &= \{p_1, p_3, p_4, p_5, p_6\}, \\ F(e_2) &= \{p_1, p_2\}, \end{aligned}$$

$$\begin{aligned}
 F(e_3) &= \{p_1, p_2, p_4\}, \\
 F(e_4) &= \{p_1\}, \\
 F(e_5) &= \{p_3, p_4\}, \\
 F(e_6) &= \{p_2, p_4\}, \\
 F(e_7) &= \{p_1, p_3, p_5, p_6\}.
 \end{aligned}$$

Therefore, $F(e_1)$ means “patients suffer fever”, whose functional value is the set $\{p_1, p_3, p_4, p_5, p_6\}$. Thus, we can view the soft set (F, E) as a collection of approximations as below

$$(F, E) = \left\{ \begin{array}{l} \text{fever} = \{p_1, p_3, p_4, p_5, p_6\}, \\ \text{respiratory} = \{p_1, p_2\}, \\ \text{nasal discharges} = \{p_1, p_2, p_4\}, \\ \text{cough} = \{p_1\}, \\ \text{head ache} = \{p_3, p_4\}, \\ \text{sore throat} = \{p_2, p_4\}, \\ \text{lethargic} = \{p_1, p_3, p_5, p_6\} \end{array} \right\}$$

Table 2. Tabular representation of a soft set in the above example

U	e_1	e_2	e_3	e_4	e_5	e_6	e_7
p_1	1	1	1	1	0	0	1
p_2	0	1	1	0	0	1	0
p_3	1	0	0	0	1	0	1
p_4	1	0	1	0	1	1	0
p_5	1	0	0	0	0	0	1
p_6	1	0	0	0	0	0	1

Each approximation has two parts, a predicate e and an approximate value set p . For example, for the approximation “fever = $\{p_1, p_3, p_4, p_5, p_6\}$ ”, we have the predicate name of patients with fever and its value set is $\{p_1, p_3, p_4, p_5, p_6\}$.

The relation between a soft set and a Boolean-valued information system is given in the following proposition.

Proposition 3. *If (F, E) is a soft set over the universe U , then (F, E) is a binary-valued information system $S = (U, A, V_{\{0,1\}}, f)$.*

Proof. Let (F, E) be a soft set over the universe U , we define a mapping

$$F = \{f_1, f_2, \dots, f_n\},$$

where

$$f_i : U \rightarrow V_i \text{ and } f_i(x) = \begin{cases} 1, & x \in F(e_i) \\ 0, & x \notin F(e_i) \end{cases}, \text{ for } 1 \leq i \leq |A|.$$

Hence, if $A = E$, $V = \bigcup_{e_i \in A} V_{e_i}$, where $V_{e_i} = \{0,1\}$, then a soft set (F, E) can be considered as a binary-valued information system $S = (U, A, V_{\{0,1\}}, f)$. □

From Proposition 3, it is easily to understand that a binary-valued information system can be represented as a soft set. Thus, we can make a one-to-one correspondence between (F, E) over U and $S = (U, A, V_{\{0,1\}}, f)$.

4 Parameterization Reduction and Decision Making

In this section, we present the existing techniques on soft reduction and decision making techniques in [5], [6] and [7]. The purpose of this analysis is to present the comparison in the previous techniques and how our proposed technique will provide an alternative way for soft decision making. Suppose we have a soft set (F, E) over universe U with the Boolean representation as displayed in Table 2. Let $f_E(p_i) = \sum_j p_{ij}$, where h_{ij} are the entries of symptoms in the Boolean-table of (F, E) . Obviously, based on the Table 3 below, p_1 will be the patient of first choice for selecting since it is clearly shown that the total symptoms, $f_E(p_1) = 5$, is the maximum choice value. In this case, we say that p_1 is the optimal decision of patient having influenza. Patients p_4 , (p_2, p_3) and (p_5, p_6) then referred as first, second and third sub-optimal decisions, respectively.

Table 3. An example of Boolean table of representation of a soft set (F, E)

U	e_1	e_2	e_3	e_4	e_5	e_6	e_7	$f_E(p_i)$
p_1	1	1	1	1	0	0	1	5
p_2	0	1	1	0	0	1	0	3
p_3	1	0	0	0	1	0	1	3
p_4	1	0	1	0	1	1	0	4
p_5	1	0	0	0	0	0	1	2
p_6	1	0	0	0	0	0	1	2

The problem of parameterization reduction and decision making through this view, then is to reduce the number of parameter (symptoms) which preserved the consistency of optimal and sub-optimal decisions.

4.1 Parameterization Reduction of Maji *et al.* [5]

Maji *et al.* presented a reduction of soft sets and its applications in a decision making problem, which can be briefly explained as follows. The most optimal decision derived by Maji will be only be deduced by identifying the rough set-based reduction set first. This can by maintained using a partition on U based on the indiscernibility relation on a set of attributes in rough set theory [8]. As for our example based on Table 3, the partition induced by the set of all attributes $E = \{e_1, e_2, e_3, e_4, e_5, e_6, e_7\}$, denoted by U/E is given by

$$\{\{p_1\}, \{p_2\}, \{p_3\}, \{p_4\}, \{p_5, p_6\}\}.$$

Therefore, if any attribute to be deleted, then the partition induced by any subset $Q \subset E$ can only be considered as a reduct if partition induced by Q/E is still $\{\{p_1\}, \{p_2\}, \{p_3\}, \{p_4\}, \{p_5, p_6\}\}$. For example, $Q \subset E$, where $Q = \{e_1, e_5, e_6, e_7\}$ and Q/E is still

$$\{\{p_1\}, \{p_2\}, \{p_3\}, \{p_4\}, \{p_5, p_6\}\}.$$

Therefore, for Maji’s techniques, the set of attributes, Q can also be considered as a reduct of E .

Table 4. A Boolean table after reduction of $Q \subset E$

U	e_1	e_5	e_6	e_7	$f_E(p_i)$
p_1	1	0	0	1	2
p_2	0	0	1	0	1
p_3	1	1	0	1	3
p_4	1	1	1	0	3
p_5	1	0	0	1	2
p_6	1	0	0	1	2

However, for Maji’s techniques, it can be seen from Table 4, that the optimal decision of patients is no longer p_1 . Patient p_1 has instead become sub-optimal patient together with that of p_5 and p_6 . The maximum value of parameters after reduction, i.e., $f_E(p_3) = f_E(p_4) = 3$, is now the maximum choice. Thus, the new optimal decision of patient in this case will be both p_3 and p_4 . It is very obvious, selection for optimal patient will yield inconsistent result as shown from one example of reduct $Q \subset E$, as in the above example. Even if we are to consider the selection of the next optimal patient, in the case of Q , it will also provide inconsistent result. If for another example, $G \subset E$, $G = \{e_1, e_3, e_5, e_6, e_7\}$ and G/E is still producing the partition like $\{\{p_1\}, \{p_2\}, \{p_3\}, \{p_4\}, \{p_5, p_6\}\}$. Therefore G/E can also be considered as a reduct of E .

Table 5. A Boolean table after reduction of $G \subset E$

U	e_1	e_3	e_4	e_5	e_7	$f_E(p_i)$
p_1	0	0	1	1	1	3
p_2	1	0	0	0	1	2
p_3	0	1	1	1	0	3
p_4	1	1	0	1	1	4
p_5	0	0	1	1	0	2
p_6	0	0	1	1	0	2

As can be seen clearly again, now the optimal decision of patients has now changed to p_4 , for this new set of reduct G which is $G \subset E$. Again the issue of changing the value of the optimal decision of patient due to the change of maximum value of different set of reduct will is very critical and of our main concern here. This will not be an ideal technique of Maji's proposal, since it does not provide consistent result in terms of different optimal decision of patients that is suggested. Furthermore, Maji did not discuss the sub-optimal patient that can be derived from the reduced set. And as we have explored, the sub-optimal patient if addressed by Maji, then again different set of reduct will yield different sub-optimal patient, thus again not able to provide consistency in selecting sub-optimal patient. As pointed out earlier in this paper, that the main intention of reduction is to reduce complexity in decision making but at the same time, still able to provide consistencies in decision making.

4.2 Parameterization Reduction of Chen *et al.* [6]

For the sake of comparison results from different proposed techniques, let us again use the Table 3 as example for analysis purposes. The main purpose of Chen techniques, was to maintain consistency in optimal value of decision thus improving the process of decision making from Maji. Chen has defined $f_E(p_i) = \sum_j p_{ij}$ where h_{ij} are the entries in the Boolean-table of (F, E) and M_E denoted for collection of

Table 6. A Boolean table after parameterization reduction

U	e_1	e_2	e_3	$f_E(p_i)$
p_1	1	1	1	3
p_2	0	1	1	2
p_3	1	0	0	1
p_4	1	0	1	2
p_5	1	0	0	1
p_6	1	0	0	1

objects in U which has the maximum value of f_E . Chen has in fact defined a *dispensable set* $A \subset E$ if only if $M_{E \setminus A} = M_E$. Clearly, parameter reduction of Chen has been able to provide consistency in optimal object's decision.

Let assume J , where $J \subseteq E$ and $J = \{e_4, e_5, e_6, e_7\}$. Hence, $M_{E \setminus J} = M_E$ is true, and therefore J is dispensable and we are left with $S = \{e_1, e_2, e_3\}$ and $S \subseteq E$. In the Table 6 above, p_1 is still the optimal decision of patient, thus parameter reduction thus still up-hold consistency in selecting optimal patient as opposed to Maji's attribute reduction. However, Chen's parameter reduction is only able to provide in consistency in selecting optimal decision but not in the sub-optimal decision. As can be seen from the Table 6 above, that the sub-optimal patients are p_2 and p_4 . Meanwhile, from Table 3, p_4 is the only patient to be referred clearly as the sub-optimal patient. Therefore parameter reduction of Chen does provide consistency in selecting optimal patient but falls short in providing consistency for the selection of sub-optimal patients.

4.3 Normal Parameterization Reduction of Kong et al. [7]

The main patientives of Kong's normal parameter reduction is to provide consistency in selecting optimal and sub-optimal objects for any reduced set that conforms to the original decision provided. Kong has maintained the same the partitions of objects by defining *indiscernibility relation* $IND(A)$, for $A \subset E$ as follow

$$IND(A) = \{(p_i, p_j) \in U \times U : f_A(p_i) = f_A(p_j)\}.$$

The *decision partition* of U generated by $IND(E)$ is defined as

$$C_E = \{\{p_1, \dots, p_i\}_{f_1}, \{p_{i+1}, \dots, p_j\}_{f_2}, \dots, \{p_k, \dots, p_n\}_{f_n}\}.$$

In the case that for $A \subset E$, if $f_A(p_1) = f_A(p_2) = \dots = f_A(p_n)$ implies $C_E = C_{E \setminus A}$, then A is called *dispensable set*. For this definition, Kong has termed $E - A$ as normal parameter reduction.

Based on on the Boolean table in Table 3, the decision partition induced will be $C_E = \{\{p_1\}_{f_5}, \{p_4\}_{f_4}, \{p_2, p_3\}_{f_3}, \{p_5, p_6\}_{f_2}\}$, and p_1 is the optimal decision and p_4 will be

Table 7. A Boolean table after normal parameter reduction

U	e_1	e_2	e_3	e_4	e_5	$f_E(p_i)$
p_1	1	1	1	1	0	4
p_2	0	1	1	0	0	2
p_3	1	0	0	0	1	2
p_4	1	0	1	0	1	3
p_5	1	0	0	0	0	1
p_6	1	0	0	0	0	1

the first sub-optimal decision following by p_2, p_3 and p_5, p_6 as the second and third sub-optimal patients, respectively. To do a normal parameter reduction from Table 3, since $f_{\{e_6, e_7\}}(h_1) = f_{\{e_6, e_7\}}(h_2) = \dots = f_{\{e_6, e_7\}}(h_6) = 1$, thus $Z = \{e_6, e_7\} \subset E$, is dispensable since the decision partition has not changed, that is $C_E = C_{E \setminus Z}$. By deleting parameter $Z = \{e_1, e_4\} \subset E$, we will have as what is been displayed in Table 7. And Kong has successfully shown that the optimal and sub-optimal decisions, p_1 and p_4 , $\{p_2, p_3\}$, $\{p_5, p_6\}$ respectively, thus maintaining consistency in patients decision after the reduction.

5 Soft Decision Making Using Maximal Supported Objects

Throughout this sub-section the pair (F, E) refers to the soft set over the universe U representing a Boolean-valued information system $S = (U, A, V_{\{0,1\}}, f)$.

Definition 4. Let (F, E) be a soft set over the universe U and $u \in U$. A parameter co-occurrence set of an object u can be defined as

$$\text{coo}(u) = \{e \in E : f(u, e) = 1\}.$$

Obviously, $\text{Coo}(u) = \{e \in E : F(e) = 1\}$.

Definition 5. Let (F, E) be a soft set over the universe U and $u \in U$. Support of an object u is defined by

$$\text{supp}(u) = \text{card}(\{e \in E : f(u, e) = 1\}).$$

Definition 6. Let (F, E) be a soft set over the universe U and $u \in U$. An objects u is said to be maximally supported by a set of all parameters E , denoted by $\text{Msupp}(u)$ if

$$\text{supp}(u) > \text{supp}(v), \quad \forall v \in U \setminus \{u\}.$$

Based on Definition 6, we can make supported (ranked) ordered objects according their support value as

$$U_1 > U_2 > \dots > U_n,$$

where $U_i \subseteq U$ and $U_i = \{u \in U : u \text{ is } i\text{-th maximal supported by } E\}$, for $1 \leq i \leq n$.

Thus, U_i is a collection of objects in U having the same support, i.e., objects of the same support of are grouped into the same class. Obviously $U = \bigcup_{1 \leq i \leq n} U_i$ and $U_i \cap U_j = \emptyset$, for $i \neq j$. In other word, a collection of $U/E = \{U_1, U_2, \dots, U_n\}$ is a decision partition of U , so called *cluster decision* of U .

Definition 7. Let (F, E) be a soft set over the universe U and $A \subset E$. A is said to be indispensable if $U/A = U/E$. Otherwise, A is said to be dispensable.

Based on Definition 7, we can reduce the number of parameters without changing the optimal and sub-optimal decisions.

Definition 8. For soft set (F, E) over the universe U and $A \subseteq E$. A is reduction of E if only if A is indispensable and $\text{supp}_A(u) = \text{supp}(v)$, for every $u, v \in U$.

Definition 9. For soft set (F, E) over the universe U and $u \in U$. An object u will be the optimal decision if u is maximally supported by E .

Example 10. As for example, the following will be the co-occurrence objects derived from Table 3.

$$\begin{aligned} \text{coo}(p_1) &= \{e_1, e_2, e_3, e_4, e_7\}, \text{coo}(p_2) = \{e_2, e_3, e_5\}, \text{coo}(p_3) = \{e_1, e_5, e_7\}, \\ \text{coo}(p_4) &= \{e_1, e_3, e_5, e_7\}, \text{coo}(p_5) = \{e_1, e_7\} \text{ and } \text{coo}(p_6) = \{e_1, e_7\}. \end{aligned}$$

Thus, support of each object is given as follow

$$\text{supp}(p_1) = 5 > \text{supp}(p_4) = 4 > \text{supp}(p_2) = \text{supp}(p_3) = 3 > \text{supp}(p_5) = \text{supp}(p_6) = 2.$$

Based from the Definition 7, the partition contain 4 clusters, i.e.,

$$\{\{p_1\}, \{p_4\}, \{p_2, p_3\}, \{p_5, p_6\}\},$$

where it is arrange in descending order of support value.

As noted that, the first maximal supported patient is p_1 , where $\text{supp}(p_1) = 5$ and as been defined by Definition 9, p_1 is the optimal decision. p_4 can be considered as the second maximal supported patient based on the support which the next highest, i.e., $\text{supp}(p_4) = 4$. And it is also noted that p_2, p_3 and p_5, p_6 can be the second and third maximal supported patients, since $\text{supp}(p_2) = \text{supp}(p_3) = 3$ and $\text{supp}(p_5) = \text{supp}(p_6) = 2$, respectively.

To elaborate Definition 8, let $A = \{e_2, e_3, e_5, e_6, e_7\}$. Then we will obtain

$$U / A = \{\{p_1\}, \{p_4\}, \{p_2, p_3\}, \{p_5, p_6\}\} = U / E,$$

and since $E \setminus A = \{e_1, e_4\}$, then

$$\text{supp}_{E \setminus A}(p_1) = \text{supp}_{E \setminus A}(p_2) = \dots = \text{supp}_{E \setminus A}(p_6) = 1.$$

Therefore, A is parameter reduction of E and we can now delete attribute $A = \{e_1, e_4\}$. Note that, by deleting A , we now have $U / A = \{\{p_1\}, \{p_4\}, \{p_2, p_3\}, \{p_5, p_6\}\}$, which is still the same decision partition as in U / E . Also in this case, the maximum supported objects are still maintained. As in [5] and [6], only the issue of optimal was addressed, but in our paper, any set of reduct that conforms to our rule of reduct will still provide the same optimal and sub optimal. By comparing optimal and sub-optimal decision from our proposed technique with normal parameterization reduction from [7], also giving the result of the same the optimal and sub-optimal decisions.

Furthermore, our proposed technique confirming that the reduction also provide the right objects for decision making.

6 Application for Patients Suspected Influenza

Yearly influenza epidemics can seriously affect all age groups, but the highest risk of complications occur among children younger than age two, adults age 65 or older, and people of any age with certain medical conditions, such as chronic heart, lung, kidney, liver, blood or metabolic diseases (such as diabetes), or weakened immune systems [9].

6.1 Artificial Dataset

The proposed technique as proposed in Section 5 is elaborate for patient decision making from an artificial dataset of patients suspected influenza. The dataset is shown in Table 8. The table consists of 50 patients with 7 symptoms, i.e., fever (e_1), cough (usually dry) (e_2), headache (e_3), muscle and joint pain (e_4), severe malaise (feeling unwell) (e_5), sore throat (e_6) and runny nose (e_7).

Table 8. An artificial data of patients suspected influenza

Patients	e_1	e_2	e_3	e_4	e_5	e_6	e_7
1	1	0	1	0	0	0	1
2	1	1	0	1	0	1	1
3	1	1	0	1	0	1	1
4	1	0	1	1	1	0	0
5	1	0	0	0	0	0	1
6	1	0	0	0	0	0	1
7	1	0	1	1	1	0	0
8	1	0	0	0	0	0	1
9	1	0	0	0	0	0	1
10	1	0	0	0	0	0	1
11	0	0	0	1	1	1	0
12	1	0	1	1	1	0	0
13	1	0	1	0	0	0	1
14	1	0	1	0	0	0	1
15	1	1	0	1	0	1	1
16	1	1	0	1	0	1	1
17	1	0	0	0	0	0	1
18	1	1	0	1	0	1	1
19	1	0	1	1	1	0	0
20	1	0	0	0	0	0	1
21	1	0	0	0	0	0	1
22	1	0	0	0	0	0	1
23	0	0	0	1	1	1	0
24	0	0	0	1	1	1	0
25	1	0	0	0	0	0	1

Table 8. (continued)

26	0	0	0	1	1	1	0
27	1	0	0	0	0	0	1
28	0	0	0	1	1	1	0
29	1	0	0	0	0	0	1
30	0	0	0	1	1	1	0
31	0	0	1	1	0	0	0
32	1	0	1	0	1	1	0
33	0	0	1	1	0	1	0
34	1	0	0	0	1	0	1
35	0	0	1	1	0	0	0
36	0	0	1	1	0	1	0
37	0	0	1	1	0	0	0
38	1	0	1	0	1	1	0
39	0	0	1	1	0	0	0
40	0	0	1	1	0	0	0
41	0	0	1	1	0	0	0
42	1	0	1	0	1	1	0
43	1	0	1	0	1	1	0
44	1	1	1	1	0	0	1
45	0	0	1	1	0	1	0
46	1	1	1	1	0	0	1
47	1	1	1	1	0	0	1
48	0	0	1	1	0	1	0
49	1	0	0	0	1	0	1
50	0	0	1	1	0	1	0

6.2 Result

After calculating co-occurrences of symptoms in each patient, figure 1 lists the detailed results of dispensable symptoms and reduct obtained. It is clear that symptoms severe malaise (feeling unwell) (e_5) and runny nose (e_7) are indicated as dispensable symptoms. Therefore, reduct of the symptoms is fever (e_1), cough (usually dry) (e_2), headache (e_3), muscle and joint pain (e_4), and sore throat (e_6).

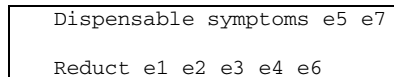


Fig. 1. Reduction of symptoms

After reduction, we can easily make decisions. The decisions are based on the support symptoms of each patient. Patients having the same support are clustered in the same class. There are 9 clusters decision we can made and they are given in

Cluster 1 (fever)	5	6	8	9	10	17	20	21	22	25	27	29	34	49
Cluster 2 (fever headache)	1	13	14											
Cluster 3 (muscleandjoinpain sorethroat)	11	23	24	26	28	30								
Cluster 4 (headache muscleandjoinpain)	31	35	37	39	40	41								
Cluster 5 (fever headache)	4	7	12	19										
Cluster 6 (headache sorethroat)	32	38	42	43										
Cluster 7 (headache muscleandjoinpain sorethroat)	33	36	45	48	50									
Cluster 8 (fever cough muscleandjoinpain sorethroat)	2	3	15	16	18									
Cluster 9 (fever cough headache muscleandjoinpain)	44	46	47											

Fig. 2. Decision clusters

Figure 2. Patient with higher supported symptoms are indicated having flu rather than lower supported symptoms.

7 Conclusion

In this paper, we have presented an application of soft decision making technique for decision making of patients suspected influenza. We have proposed an alternative approach for parameterization reduction and decision making under soft set theory. We have used the co-occurrence of parameters concept in a Boolean-valued information system for defining support of each object by parameters. Based on maximal supported objects, we have defined the notion of indispensable subset and reduct. We have presented an artificial dataset of 50 patients suspected influenza to show that the proposed technique can be used to make a correct decision. We have shown that the results obtained are equivalent with that normal parameterization and decision making problem under soft set theory. For future research, we will investigate the possibility of parameterization reduction and decision making problem under a certain criteria of parameter, e.g., under spatial ordered parameter. With this approach, it is possible that the reduct and decision obtained are different with the proposed technique in this paper.

Acknowledgement

This work was supported by the FRGS under the Grant No. Vote 0402, Ministry of Higher Education, Malaysia.

References

1. Zhao, Y., Luo, F., Wong, S.K.M., Yao, Y.Y.: A General Definition of an Attribute Reduct. In: Yao, J., Lingras, P., Wu, W.-Z., Szczuka, M.S., Cercone, N.J., Ślezak, D. (eds.) RSKT 2007. LNCS (LNAI), vol. 4481, pp. 101–108. Springer, Heidelberg (2007)

2. Molodtsov, D.: Soft Set Theory-First Results. *Computers and Mathematics with Applications* 37, 19–31 (1999)
3. Maji, P.K., Biswas, R., Roy, A.R.: Soft Set Theory. *Computers and Mathematics with Applications* 45, 555–562 (2003)
4. Abbas, K., Mikler, A.R., Gatti, R.: Temporal Analysis of Infectious Diseases: Influenza. In: *The Proceeding of ACM Symposium on Applied Computing (SAC)*, pp. 267–271 (2005)
5. Maji, P.K., Roy, A.R., Biswas, R.: An Application of Soft Sets in a Decision Making Problem. *Compututer and Mathematics with Application* 44, 1077–1083 (2002)
6. Chen, D., Tsang, E.C.C., Yeung, D.S., Wang, X.: The Parameterization Reduction of Soft Sets and its Applications. *Computers and Mathematics with Applications* 49, 757–763 (2005)
7. Kong, Z., Gao, L., Wang, L., Li, S.: The Normal Parameter Reduction of Soft Sets and Its Algorithm. *Computers and Mathematics with Applications* 56, 3029–3037 (2008)
8. Pawlak, Z., Skowron, A.: Rudiments of Rough Sets. *Information Sciences* 177(1), 3–27 (2007)
9. <http://www.who.int/mediacentre/factsheets/fs211/en/index.html>

Personal Identification by EEG Using ICA and Neural Network

Preecha Tangkraingki¹, Chidchanok Lursinsap¹, Siripun Sanguansintukul¹,
and Tayard Desudchit²

¹ Advanced Virtual and Intelligent Computing (AVIC) Research Center,
Department of Mathematics, Faculty of Science,
Chulalongkorn University Bangkok 10330, Thailand

² The Chulalongkorn Comprehensive Epilepsy Program (CCEP),
Faculty of Medicine, Chulalongkorn University, Bangkok 10330, Thailand

Abstract. The problem of identifying a person using biometric data is interesting. In this paper, the uniqueness of EEG signals of individuals is used to determine personal identity. EEG signals can be measured from different locations, but too many signals can degrade the recognition speed and accuracy. A practical technique combining Independent Component Analysis (ICA) for signal cleaning and a supervised neural network for classifying signals is proposed. From 16 EEG different signal locations, four truly relevant locations F_7 , C_3 , P_3 , and O_1 were selected. This selection can identify a group of 20 persons with high accuracy.

Keywords: Electroencephalogram, Independent Component Analysis, Neural Network, Pattern-recognition.

1 Introduction

Biometrics such as fingerprints, retinal or iris scanning, face recognition are actively used for identifications [1]. Cognitive biometrics using brain signals have become interesting as identification tools. To understand how the brain functions, many techniques such as electroencephalography (EEG), magnetoencephalography (MEG), function magnetic resonance imaging (fMRI), and positron emission tomography (PET) have been utilized. Each technique has its own strengths and weaknesses. In this paper EEG has been used to analyze patterns because EEG has a desirable property for excellent time resolution and low cost of instrumentation.

The objectives of this work strive to acquire locations on the scalp (channels) that are the most promising locations for personal identification. Furthermore, the minimum numbers of channel necessary for the identifications would also be explored. Identifications by EEG using Independence Component Analysis technique has been used to identify individual signals from different areas of the brain. Neural network technique was then performed on pattern recognition for identification. Our study is organized as follows. Section 1 is the introduction, using EEG for identifications. Section 2 presents the background work; techniques employed in the experiment such as ICA and neural networks and the

experimental process. Then, the results are displayed and discussed. Finally, the conclusion is included in Section 3.

2 Background Methodology and Experiments

2.1 Background

Electroencephalography (EEG). EEG is the measurement of electrical activity produced by the brain as recorded from electrodes placed on the scalp. These signals are small and combined with other signals such as eye tracking and EMG (electromyography). EEG is surrounded by large electrical potentials from the environment. Brain waves have been categorized into five basic groups according to its frequency 1) Delta (1-4Hz), 2) Theta (4-8Hz), 3) Alpha (8-12 Hz), 4) Beta (12-30Hz), and 5) Gamma (30-50Hz) waves. Although none of these waves are ever emitted alone, the state of consciousness of the individual may result in one frequency being more pronounced than others. Brain waves were first recorded in 1874 by Richard Caton, who connected equipment directly to the cerebral cortex of a rabbit. In 1929, Hans Berger published the first information on scalp-recorded brain waves in humans. The differential input was first amplified by B. H. C. Matthews in 1934 and revolutionized the high-gain amplification of biologic electrical signals, including brain waves. In 1935 Frederic Gibbs, Hallowell Davis, and William Lennox, published the first EEG paper in English dealing with epilepsy in humans. EEG is now extensively used in diagnosing epilepsy and in the study of how the brain functions in both animals and humans.

Independent Component Analysis (ICA). ICA is a member of a class of Blind Source Separation (BSS). The aim of source separation is to recover original signals from known observations where each observation is an unknown mixture of the original signals [2]. In this paper ICA was used to solve the problem of separating multi-channels EEG data into independent sources. We tested the potential usefulness of the ICA algorithm for EEG source decomposition by applying the algorithm to simulated EEG data.

Neural Networks. Neural network is a process paradigm that mimics the structures and functions of the human nervous system. Pattern recognition is an important application which can be implemented using a feed-forward neural network that has been trained accordingly. During the training, the network learns to associate output with input patterns. When the network is used, it identifies the input pattern and tries to output the associated output pattern similar to the way the human brain works.

2.2 Methodology and Experiments

The experimental methodology in this study consists of five processes as shown in Figure 1: Collect EEG signal; Feature extraction using ICA 22 algorithms;

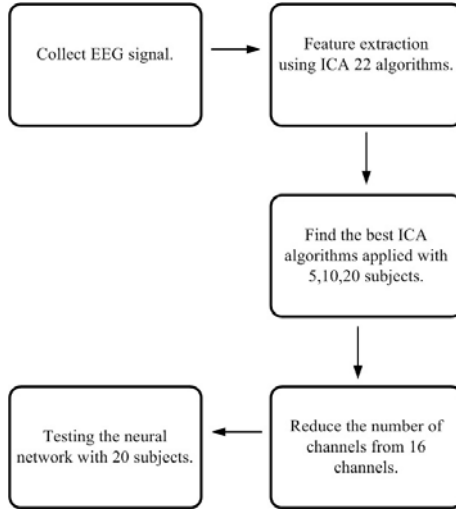


Fig. 1. Experimental diagram process

Find the best ICA algorithms applied with 5,10,20 subjects; Reduce the number of channels from 16 channels; Test the neural network with 20 subjects.

Collect EEG signal. EEG signals were collected from 20 normal patients (eight men and twelve women) from the Chulalongkorn Comprehensive Epilepsy Program (CCEP). The age range of subjects is between 12 and 40 years. The EEG signal was recorded while subjects were resting with their eyes open. The electrodes using 16 gold cups were placed on the scalp. Electrode placement are used to record EEG signals from the locations $FP_1, F_7, T_3, T_5, FP_2, F_8, T_4, T_6, F_3, C_3, P_3, O_1, F_4, C_4, P_4, O_2$ according to 10-20 system shown in Figure 2. In this paper, channels are labeled as ch1, ch2, ..., ch16 representing EEG from location FP_1, F_7, \dots, O_2 as mentioned earlier.

Recording sessions used monopolar montage with reference at the mastoid area. The EEG amplifier was Grass model 8 plus. The sampling rate was 200 Hz. The data was digitized by BMSI board using Stellate Harmony EEG software. The data was exported as EDF (European Data Format). Electromyography (EMG) and electrocardiogram (ECG) signals were initially removed from the EEG signals since they were irrelevant. Finally, there were 3,000 sample data sets of 16 channels collected from 20 subjects. A sample of the raw EEG signals is shown in Figure 3(a).

Feature extractions using ICA. The purpose of this step is to identify the individual signal coming from each channel using 22 different ICA algorithms. The signal data obtained from the previous step were processed using ICALAB. The ICALAB for signal processing is the package for MATLAB that implements

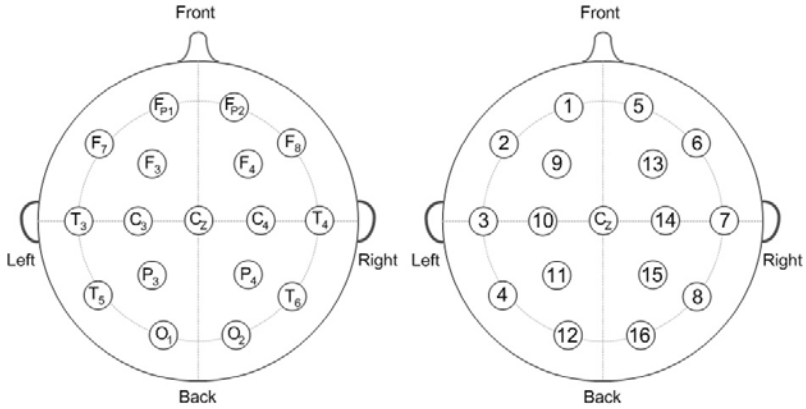


Fig. 2. The locations of electrode placements on the scalp using 10-20 system [3]

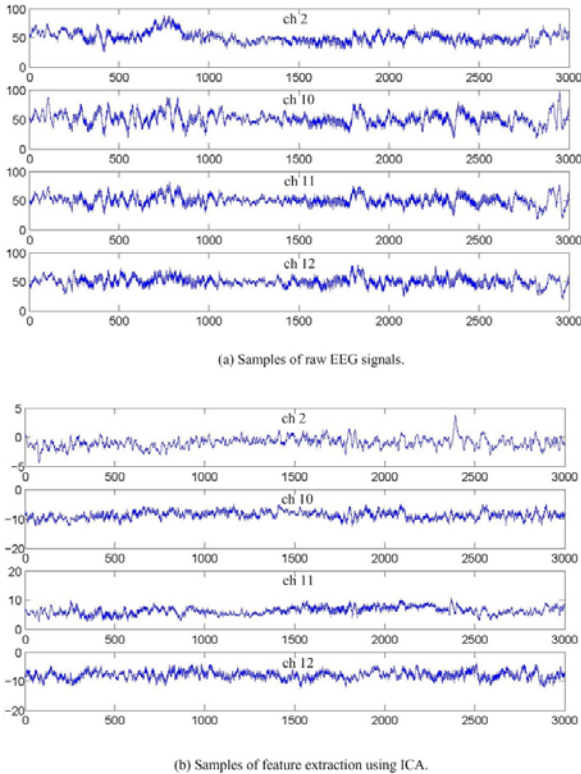


Fig. 3. The samples of raw and feature extraction using ICA. (a) Raw EEG signals from four channels (ch2, ch10, ch11, ch12) of a subject are displayed. The x-axis of each channel denotes number of samples, whereas y-axis represents the amplitude of the EEG signal. (b) The EEG signals after being processed by ICA algorithm (ERICA).

a number of efficient algorithms for ICA [4]. The list of these 22 different ICA algorithms used in the experiment is shown in Table 1. The sample of raw EEG signals in Fig. 3(a) after being processed by ICA algorithm for feature extraction is shown in Fig. 3(b).

Find the best ICA algorithms using a neural network. The objective of this part is to explore which ICA algorithm performs the best among 22 different ICA algorithms. The experiments were conducted by using pattern recognition feature of the neural network on MATLAB 2008a. All EEG data processed from feature extractions were applied to 5, 10, and 20 subjects, respectively. Then, the results were compared as shown in Table 1. The first column shows the name of 22 ICA algorithms. The second column shows accuracy, i.e. the percentage of correctly classified cases from training, validation and testing data sets for five subjects. Note that there are several algorithms achieving 100 percent. To see whether this accuracy is still retained for the other sizes of experimented subjects, the numbers of study subjects were increased to 10 and 20 as shown in the third and fourth columns, respectively. The experimental results of 5, 10 and 20 subjects were studied. ERICA, EWASOBI, JADEop, SIMBEC, SOBI, SOBI-BPF, SOBIRO, and WASOBI algorithm gave high percentage of accuracy. Therefore, one of these algorithms was chosen to employ in the next step to reduce the number of channels. Here, ERICA is the selected algorithm.

Reduce the number of channels from 16 channels. In order to find the minimum number of relevant channels for personal identification, the experiments started with a group of two channels. The total number of selected groups of two channels can be computed by simple binomial coefficient ($^{16}C_2$), which is equal to 120 combinations. In other words, this is the same as counting the total number of possible ways of selecting 2 channels from 16 channels. The EEG data were processed for feature extraction using ERICA algorithm. 500 sample signals of each subject were used for testing. The identification process is based on the extracted features using a 2-layer neural network (one hidden and one output layer) with 20 hidden neurons and some output neurons. The number of output neurons is set according to the number of subjects. For example, there are five output neurons in case of five subjects.

The networks were trained with scaled conjugate gradient backpropagation. The learning algorithm performance was evaluated using mean square error and confusion matrices. All signals of the selected channels were systematically arranged and grouped as the input vectors for training, validation, and testing sets. In the case of two channels, data set consists of two channels of 5, 10 and 20 subjects with each subject using 500 samples. The 500 samples data were divided into 60% (300 samples) for training set, 20% (100 samples) for validation set, and 20% (100 samples) for the testing set. In case of two selected channels, each training pattern is a 4-element vector formed by the signal value from each channel occurring at the same time. The training patterns are generated

Table 1. The accuracy percentage of each ICA algorithm when applied with 5, 10, 20 subjects. The details of all algorithms can be seen in ICALAB [2], [4].

Algorithms	5 Subjects (%)	10 Subjects (%)	20 Subjects (%)
AMUSE	91.44	78.54	57.16
ERICA	100.00	100.00	100.00
EVD2	98.52	84.78	55.66
EWASOBI	100.00	100.00	100.00
FAJDC4	100.00	100.00	95.00
FJADE	100.00	100.00	95.00
FOBI-E	96.20	81.02	44.85
JADEop	100.00	100.00	100.00
JADETD	94.85	75.96	51.16
MULCOMBI	100.00	100.00	95.00
POWERICA	100.00	100.00	95.00
QJADE	100.00	100.00	90.00
SAD	90.88	70.34	45.67
SIMBEC	100.00	100.00	100.00
SOBI	100.00	100.00	100.00
SOBI-BPF	100.00	100.00	100.00
SOBIRO	100.00	100.00	100.00
SONS	94.40	74.38	50.43
SYMMETRIC	100.00	100.00	89.62
THINICA	100.00	100.00	95.00
UNICA	100.00	100.00	95.00
WASOBI	100.00	100.00	100.00

by using the first six signal values from each channel. The next two signal values are used for forming the validation patterns and another next two signal values next to validation patterns are formed as testing patterns. The network was trained with five initial weight sets to test whether the training, validation, and testing sets are properly formed. If they are properly formed, then the network will always be converged. Otherwise, it may be diverged by some initial weight sets.

Table 2 displays the results of groups of two channels when applied with 5, 10 and 20 subjects. When 5 subjects are used in experiments, nine combinations have 100% accuracy percentage. When 10 and 20 subjects are used in the experiment, the best accuracy percentage in the group of Ch2 and Ch7 is dropped to 91.96% and 61.67%, respectively.

The results in Table 2, suggested that only two channels may not be enough to use for identification since the accuracy percentage with 20 subjects is about 60%. Therefore, the number of channel in a group of three is investigated to increase the accuracy percentage. Again, the possible combinations of 3 out of 16 channels are considered (${}^{16}C_3$), which is equal to 560 combinations. From the three-channel experiment, the training set consists of three channels of 5, 10,

Table 2. The accuracy percentage of 2 channels when applied with 5 and 10 subjects using ERICA algorithm

Channel (I)	Channel (II)	Channel 5 subjects (%)	Channel 10 subjects (%)	Channel 20 subjects (%)
ch 2	ch 7	100.00	91.96	61.67
ch 2	ch 9	100.00	91.12	59.04
ch 2	ch 14	100.00	72.62	47.76
ch 2	ch 15	100.00	74.28	48.44
ch 5	ch 10	100.00	71.90	39.46
ch 5	ch 15	100.00	83.76	20.34
ch 9	ch 10	100.00	88.42	44.07
ch 10	ch 16	100.00	88.10	54.05
ch 12	ch 15	100.00	64.96	54.62

and 20 subjects with each subject using 500 samples. The ratio of the training set, validation set and testing set is 60:20:20 as in the case of two channels. The results of three channels are shown in Table 3.

Table 3. The accuracy percentage of 3 channels when applied with 5,10,20 subjects using ERICA algorithm

Channel (I)	Channel (II)	Channel (III)	Channel 5 subjects (%)	Channel 10 subjects (%)	Channel 20 subjects (%)
ch 4	ch 5	ch 10	100.00	100.00	90.52
ch 4	ch 10	ch 13	100.00	99.96	68.08
ch 4	ch 13	ch 15	100.00	100.00	75.16
ch 5	ch 9	ch 10	100.00	99.98	55.34
ch 9	ch 13	ch 15	100.00	99.98	83.41

Table 3 shows only partial of the three channels combinations. The accuracy percentages with 5 subjects reach 100%. With 10 subjects, the accuracy percentages reach 100% in the groups of Ch4, Ch5, Ch10 and Ch4, Ch13, Ch15. When numbers of subjects were increased to 20 subjects, the result shows that the highest percentage is 90.52% in the group of Ch4, Ch5, Ch10.

The accuracy percentage of the experiment using three channels for 20 subjects is still not high enough (the ideal accuracy should get close to 100%). So the number of channels was increased to a group of four channels to increase the accuracy percentage of the 20 subjects.

In case of four channels, total combinations of 4 out of 16 channels are considered (${}^{16}C_4$), which is equal to 1,820 combinations. The training set consists of four channels of 5, 10, and 20 subjects, with each subject employing 500 samples. The ratio of the training set, validation set and testing set of four channels is the same as in the case of two and three channels. The results four channels groups are shown in Table 4.

Table 4. The accuracy percentage of 4 channels when applied with 5,10,20 subjects using ERICA algorithm

Channel (I)	Channel (II)	Channel (III)	Channel (IV)	5 subjects (%)	10 subjects (%)	20 subjects (%)
ch 2	ch 10	ch 11	ch 12	100.00	100.00	98.71
ch 2	ch 4	ch 5	ch 5	100.00	100.00	98.03
ch 4	ch 5	ch 10	ch 15	100.00	100.00	97.34
ch 2	ch 4	ch 7	ch 11	100.00	100.00	96.96
ch 2	ch 4	ch 13	ch 15	100.00	100.00	96.43

Table 4 shows only a part of the four channels combinations. The accuracy percentages with 5 and 10 subjects reach 100%. When the numbers of subjects were increased to 20 subjects, the result shows that the highest percentage is 98.71% in the group Ch2, Ch10, Ch11, Ch12. The results of the experiments for 5, 10, and 20 subjects are acceptable. Therefore, four channels are the lowest possible channels for this experiment.

Table 5. Reducing the number of channels from 16 to 4 channels

Algorithms	The best accuracy for 5 subjects (combination)	The best accuracy for 10 subjects (combination)	The best accuracy for 20 subjects (combination)
ERICA	1,502 combinations reach 100%	185 combinations reach 100%	ch2 ch10 ch11 ch12 98.71%
EWASOBI	862 combinations reach 100%	34 combinations reach 100%	ch4 ch 5 ch 8 ch13 91.17%
JADEop	40 combinations reach 100%	4 combinations reach 100%	ch9 ch11 ch14 ch15 93.48%
SIMBEC	110 combinations reach 100%	14 combinations reach 100%	ch8 ch12 ch13 ch16 93.48%
SOBI	110 combinations reach 100%	4 combinations reach 100%	ch6 ch 7 ch12 ch13 96.22%
SOBI-BPF	1,511 combinations reach 100%	218 combinations reach 100%	ch3 ch 8 ch 9 ch12 98.57%
SOBIRO	170 combinations reach 100%	22 combinations reach 100%	ch5 ch10 ch13 ch15 97.87%
WASOBI	1,814 combinations reach 100%	65 combinations reach 100%	ch3 ch 6 ch 9 ch10 94.26%

The result of Table 1 illustrated that 8 algorithms; namely ERICA, EWA-SOBI, JADEop, SIMBEC, SOBI, SOBI-BPF, SOBIRO, and WASOBI algorithm, give the 100% accuracy. Thus, all these eight algorithms will be applied with the groups of four channels to further investigate the best algorithm for

identification. In the experiments using 8 different ICA algorithms, the total combination of 4 out of 16 channels is equal to 1,820 combinations. The first algorithm (ERICA) starts with all 1,820 combinations and tests with 5 subjects. There are 1,502 combinations reaching 100% accuracy. Then, these 1,502 combinations are experimented with 10 subjects. The result shows that there are 185 combinations that give 100% accuracy. These 185 combinations are then applied with 20 subjects. The final result is the combinations of Ch2, Ch10, Ch11, and Ch12, which give the highest accuracy percentage (98.71%). The same process was applied for other algorithms. The results are shown in Table 5.

From the result of Table 5, the best algorithm for 4-channel groups is ERICA. The accuracy percentage is 98.71% in the group of channels ch2, ch10, ch11, ch12.

Test the neural network with 20 subjects. The purpose of this step is to test the accuracy of recognition. The neural network that has been trained earlier using ERICA algorithm on ch2, ch10, ch11, ch12 will be used to test the accuracy percentage. From Table 5, ERICA algorithm performs the best. In order to further investigate the performance of this recognition, 3,000 samples of the EEG signals are classified into 6 groups G1 - G6 (shown in Fig. 4). Each group contains 500 samples. The first group: G1 consists of sample data from 1 - 500, the second group: G2 consists of sample data from 501-1000, and so on until the sixth group: G6 is the sample data from 2501-3000. Note that G1 has been used for training and testing the recognition in the previous step. However, G2 to G6 are samples that have not been used to train and test the recognition performance before. Thus, these samples: G2 to G6 are then used to test the recognition performance on 20 subjects. The results are shown in Table 6.

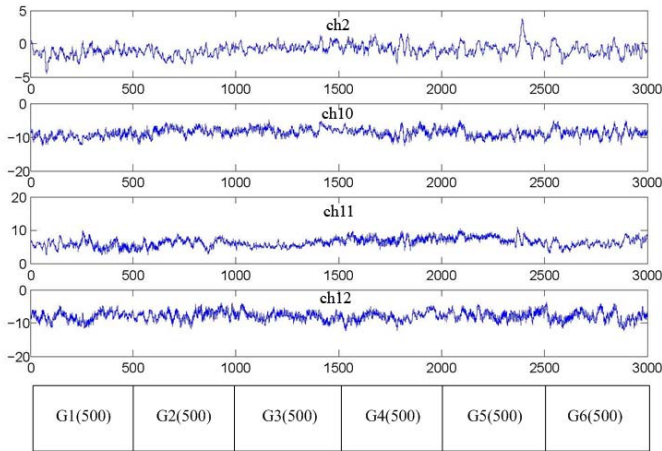


Fig. 4. The EEG signals are divided into six groups for testing

Table 6. The accuracy percentage of 20 subjects when test with the best combination of ERICA algorithm 4 channels 500 samples

Subject	G1 (%)	G2 (%)	G3 (%)	G4 (%)	G5 (%)	G6 (%)
1	100.00	100.00	99.20	98.40	93.40	98.80
2	100.00	100.00	100.00	98.20	100.00	100.00
3	100.00	100.00	100.00	100.00	100.00	100.00
4	100.00	100.00	100.00	100.00	100.00	100.00
5	100.00	99.60	100.00	99.80	100.00	100.00
6	97.00	99.60	99.80	99.60	97.60	99.40
7	99.60	97.80	95.00	99.00	94.60	94.00
8	100.00	99.40	92.40	100.00	100.00	99.60
9	100.00	100.00	100.00	99.80	100.00	100.00
10	97.40	96.40	93.40	97.80	96.20	90.00
11	100.00	99.00	95.60	99.00	98.80	99.60
12	100.00	100.00	100.00	99.20	99.80	95.00
13	99.80	98.60	99.20	100.00	100.00	100.00
14	95.40	81.00	81.80	91.60	82.40	81.20
15	99.40	92.80	95.40	97.00	92.60	94.20
16	96.80	96.20	97.60	94.80	89.00	93.00
17	94.60	92.60	94.60	94.40	98.00	96.40
18	95.80	89.40	93.20	94.80	86.60	90.60
19	98.40	98.40	99.40	98.40	99.60	90.60
20	100.00	99.80	100.00	100.00	99.80	100.00
Average						
G(I)	98.71	97.03	96.83	98.09	96.42	96.59

Numbers in Table 6 are the accuracy percentage for each subject in a particular group. In the first row, for example, the accuracy percentage of subject 1 was evaluated with all groups of data (G1 to G6) . The accuracy percentage of subject 1 using G1 data is 100%. The accuracy percentages using G2 - G6 data are 100%, 99.20%, 98.40%, 93.40%, 98.80%, respectively. The last row of the Table 6 shows the average accuracy percentage in each test data group. It can be seen that the average accuracy of these six groups is in the range of 96% - 99%. The results show that identifications on 20 subjects are very encouraging.

3 Conclusion and Future Work

In this paper, a practical technique for identifying from standard 16 EEG locations of EEG signals is proposed. The signals are cleaned first by applying Independent Component Analysis. Then, a supervised neural network is used to test the accuracy of the identification process. We found that minimum

location of EEG signals locations F_7 , C_3 , P_3 , and O_1 are relevant channel efficiently used for identification by ERICA algorithm. The results show that EEG signals can be used for identification for 20 subjects with high accuracy. However, the number of subjects is considered to be small and all data are collected at the same time period. It is more interesting to prove that this study still works even to higher number of subjects and can separate the subject who does not in the group. In addition, the EEG data which are still obtained from different time period on the same subjects would still give the high accuracy percentage for the identifications. This hypothesis will be further studied.

Acknowledgments. This work was supported by The Centre of Excellence in Mathematics, PERDO, Commission on Higher Education, Thailand and THE 90th ANNIVERSARY OF CHULALONGKORN UNIVERSITY FUND (Ratchadaphiseksomphot Endowment Fund).

References

1. Jain, A.K., Ross, A., Prabhakar, S.: An Introduction to Biometric Recognition. *IEEE Trans. Circuits Syst. Video Technol.* 14(1), 4–20 (2004)
2. Cichocki, A., Amari, S.: Adaptive blind signal and image processing: learning algorithms and applications. Wiley, New York (2003)
3. Niedermeyer, E., da Silva, F.L.: Electroencephalography: Basic Principles. In: *Clinical Applications, and Related Fields*, pp. 139–141. Lippincott Williams, Wilkins (2004)
4. Cichocki, A., Amari, S., Siwek, K., Tanaka, T., et al.: ICALAB toolboxes, <http://www.bsp.brain.riken.jp/ICALAB>
5. Poulos, M., Rangoussi, M., Alexandries, N.: Neural network based person identification using EEG features. In: *Proceeding IEEE of the international conference on Acoustic, Speech and Signal Processing, ICASSP 1999, Arizona, USA*, pp. 1117–1120 (1999)
6. Poulos, M., Rangoussi, M., Kafetzopoulos, E.: Person identification via the EEG using computational geometry algorithms. In: *Proc. Intl. Conf. EUSIPCO 1998, Rhodes, Greece (September 1998)*
7. Marcel, S., Millan, J.: Person authentication using brain-waves (EEG) and maximum a posteriori model adaptation, Tech. Rep. 81, p. 11, IDIAP Research Report, Valais, Switzerland (2005)
8. Paranjape, R.B., Mahovsky, J., Benedicenti, L., Koles, Z.: The electroencephalogram as a biometric. In: *Proceedings of Canadian Conference on Electrical and Computer Engineering, Toronto, Canada, May 2001, vol. 2*, pp. 1363–1366 (2001)
9. Mohammadi, G., Shoushtari, P., Ardekani, B., Shamsollahi, M.: Person identification by using AR model for EEG signals. In: *Proceedings of the 9th International Conference on Bioengineering Technology (ICBT 2006), Czech Republic, vol. 5* (2006)
10. Daly, D.D., Markand, O.N.: Focal brain lesions. In: Daly, D.D., Pedley, T.A. (eds.) *Current Practice of Clinical Electroencephalography*, 2nd edn. Raven Press, New York (1990)
11. Poulos, M., Rangoussi, M., Alexandris, N., Evangelou, A.: On the use of EEG features towards person identification via neural networks. *Medical Informatics and the Internet in Medicine* 26(1), 35–48 (2001)

12. Poulos, M., Rangoussi, M., Alexandris, N., Evangelou, A.: Person identification from the EEG using nonlinear signal classification. *Methods of Information in Medicine* 41(1), 64–75 (2002)
13. Palaniappan, R., Ravi, K.V.R.: A New Method to Identify Individuals Using Signals from the Brain. In: *Proc. Fourth Int'l Conf. Information Comm. and Signal Processing*, pp. 15–18 (2003)
14. van Drongelen, W.: *Signal Processing for Neuroscientists An introduction to the analysis of physiological signals*. Academic Press, London (2007)
15. Saeid Sanei, J.A.: *Chambers, EEG Signal Processing*. Wiley, Chichester (2007)

A Formal Concept Analysis-Based Domain-Specific Thesaurus and Its Application in Document Representation

Jihn-Chang Jehng¹, Shihchieh Chou², and Chin-Yi Cheng²

¹Institute of Human Resource Management, National Central University,
No.300, Jhongda Rd., Jhongli City, Taoyuan County 32001, Taiwan
jehng@mgmt.ncu.edu.tw

²Department of Information Management, National Central University,
No.300, Jhongda Rd., Jhongli City, Taoyuan County 32001, Taiwan
scchou@mgmt.ncu.edu.tw,
yoda0612@seed.net.tw

Abstract. Many techniques in the process of document retrieval and clustering, based on the vector space model, represent documents by vectors. They ignore the conceptual relationships of terms such as synonyms, hypernyms and hyponyms and, especially, treat terms as a *bag of terms*. The application of conceptual relationships of terms has been proved by generating improved results for document clustering in previous studies. For those studies, thesauri like Word-Net were used to provide the information of relationships between terms. However, some domain-specific terms like "query expansion" and "document clustering" cannot be found in these thesauri. These terms are thought of as important features in domain-specific documents. In this paper, we propose an automatic domain-specific thesaurus building approach based on Formal Concept Analysis (FCA) dealing with the problem with general thesauri. We also apply the domain-specific thesaurus as background knowledge to represent documents by concept dimension vectors. In the evaluation, an improved result by our method compared to traditional approaches is shown.

Keywords: Information Retrieval, Formal Concept Analysis, Concept Lattice, Vector Space Model.

1 Introduction

With the growth of internet technology, ever more information is becoming available online. However, this causes the problem of information overload, in that users have to make greater efforts to find the information they need from web sites or document bases. Many techniques, such as information retrieval, personalized portal or document management are developed to help users solve this problem, by organizing the search results or document bases in hierarchical structures, or providing users additional vocabularies to reformulate their query text.

Most of the techniques in the processing of document retrieval and clustering represent documents by vectors. Usually, the dimensions of the vectors are the terms extracted from documents. This approach has been widely applied to information retrieval and related techniques, such as document retrieval, document clustering, and document classification.

The main problem with this *bag of terms* approach is that it ignores the conceptual similarity of terms [1][2][3]. Suppose there are two documents s_1 and s_2 ; term "query expansion" appears in the topic of s_1 and "document clustering" appears in s_2 . In the traditional approach, the two documents will be identified irrelevant, because there is no hypernym like "information retrieval" to make them relevant.

The application of conceptual relationships of terms to document retrieval and clustering in some cases can improve the precision and recall rate [1][2]. Usually, the authors of documents use different words to describe the same thing. However, traditional approaches do not take these relevant words into synonym sets. A thesaurus like WordNet [4] can provide the information of conceptual relationships of terms such as synonyms, hypernyms or hyponyms. Previous studies, such as Hotho et al. [2] and Recuperio [1], used WordNet as background knowledge to improve the results of document clustering.

A thesaurus like WordNet is a general thesaurus, and does not aim at the use of specific domains. Therefore, some domain-specific words or phrases cannot be found in WordNet. These words or phrases are always thought of as important features in domain-specific documents. For example, phrases such as "information retrieval", "query expansion" and "document clustering" are important features in information retrieval-related documents; however, they cannot be found in WordNet. To deal with the problem of general thesauri, some researches construct domain-specific thesaurus by ontology [3][5]. However, the building of domain-specific ontology always requires the participation of humans and is manual or semi-automatic.

To tackle the problems mentioned above, an automatic domain-specific thesaurus building methodology is needed. Formal concept analysis (FCA) is mainly used for data analysis and in particular for extracting inherent relationships between objects in a given context [6]. This feature helps us to extract the relationships of terms and form the terms into a thesaurus. Previous researches applied FCA to identify the related terms in a given document set and used these terms to generate expanding terms for query expansion [7].

In this paper, we propose an automatic domain-specific thesaurus building approach by FCA. Each concept of the domain-specific thesaurus consists of related terms, and the relations between concepts are inherent. This thesaurus provides the relations of terms for a given specific domain. We also use it as background knowledge to represent documents by concept dimension vectors.

The rest of this paper is organized as follows: In section 2, we introduce the FCA and its applications in information retrieval. Major definitions of FCA are also described. In section 3, the proposed method of constructing the domain-specific thesaurus by FCA is explained. In this section, we first introduce the procedures of document preprocessing, then build our domain-specific thesaurus by FCA. Finally, we explain the method of representing documents by concept dimension vectors. In section 4, evaluation of our method is conducted. We make a comparison with traditional approaches to show the improved results of the use of the domain-specific thesaurus. Finally, we come to the conclusions of this paper.

2 Formal Concept Analysis

FCA, first introduced by Rudolf Wille in 1982 [8], is a theory of data analysis which identifies conceptual structures among data sets [9][10]. It is based on the *Lattice Theory* [11], which has been widely applied within many different realms. FCA forms a given application domain in terms of objects (extensions) and the attributes (intensions) that are used to describe those objects, and provides a meaningful, comprehensible interpretation of the given context according to *Concept Lattices* [10].

FCA has been successfully used in a wide range of applications in information retrieval. Groojen and Weide use FCA to discover the number of related term combinations for query expansion [7]. Rajapakse and Denham make up a learning strategy with FCA to improve the results of document retrieval [12]. Cimiano et al. apply FCA to generate concept hierarchies from text corpora for certain domain applications [6].

One of the most important features of FCA is that it can discover the inherent relationships between objects from a given context. In the application of information retrieval, some hierarchical document clustering techniques have been developed; these techniques are based on the inherent relationships between documents. For example, Carpineto and Romano [13] apply a *Concept Lattice* generated from FCA to support conceptual document clustering and provide a browsing retrieval by the system they implemented named GALOIS. CREDO, another case, is a system that organizes web search results in a hierarchical structure by *Concept Lattice*, is also developed by Carpineto and Romano [14].

In this paper, we apply FCA to generate a domain-specific thesaurus to extract the relations of terms. The following are the major notions of FCA.

Definition 1

A *formal context* is defined as a triple (G, M, I) , where G and M are two sets and $I \subseteq G \times M$ is a binary relation between G and M . The elements of G are called objects and those of M are called attributes. I is the relation defined between G and M . To represent a relation I between object g and attribute m , we write gIm or $(g, m) \in I$.

In Table 1, we provide an example of a formal context which contains four objects from *Document1* to *Document4*, and attributes are the words appearing in the given documents. The elements of Table 1 represent the binary relation between G and M that indicates whether terms appeared in the documents.

Definition 2

For $A \subseteq G$, we define

$$A' = \{m \in M \mid (g, m) \in I \text{ for all } g \in A\}$$

denotes the set attributes common to all the objects in A .

Table 1. An example of formal context

Document/Term	<i>Information retrieval</i>	<i>Query expansion</i>	<i>Document clustering</i>	<i>Local analysis</i>	<i>Global analysis</i>
<i>Document 1</i>	X	X		X	
<i>Document 2</i>	X		X		
<i>Document 3</i>	X	X			X
<i>Document 4</i>	X		X		

generating *Concept Lattices* can be found in [15]. In this paper, we use the tool of *Concept Explorer*¹ for producing *Concept Lattices*.

3 Domain-Specific Thesaurus for Document Representation

Traditionally, document representation is based on the use of *bag of terms*. Given a document set D , $T=\{t_1, t_2, \dots, t_n\}$ is the set of different terms appearing in D . Then, a document $d_i \in D$ can be represented as a term dimension vector

$$\vec{d}_i = (tfidf(d_i, t_1), tfidf(d_i, t_2), \dots, tfidf(d_i, t_n)),$$

where $tfidf(d_i, t_n)$ is a weighting schema for measuring the importance of term t_n in document d_i . The $tfidf$ weighting schema can be defined as follows.

Definition 5: $tfidf$ (term frequency-inverted document frequency) term weighting schema

$$tfidf(d_i, t_n) = tf(d_i, t_n) \times \log\left(\frac{|D|}{df(t_n)}\right)$$

where $tf(d_i, t_n)$ is the frequency of term t_n in document d_i ; $df(t_n)$ is the document frequency that indicates the number of documents containing term t_n . This weighting schema is used to measure the importance of terms in a given document. Some further applications, such as document retrieval, document clustering, or document classification, apply *Cosine* or *Jaccard* measure to compute the distance between two vectors which are represented documents or queries.

This very typical method of document representation, called the vector space model, was first introduced by Salton [16] and has been successfully and widely applied to information retrieval. However, one of the main drawbacks of this approach is its ignoring of conceptual relationships of terms. To tackle this problem, we proposed an automatic domain-specific thesaurus building approach based on FCA. This domain-specific thesaurus is mainly used to represent the relations between terms via the relations between formal concepts. We also applied this thesaurus as background knowledge to represent documents by concept dimension vectors.

In our method, FCA is applied to construct the domain-specific thesaurus. The *Concept Lattice* generated by FCA for a given document set consists of a set of formal concepts in which related terms are formed. By the definitions introduced in Section 2, all relations between formal concepts in a *Concept Lattice* are inherent, and these relations are called super-concept and sub-concept relations. Therefore, we assume that the terms in the same formal concept can be regarded as synonyms; and the terms in super-concept are hypernyms for those terms in their sub-concepts; on the other hand, terms in sub-concept can be thought as hyponyms for those terms in their super-concepts. In the following sections, we introduce the procedure for our method in detail.

¹ *Concept Explorer* can be download at <http://sourceforge.net/projects/conexp/>

3.1 Document Preprocessing

First, all documents are broken down into sentences and then these sentences are part-of-speech (POS) tagged, using Stanford POS Tagger². The tags help us to identify the POS of each word in a sentence, such as nouns, verbs, and adjectives. The POS tags are also important in extracting the phrases from documents. These extracted phrases, we thought, were critical features of domain-specific documents. The following is an example of POS tagged sentence.

Original example sentence:

In this paper, we propose an automatic domain-specific thesaurus building approach based on Formal Concept Analysis (FCA) dealing with the problem with general thesauri.

POS tagged example sentence:

In/IN this/DT paper,/NN we/PRP propose/VBP an/DT automatic/JJ domain-specific/JJ thesaurus/NN building/NN approach/NN based/VBN on/IN Formal/NNP Concept/NNP Analysis/NNP (FCA)/NNP dealing/VBG with/IN the/DT problem/NN with/IN general/JJ thesauri./NN

Each token after the slash is the POS tag of words. For example, *IN* is preposition; *DT* is determiner; and *NN* is noun. A complete list of POS tags can be found in [17].

Second, we keep all nouns, and in order to extract the domain-specific phrases from documents, we refer to a list of POS patterns constructed by Ou, et al. [18] for recognizing phrases (See Table 2). Then, stop-words are removed from remained nouns.

Table 2. Part-of-speech patterns for recognizing phrases

Part-of-speech tag			Example term
1	2	3	
NN			QE, IR, FCA
JJ	NN		unexpanded query
NN	NN		query expansion, document clustering
JJ	NN	NN	interactive query expansion

Third, we stemmed all remained nouns and phrases by Porter stemming algorithm [19].

3.2 Concept Lattice Generation

After the document preprocessing, we get a set of terms T which consists of single nouns or phrases. Then, a formal context is constructed from a training document set D , where the objects of a formal context $fc_k \in FC$ are the documents in the training set, and the attributes are the terms belong to T ; the relation between objects and

² <http://nlp.stanford.edu/software/tagger.shtml>

attributes is the appearance of $t \in T$ in $d \in D$. Then, a *Concept Lattice* can be generated by FCA, which contains a number of formal concepts within several terms Tfc_k and documents Dfc_k , and super-concept and sub-concept relations between formal concepts.

3.3 Domain-Specific Thesaurus Generation

We defined our domain-specific thesaurus as a network-like structure, which is used to represent relations of terms for a given document set. The following notions are given for our domain-specific thesaurus.

Definition 6

A domain-specific thesaurus is defined as a tuple $DST=(C, <)$ consisting of a set of concept C and inherent relation $<$ with C . These inherent relations are also called sub-concept and super-concept relations. For each concept $c_i \in C$, it contains a set of term Tc_i .

Definition 7

For $c_1, c_2 \in C$, if $c_1 < c_2$, we called c_1 is a sub-concept of c_2 ; on the other hand, c_2 is super-concept of c_1 . If c_1 is sub-concept of c_2 , then Tc_2 is a subset of Tc_1 ($Tc_2 \subseteq Tc_1$).

Definition 8

A concept c_i is the root concept in DST , if and only if there is no super-concept for c_i . We write a root concept as c_{Root} . A root concept c_i is the only concept that its Tc_i may be empty in DST .

In order to transform the *Concept Lattice CL* generated in Section 3.2 into a domain-specific thesaurus DST defined above, the following steps are executed. First, all information about objects is removed from a formal concept, while the object information is no longer needed in the following process. Second, we delete the most bottom formal concept whose object set is empty and term set contains all terms of T . Third, we eliminate the formal concept whose set of terms Tc_i equals the union of sets of terms of its all direct super-concepts, and the edges connected to this concept.

This pruned *Concept Lattice* with no circle is used to represent the domain-specific thesaurus DST defined above. We keep all relations between formal concepts and use the pruned formal concepts to represent concept c_i defined in DST .

3.4 Transform Term Vector to Concept Vector

In order to deal with the problem with traditional *bag of terms* approaches which do not take account of the conceptual similarity of terms in their applications, we apply the domain-specific thesaurus defined in Section 3.3 as background knowledge to represent a given document by concept dimension vector. A new concept dimension vector can be represented as

$$\vec{d}_i = (cw(d_i, c_1), cw(d_i, c_2), \dots, cw(d_i, c_n)),$$

where $cw(d_i, c_n)$ is the weight of concept $c_i \in C$ in the document d_i . For each document, all of its concept weights are determined by the following steps.

First, we compute all concept weights for all concepts $c_i \in C$. We define a concept weighting schema for each concept as follows.

Definition 9

For all $c_i \in C$,

$$cw(d_i, c_n) = \sum_{t \in Tc_i'} f(d_i, t)$$

where Tc_i' is term set in c_i excluded the terms appeared in the direct super-concepts of concept c_i .

Second we update concept weights for all concepts, except the concepts at most bottom level, from its sub-concepts. In our definitions of *DTS*, a concept c_i may inherit from multiple super-concepts. This is significantly different to hierarchical structure. In hierarchical structure, one concept only has one super-concept. Therefore, we developed a bottom-up approach which updates the weight of concepts by moving the weight of a given concept up to all its super-concepts.

A recursive algorithm called *UpdateConceptWeights* (Fig. 2) is used to update the concept weights for all super-concepts for a given concept. A given concept c and its weight *SubWeight* will be the inputs of *UpdateConceptWeights*. If concepts c is not root concept, this algorithm will get the summation value *NumberOfSuperTerms* of the number of terms of c 's direct super-concepts (lines 5-7). For each direct super-concept c_{super} of c , their weight is determined by the maximum function (line 9). Two values will be the inputs of this maximum function, one is the original concept weight of c_{super} and another is distributed value of *SubWeight*. The distributed value of *SubWeight* for c_{super} is defined as

$$\frac{NumberOfTerms(c_{super})}{NumberOfSuperConceptTerms} \times \omega,$$

Algorithm: UploadConceptWeights

Input: Concept c , SubConceptWeights *SubWeight*

Output: Concept weights for all concepts in *DST*

1. **if** $c = c_{Root}$
2. **return**
3. **else**
4. $NumberOfSuperConceptTerms = 0$
5. **foreach** $c_{super} \in SuperConcept(c)$
6. $NumberOfSuperConceptTerms += NumberOfTerms(c_{super})$
7. **end foreach**
8. **foreach** $c_{super} \in SuperConcept(c)$
9. $cw(d_i, c_{super}) = \max(cw(d_i, c_{super}),$
 $\frac{NumberOfTerms(c_{super})}{NumberOfSuperConceptTerms} \times \omega \times SubWeights)$
10. $UploadConceptWeights(c_{super}, cw(d_i, c_{super}))$
11. **end foreach**
12. **end if**

Fig. 2. Pseudo code of *UpdateConceptWeights*

where $NumberOfTerms(c_{super})$ represents the number of terms in c_{super} and ω is a constant for tuning. The Updated concept weight of c_{super} will be the maximum value of the two input values of maximum function at line 9. Finally, all direct super-concepts of c and their updated concept weight $cw(d_i, c_{super})$ will be the inputs of algorithm *UpdateConceptWeights* (line 10). This algorithm will end if input concept is root concept (lines 1-2).

All concepts, except root concepts, with their weight in *DTS* will be the inputs of the algorithm *UpdateConceptWeights*. The execution sequence of this algorithm is from the concepts at the most bottom level to concepts at topmost level, level by level. For each concept it will move its weight up to all its direct and indirect super-concepts by this algorithm.

4 Evaluation

In this section, we demonstrate our method introduced in Section 3 through a simple example given as follows: We first generate a domain-specific thesaurus from an example training document set. Furthermore, this thesaurus is used as background knowledge for representing the given document set by concept dimension vectors. Finally, we analyze and compare the results of the concept dimension vectors with traditional term dimension vectors in the measuring of document similarities.

The test data set we use in this evaluation are the sentences selected from [20], which is a research paper mainly focused on the topic of query expansion. We select sentences, which describe the techniques about local and global query expansion, as training and test document sets.

A training document set which consists of four documents from *TrainDoc1* to *TrainDoc4* is given as follows:

- TrainDoc1:* Existing techniques for automatic query expansion can be categorized as either global or local.
- TrainDoc2:* A global technique requires some corpuswide statistics that take a considerable amount of computer resources to compute.
- TrainDoc3:* Local technique processes a small number of top-ranked documents retrieved for a query to expand that query.
- TrainDoc4:* A local technique may use some global statistics such as the document frequency of a term.

Highlighted terms, including single nouns and phrases, are selected by the procedure of document preprocessing introduced in Section 3.1. All selected terms are stemmed by Porter's algorithm. We use these terms as the attributes and training documents as objects to form the formal context. Then, according to the method described in Section 3.2 and Section 3.3, a domain-specific thesaurus for this context can be generated as Fig. 3.

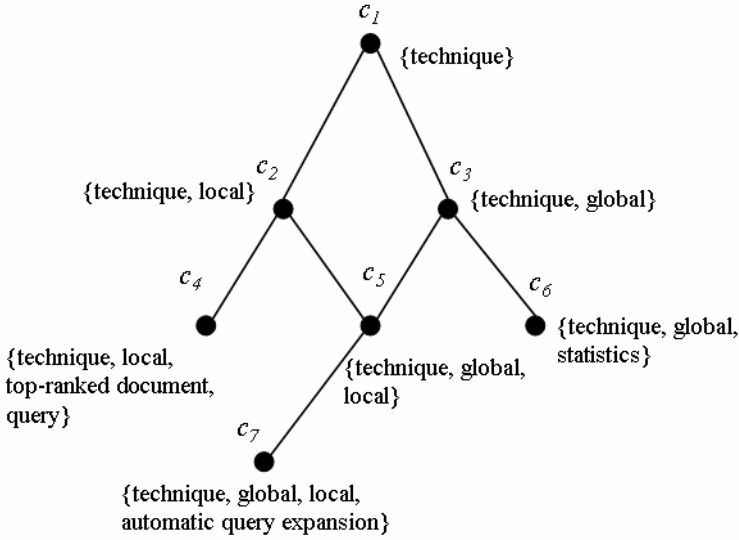


Fig. 3. Domain-specific thesaurus for test document set

Three test documents are given as follows: The important terms in documents are selected and highlighted.

- TestDoc1:* One of the earliest global techniques is term clustering.
- TestDoc2:* Information from the top-ranked documents is used to reestimate the probabilities of query terms in the relevant set for a query.
- TestDoc3:* Local techniques expand a query based on the information in the set of top-ranked documents retrieved for the query.

According to the procedures described in Section 3.4, each concept in domain-specific thesaurus, as shown in Fig. 3, is the dimension of vectors which will be used to represent test documents. Then, we compute concept weights for all concepts in domain-specific thesaurus for each test document by concept weighting schema defined in definition 9, and execute the algorithm of *UpdateConceptWeights* to update the concept weights for concepts which have sub-concepts. In this example, the value of ω is set to 0.5. The values of concept weights for each test document are used as the dimension values for vectors. The results of document representation are shown in Table 3.

Table 3. Vectors of test documents

	Dimensions						
	c_1	c_2	c_3	c_4	c_5	c_6	c_7
$\vec{d}_{TestDoc1}$	1	0	1	0	0	0	0
$\vec{d}_{TestDoc2}$	0.5	1	0	2	0	0	0
$\vec{d}_{TestDoc3}$	1	1.5	0	3	0	0	0

In evaluating the method, we are especially interested that the sub-concept and super-concept relationships between terms extracted by our method may create a chance to make two documents relevant. These two documents may be identified irrelevant by traditional approaches. Table 4 shows the document representations by traditional *bag of terms* approach with *tf-idf* weights value for the test documents. After *Cosine* measuring, the similarity between test documents *TestDoc1* and *TestDoc2* is zero, because, there are no common terms between these two documents. However, according to the results of our method shown in Table 3, *TestDoc1* and *TestDoc2* are accounted with small similarity by *Cosine* measuring, because they have a shared concept in c_1 . The terms {"top-ranked documents", "query"} in c_4 are sub-concept of term {"local"} in c_2 , it seems reasonable, because selecting useful expansion terms from top-ranked documents is the basic idea of local based query expansion. The term {"local"} in c_2 and term {"global"} in c_3 inherit root concept c_1 which contains term {"technique"}; in this context, it may mean the technique of query expansion. Therefore, although there is no term in *TestDoc2* matching the terms in concept c_1 , *TestDoc2*, however, contains terms in c_4 which indirectly inherits c_1 , and therefore c_1 becomes a shared concept between *TestDoc1* and *TestDoc2* and makes these two documents relevant.

Table 4. *tf-idf* weights of test documents

	global	technique	local	information	Term dimensions					term cluster
					top-ranked documents	probability	query	query terms		
<i>TestDoc1</i>	0.48	0.18	0	0	0	0	0	0	0.48	
<i>TestDoc2</i>	0	0	0	0.18	0.18	0.48	0.18	0.48	0	
<i>TestDoc3</i>	0	0.18	0.48	0.18	0.18	0	0.35	0	0	

5 Conclusion

In this paper, we provide an automatic domain-specific thesaurus building approach, and use it as background knowledge to represent documents by concept dimension vectors. This thesaurus built by FCA contains several concepts with the relations of sub-concept and super-concept, and each concept consists of a set of terms. According to the concept the terms belong to, the relationships between terms, including synonyms, hypernyms and hyponyms, can be identified. We apply this domain-specific thesaurus and develop the methodology of concept weighting for document representation. The demonstration shows the improved results on the measuring of similarities between documents. Compared to traditional approaches, our method successfully makes documents with shared concepts relevant. Our work is ongoing and will be extended to the information retrieval-related techniques, including document clustering, document classification and query expansion. Furthermore, given that *Concept Lattice* building is a computationally difficult task [13][12], we also plan to develop a learning strategy for constructing domain-specific thesauri incrementally.

References

1. Recuperó, D.R.: A New Unsupervised Method for Document Clustering by Using WordNet Lexical and Conceptual Relations. *Information Retrieval* 10(6), 563–579 (2007)
2. Hotho, A., Staab, S., Stumme, G.: Wordnet Improves Text Document Clustering. In: Proc. of the SIGIR 2003 Semantic Web Workshop (2003)
3. Hotho, A., Staab, S., Stumme, G.: Ontologies Improve Text Document Clustering. In: Proc. of the IEEE International Conference on Data Mining (ICDM), pp. 541–544. IEEE Press, New York (2003)
4. Miller, G.: WordNet: A Lexical Database for English. *CACM* 38(11), 39–41 (1995)
5. Hotho, A., Staab, S., Maedche, A.: Ontology-based Text Clustering. In: IJCAI 2001, Workshop on Ontology Learning. Seattle, Washington, USA (2001)
6. Cimiano, P., Hotho, A., Staab, S.: Learning Concept Hierarchies from Text Corpora Using Formal Concept Analysis. *J. of Artificial Intelligence Research* 24, 305–339 (2005)
7. Grooijten, F.A., Weide, T.P.: Conceptual Query Expansion. *Data & Knowledge Engineering* 56(2), 174–193 (2006)
8. Wille, R.: Restructuring Lattice Theory: an Approach Based on Hierarchies of Concepts. In: Rival, I. (ed.) *Ordered sets*, pp. 445–470. D. Reidel Publishing Company, Dordrecht (1982)
9. A Formal Concept Analysis Homepage,
<http://www.upriss.org.uk/fca/fca.html>
10. Ganter, B., Wille, R.: *Formal Concept Analysis - Mathematical Foundations*. Springer, Heidelberg (1999)
11. Birkhoff, G.: *Lattice Theory*. American Mathematical Society, Providence (1967)
12. Rajapakse, R.K., Denham, M.: Text Retrieval with more Realistic Concept Matching and Reinforcement Learning. *Information Processing and Management* 42(5), 1260–1275 (2006)
13. Carpineto, C., Romano, G.: A Lattice Conceptual Clustering System and its Application to Browsing Retrieval. *Machine Learning* 24(2), 95–122 (1996)
14. Carpineto, C., Romano, G.: Exploiting the Potential of Concept Lattices for Information Retrieval with CREDO. *J. of Universal Computer Science* 10(8), 985–1013 (2004)
15. Carpineto, C., Romano, G.: *Concept Data Analysis: Theory and Applications*. John Wiley & Sons, Ltd., Hoboken (2004)
16. Salton, G.: A Vector Space Model for Automatic Indexing. *CACM* 18, 613–620 (1975)
17. Marcus, M.P., Santorini, B., Marcinkiewicz, M.A.: Building a Large Annotated Corpus of English: The Penn Treebank. *Computational Linguistics* 19(2), 313–330 (1993)
18. Ou, S., Khoo, C.S., Goh, D.H.: Design and Development of a Concept-Based Multi-Document Summarization System for Research Abstracts. *J. of Information Science* 34(3), 308–326 (2008)
19. Porter, M.F.: An algorithm for suffix stripping. *Program*. 14(3), 130–137 (1980)
20. Xu, J., Croft, W.B.: Improving the Effectiveness of Information Retrieval with Local Context Analysis. *ACM Tran. on Information Systems* 18(1), 79–112 (2000)

On the Configuration of the Similarity Search Data Structure D-Index for High Dimensional Objects

Arnoldo José Müller-Molina and Takeshi Shinohara

Department of Artificial Intelligence
Kyushu Institute of Technology
Kawazu 680-4 Iizuka 820-8502, Japan
arnoldo@daisy.ai.kyutech.ac.jp,
shino@ai.kyutech.ac.jp

Abstract. Among similarity search indexes, the *D-index* introduced by Gennaro et al. in 2001 is regarded as an efficient metric access method. The performance of this index depends on several parameters, and their optimal configuration remains an open problem. We study two performance issues that occur when the D-index handles high dimensional objects. To solve these problems, we introduce an optimization that simplifies the D-index. By doing this, we remove two configuration parameters and improve performance.

Keywords: similarity search, digital library, multimedia system, information retrieval, spatial index.

1 Introduction

New similarity search applications are constantly being developed, ranging from music similarity search to open source license violation detectors [1]. To support such development, metric access methods (MAMs) are valuable tools as they can be used by researchers to find suitable similarity measures for objects. Among centralized metric access methods, the D-index [2] is regarded as one of the fastest MAMs available. The index implements an *addressing scheme* based on the relationship between each object and a set of objects called *pivots* taken from the database. Pivoting techniques have been traditionally used to estimate a *lower bound* of the distance function employed [3,4,5]. The combination of these techniques is what makes this index unique.

The D-index requires several parameters whose configuration severely affects performance. The optimal configuration of these parameters still remains an open problem to date. In this paper, we study the structure of the D-index in the context of costly distance functions. Our main objective is therefore to minimize distance computations.

The D-index employs a *clustering* technique where data is divided into levels composed of buckets. Additionally, it employs a *pivot filtering* technique [5] used

to reduce dimensionality. When a large number of pivots is employed, one of the following issues arise:

- The D-index tries to distribute evenly data among a set of buckets however, in practice, about 80% to 99% of the data stays in one bucket.
- The search algorithm can access an exponential number of buckets with respect to the cardinality of the pivot set employed.

To correct these problems, we propose a simplification of the D-index that we call the SD-index (Simple D-index). Our optimization distributes the data into a larger number of buckets. Additionally, we propose an heuristic that reduces the number of buckets that need to be accessed.

Our main conclusion is that it is possible to increase performance by reducing the structure of the D-index to only one level. Additionally, our technique brings two benefits. First, two parameters are removed, and this makes the D-index more accessible to a wide range of users. Second, index creation becomes faster as bucket levels have been removed.

In Section 2, we review key concepts related to the D-index. In Section 3, our technique is discussed. In Section 4, we experimentally demonstrate the effectiveness of our method. Finally in Section 5, we point out future research objectives and the concluding remarks.

2 Background

2.1 Metric Space

Let $\mathcal{M} = (\mathcal{D}, d)$ be a *metric space* for a domain of objects \mathcal{D} and $d : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$, a total *distance function* that satisfies the following properties:

$$\begin{aligned} \forall x, y \in \mathcal{D}, d(x, y) &\geq 0 \\ \forall x, y \in \mathcal{D}, d(x, y) &= d(y, x) \\ \forall x, y \in \mathcal{D}, d(x, y) = 0 &\iff x = y \\ \forall x, y, p \in \mathcal{D}, d(x, y) &\leq d(x, p) + d(p, y) \end{aligned}$$

Given a collection $X \subseteq \mathcal{D}$, a *range query* for $q \in \mathcal{D}$ and a range r retrieves all the objects within distance r , that is the set $\{x \in X | d(x, q) \leq r\}$. A *nearest neighbor query* returns the k closest elements to the query object q , or the set $R \subseteq X$ such that $|R| = k$, for any $x \in R$ and any $y \in X - R : d(q, x) \leq d(q, y)$. In this paper, we focus on queries that combine both range and nearest neighbor queries.

2.2 Pivot Filtering

It is possible to reduce the number of distance computations executed by employing the well established pivot filtering technique [5]. This embedding helps

to reduce distance computations by exploiting the triangle inequality. To understand this, a well known property of the triangle inequality is useful. For a metric space $\mathcal{M} = (\mathcal{D}, d)$ and three objects $x, y, p \in \mathcal{D}$:

$$|d(p, x) - d(p, y)| \leq d(x, y). \tag{1}$$

The technique selects n *pivot* objects, and obtains the distances of all the objects to these pivots. These pre-computed distances can be used to obtain a lower bound of the real distance. Formally, for a metric space $\mathcal{M} = (\mathcal{D}, d)$ and a set of pivots $P = \{p_1, \dots, p_n\}$, the mapping $\Psi : (\mathcal{D}, d) \rightarrow (\mathbb{R}^n, L_\infty)$ can be defined:

$$\Psi(x) = (d(x, p_1), d(x, p_2), \dots, d(x, p_n)).$$

The resulting vector is called the *pivot vector*. A lower bound of d can be obtained with:

$$L_\infty(\Psi(x), \Psi(y)) = \max_{i=1}^n |d(x, p_i) - d(y, p_i)| \leq d(x, y).$$

This lower bound can be used in queries to eliminate objects if $L_\infty(\Psi(x), \Psi(y)) > r$, thus avoiding costly distance computations. Naturally, as the number of pivots increases the pruning effect improves. The technique has been studied by different authors in the past [5][6][3][7].

2.3 D-Index

The D-index was introduced by [8] and refined by [2]. The index partitions data into h levels and $2^{n_i} - 1$ buckets per level i . A function decides in which level and bucket any given object must be placed. The main characteristic of the index is that when a range query has a search radius up to some predefined ρ , only one bucket must be accessed per level. Figure 1 shows how the structure would look like for a particular example. Each square drawn with solid lines represents a bucket that is used to store objects. For levels 0 to 2, these buckets are called *separable buckets*. Dashed squares represent *exclusion buckets* that hold objects that are recursively partitioned to the next level. At the very bottom, the exclusion bucket E is employed to capture all the objects that cannot be stored in all the previous levels.

To construct the D-index, for each level i , n_i pivots p_{ik} ($0 \leq k < n_i$) are selected, and the corresponding median distances m_{ik} of all the objects that entered the level and each pivot p_{ik} are calculated. The *excluded middle partitioning* scheme [9] is employed to split the data. For $x \in \mathcal{D}$, the *split* function bps implements this partitioning scheme:

$$bps_i^{k,\rho}(x) = \begin{cases} 0 & \text{if } d(x, p_{ik}) \leq m_{ik} - \rho \\ 1 & \text{if } d(x, p_{ik}) > m_{ik} + \rho \\ - & \text{otherwise} \end{cases} \tag{2}$$

Objects within the interval $]m_{ik} - \rho, m_{ik} + \rho]$, belong to the *exclusion set*. In this case, the function returns the value “-”. *Separable sets* are created from objects

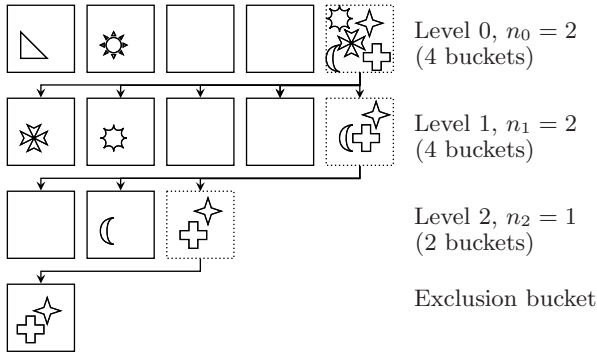


Fig. 1. Example of the structure of the D-Index.

for which the *bps* function returns 0 or 1. An example of the split function *bps* is illustrated in Figure 2. The exclusion set of width 2ρ contains the triangle-shaped objects. The innermost circle represents one separable set, containing circle-shaped objects. Finally, the square-shaped objects belong to the second separable set. Given a range query such that $r \leq \rho$, the function *bps* guarantees that only one separable set is accessed. This property is called the *separable property*.

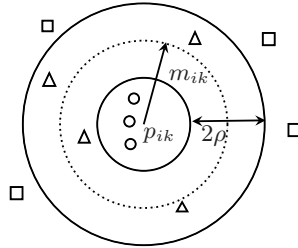


Fig. 2. Example of the *bps* function

The separation of data into three sets is not sufficient for most applications. Therefore, a higher order split function can be defined by concatenating the results b_k of each function $bps_i^{k,\rho}$. The joint n -split function is denoted by $bps_i^\rho(x)$ and the return value is a string $b = (b_0, \dots, b_{n_i-1})$ where $b_k \in \{0, 1, -\}$.

To store each bucket, an addressing technique that receives a string and returns a *bucket identification* number was proposed:

$$\langle b \rangle = \begin{cases} \sum_{k=0}^{n_i-1} 2^k b_k & \text{if } \forall k \ b_k \neq - \\ 2^{n_i} & \text{otherwise} \end{cases} \quad (3)$$

When no “-” elements are returned by any of the $bps_i^{k,\rho}$ functions, a number smaller than 2^{n_i} is returned. By using the function bps_i^ρ and the $\langle \cdot \rangle$ operator,

a bucket identification number can be assigned to each object $o \in \mathcal{D}$. From the previous discussion, the following bucket structure is created:

$$\begin{aligned}
 &B_{0,0}, B_{0,1}, \dots, B_{0,2^{n_0}-1} \\
 &B_{1,0}, B_{1,1}, \dots, B_{1,2^{n_1}-1} \\
 &\vdots \\
 &B_{h-1,0}, B_{h-1,1}, \dots, B_{h-1,2^{n_{h-1}}-1} \\
 &E
 \end{aligned}$$

Each bucket $B_{i,j}$ contains a *header* that holds the pivot vectors for each object stored in the bucket. The header is sorted by the first pivot, and the objects can be stored in an alternate location. This bucket layout was called *elastic bucket* [2].

Insert. Figure 3 shows the insertion procedure. For every level, the algorithm attempts to find a bucket. If a bucket is found, the object is inserted (line 4) and the process terminates. If all the levels are tried without success, the object is inserted into the exclusion bucket E (line 8).

Range search. Figure 4 shows the range search algorithm. The algorithm evaluates each level and reads one or more buckets depending on r . If r is not greater than ρ (line 4), two possibilities arise. In the case where both the separable set and the exclusion set are intersected by the query (line 5), one separable bucket is accessed. Otherwise, no separable bucket is accessed for the current level because the query region is contained exclusively in the exclusion set of the level.

When the radius is greater than ρ , more than one separable set must be accessed per level. Recall that when invoking $bps_i^{r-\rho}$, a string $b = (b_0, b_1, \dots, b_{n_i-1})$ is returned. To explore all the possible buckets, we must replace each b_k with value “-” alternately with 0 and 1, and generate all the possible bucket identification number combinations. This is accomplished by calling function s_i in line 10.

In Figure 5, the function s_i receives a string b initialized with each $b_k = 0$. This string is used as a template to generate all the possible bucket identification numbers. The function will process each split function $bps_i^{k,r-\rho}$ until a bucket identification number can be created. When $k < n_i$, the k th split function’s value is computed (line 7). If the resulting value is “-” (line 9), the bucket identification numbers where $b_k = 0$ and $b_k = 1$ are set in b , and the next split function is recursively computed. Otherwise, the value obtained by $bps_i^{k,r-\rho}$ is replaced and the next split function $k + 1$ is processed in line 12. When $k = n_i$, a bucket identification number is already in b and bucket $B_{i,\langle b \rangle}$ can be searched.

2.4 Parameter Configuration Issues

In this Section, we discuss problems related to the parameter configuration of the D-index. Our discussion assumes that search is CPU bounded, and therefore

```

Insertion Algorithm
1: function insert( $o \in \mathcal{D}$ )
2:   for  $i = 0$  to  $h - 1$  do
3:     if  $\langle bps_i^\rho(o) \rangle < 2^{n_i}$  then ▷  $o$  is not in  $E$ ?
4:       Store  $o$  in bucket  $B_{i, \langle bps_i^\rho(o) \rangle}$ 
5:       exit
6:     end if
7:   end for
8:   Store  $o$  in bucket  $E$ 
9: end function

```

Fig. 3. Insertion algorithm of the D-index

\mathcal{Q} denotes the query region
 q Query object, r Range

```

1: function RangeSearch( $q, r$ )
2:    $res \leftarrow \{\}$  ▷ Result set
3:   for  $i = 0$  to  $h - 1$  do
4:     if  $r \leq \rho$  then ▷ Search radius up to  $\rho$ 
5:       if  $\langle bps_i^{\rho-r}(q) \rangle < 2^{n_i}$  then ▷ Not exclusive containment in exclusion bucket
6:          $res \leftarrow res \cup \mathcal{Q} \cap B_{i, \langle bps_i^{\rho-r}(q) \rangle}$ 
7:       end if
8:     else ▷ Search radius is greater than  $\rho$ 
9:        $b \leftarrow (b_0, b_1, \dots, b_{n_i-1})$  where  $b_i = 0$ 
10:       $res \leftarrow res \cup s_i(q, r, k, b)$ 
11:     end if
12:   end for
13:   return  $res \cup (\mathcal{Q} \cap E)$ 
14: end function

```

Fig. 4. Range search algorithm

our objective is to minimize distance computations. From our experience [10,11], this is very common when handling trees or other high dimensional objects that utilize complex distance functions.

Dohnal *et al.* [2] pointed out that “generally” the number of pivots (or split functions) per level is assigned in decreasing order. In the following, we assume this and therefore level $i = 0$ has the greatest number of pivots. We also assume that the pivots used in addressing are the pivots using in filtering.

The D-index structure presents two problems when a large number of pivots is used. First, as the number of pivots increases, the probability that an object will fall into the exclusion bucket of the next level becomes higher. This is because as the number of pivots for a level increases, the probability of obtaining a “-” also increases. This follows from the definition of bps_i^ρ (Section 2.3). In this setting, the D-index structure degenerates, and the exclusion bucket E tends to contain a large percentage of the data. The upper layers, with the biggest amount of pivots,

$b = (b_1, b_2, \dots, b_{n_i})$ string
 k index to be changed
 v the new value for b_k

```

1: function  $set(b, k, v)$ 
2:   Change the  $k$ th element of  $b$  to  $v$ 
3:   return  $b$ 
4: end function
    $i$ : Current level,  $q$ : Query object,  $\mathcal{Q}$  query region
    $r$  Range,  $k$  Processing the  $k$ th element of  $b$ 
    $b$  String  $b = (b_0, b_1, \dots, b_{n_i-1})$  where  $b_i \in \{0, 1\}$ 
5: function  $s_i(q, r, k, b)$ 
6:   if  $k < n_i$  then                                ▷ Still creating bucket id.
7:      $bit \leftarrow bp s_i^{k, r-\rho}(q)$ 
8:     if  $bit == \text{“-”}$  then                            ▷ Process both  $b_k$ .
9:       return  $s_i(q, r, k + 1, set(b, k, 0)) \cup$ 
10:         $s_i(q, r, k + 1, set(b, k, 1))$ 
11:     else                                            ▷ Process only one  $b_k$ .
12:       return  $s_i(q, r, k + 1, set(b, k, bit))$ 
13:     end if
14:   else                                            ▷ Bucket id is complete.
15:     return  $\mathcal{Q} \cap B_{i, (b)}$                             ▷ Search bucket  $b$ 
16:   end if
17: end function

```

Fig. 5. Combination generator for queries $r > \rho$

store a small subset of the data. This represents a waste, as pivot vectors of the upper levels are the most helpful in reducing distance computations because the pivot set size is larger.

To illustrate this problem, in Table 1 the distribution of objects per level for different ρ and $h = 9$ in two real data-sets is displayed. The description of the data-sets is given in Section 4. The percentage of objects that pass to the next level is shown for each level. In this example, the D-index will search data in a limited number of buckets. As ρ grows bigger, the percentage of data stored in the exclusion bucket increases. This behavior can be observed with relatively large ρ and pivot sets. The large number of objects stored in the exclusion bucket suggests that D-index behaves like sequential search on pivot vectors [6]. A solution is to reduce the number of pivots. For fast distance functions this is possible. Nevertheless, as the complexity of a distance function increases, it is desirable to have more pivots so that the lower bound becomes more precise.

The second problem occurs when ρ is small. In this case, the probability that an object will fall into the exclusion zone will be also reduced. Recall from Section 2.3 that when $r > \rho$, the function s_i must be called. This implies that at most $\sum_{i=0}^{h-1} 2^{n_i}$ buckets would have to be accessed. The generation of bucket identification number combinations by s_i is an expensive operation.

Matching with a large set of pivots can negatively affect the distribution of buckets. On the other hand, if distribution is improved by changing ρ , exponential I/O

Table 1. Level distribution for each level h . The column “%” contains the percentage of objects that fall into the next level $i + 1$. The last row represents the percentage of objects stored in the exclusion bucket. The column n_i contains the number of pivots for level i .

i	MTD $\rho = 8$		TED $\rho = 7$	
	n_i	%	n_i	%
0	16	95.06	30	99.99
1	15	93.26	28	99.99
2	13	92.16	26	99.99
3	12	91.65	24	99.99
4	11	90.61	22	99.99
5	10	90.11	20	99.99
6	9	89.68	19	99.99
7	8	88.88	18	99.99
8	7	88.14	17	99.99

occurs. To avoid these problems, we propose a simplification of the D-index that attempts to alleviate these issues.

3 Proposed Technique

In this Section, we introduce the SD-index, a simplification of the D-index that attempts to solve the problem described in the previous section. The main objective of our modification is to avoid the creation of levels, as they reduce the size of the pivot vectors, while maintaining the bucket addressing scheme used in the D-index. Both objectives can be achieved by setting ρ to 0.

Let us explain the implications of this change. First, as $bps_i^{k,r-\rho}$ never returns “-”, the exclusion bucket disappears. As a consequence, only one level $i = 0$ exists. The bucket identification number function (3) remains unmodified.

The insert algorithm (Figure 3) can be reduced to only line 4. Likewise, the range search algorithm of Figure 4 can be reduced to lines 9 and 10. The main problem our modification introduces is that s_i must access at most 2^{n_0} buckets. As mentioned, this issue occurs to the D-Index when $r > \rho$. To alleviate this, we implemented two heuristics that try to reduce the amount of buckets that must be accessed.

3.1 Heuristic 1: Bucket Visiting Order

The first heuristic is designed to guide s_i to buckets that are more likely to contain objects close to the query. Our current implementation works only for discrete numbers but it can be extended for real valued numbers. We pre-compute n_0 arrays f_k ($0 \leq k < n_0$) that contain the distance distribution for each pivot p_{k0} . The value $f_k[i]$ represents the number of objects that are at distance i from pivot p_{k0} . For $o \in \mathcal{D}$ and a range r , the functions $l_k(o, r)$ and $r_k(o, r)$ are defined as:

$$l_k(o, r) = \sum_{j=d(p_{0k}, o)-r}^{m_{0k}} f_k[j],$$

$$r_k(o, r) = \sum_{j=m_{0k}}^{d(p_{0k}, o)+r} f_k[j].$$

During block number generation in s_i , if $l_k(o, r) \geq r_k(o, r)$ then the recursion will proceed first with $b_k = 0$ and then $b_k = 1$ will be executed. Otherwise $b_k = 1$ will go first followed by $b_k = 0$. With this, we intend to visit first bucket identification numbers that have more objects intersecting the query rectangle on a certain “side” of the median.

3.2 Heuristic 2: Empty Bucket Pruning

A total of 2^{n_0} bucket identification numbers can exist. Nevertheless, in practice only a subset of all the possible buckets is actually used. The second heuristic is designed to avoid accessing block identification numbers that do not exist. We

i : Current level, q : Query object, \mathcal{Q} : query region
 r : Range, k : Processing the k th element of b
 b : String $b = (b_0, b_1, \dots, b_{n_i-1})$ where $b_i \in \{0, 1, -\}$

```

1: function  $s'_i(q, r, k, b)$ 
2:   if not  $(b_0, \dots, b_k) \in pr_k$  then
3:     return  $\{\}$  ▷ Heuristic 2
4:   end if
5:   if  $k < n_i$  then
6:      $bit \leftarrow bp s_i^{k, r-\rho}(q)$ 
7:     ▷ Multiple  $b_k$  to process
8:     if  $bit == \text{“-”}$  then
9:       ▷ Heuristic 1
10:      if  $l_k(q, r) \geq r_k(q, r)$  then
11:        return  $s_i(q, r, k+1, set(b, k, 0)) \cup$ 
12:           $s_i(q, r, k+1, set(b, k, 1))$ 
13:      else
14:        return  $s_i(q, r, k+1, set(b, k, 1)) \cup$ 
15:           $s_i(q, r, k+1, set(b, k, 0))$ 
16:      end if
17:    else ▷ One  $b_k$  to process
18:      return  $s_i(q, r, k+1, set(b, k, bit))$ 
19:    end if
20:  else
21:    return  $\mathcal{Q} \cap B_{i, (b)}$  ▷ Search bucket  $b$ 
22:  end if
23: end function

```

Fig. 6. Combination generator for SD-index

keep n_0 sets pr_k , each of them containing strings of size $k + 1$ ($0 \leq k < n_0$). At insertion time, for each string $b = (b_1, b_2, \dots, b_{n_0})$ generated by bps_i^ρ , we extract the sub-strings $a_0 = (b_0), a_1 = (b_0, b_1), \dots, a_{n_0-1} = (b_0, \dots, b_{n_0-1})$. Each a_k will be inserted into pr_k . In s_i , for the k th split function, we can verify if the block prefix we have constructed belongs to pr_k . If it belongs to pr_k , then we can proceed with the given combination, otherwise search should stop at this point. This heuristic can be implemented with a trie, or by using a *bloom filter* [11] for each pr_k . A better approach would be perhaps to do a sequential scan on the binary strings and use the hamming distance as an disk read order criteria. This would guarantee that the number of shared partitions is maximum.

In Figure 6, we rewrite the function s_i to add the heuristics described above. Heuristic one is located between lines 10 and 16, and heuristic two is located at line 2.

If only distance computations are taken into account, the SD-index is optimal because all the objects are mapped with an equal number of pivots n_0 . As n_0 is the largest number of pivots for any level, and only a small percent of the data is stored in the first level, the D-index is likely perform more distance computations for large ρ . On the other hand, as ρ reaches zero, more data is stored in the upper levels of the index. In this case, the D-index must access at most $\sum_{i=0}^{h-1} 2^{n_i}$ buckets. The SD-index will access at most 2^{n_0} buckets.

4 Experiments

We evaluated the SD-index on two real data-sets:

- *MTD*: 350000 trees compared by the *mtd* [12] tree distance metric ($O(n^2)$). The trees are binary program fragments [10] used by an open source license violation detector. Each tree has up to 500 nodes.
- *TED*: 245000 trees compared by the tree edit distance [13] metric ($O(n^3)$). This data-set is a subset of *MTD* and includes trees that have up to 20 nodes.

4.1 Computational Costs

We have chosen distance functions with relatively high complexity that require many pivots to become feasible to match. As a rough estimation, one distance computation takes as long as the computation of the L_2 distance on a pair of integer vectors of 6500 dimensions for *mtd* and 890000 dimensions for *ted*.

All our experiments present the averaged results of querying the data-set with an exclusive set of 400 different objects. The maximum range query radii was selected so that roughly 20% of the data is returned. The distance distributions between data and query sets are illustrated in Figure 7.

We measured disk access operations, distance computations and total execution time. We implemented the algorithms in Java and executed them on Intel Xeon computers running at 2.66GHz and with 4GB of memory on Ubuntu Linux

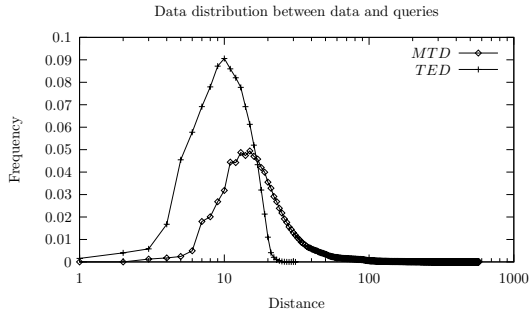


Fig. 7. Distance distribution for data-sets *TED* and *MTD*

7.10. The GPL source code along with the data-sets can be downloaded^[1]. All the indexes were built on top of the same framework^[2].

4.2 Indexing

Besides D-index and SD-index, we have employed sequential search as a baseline. Additionally, we used the *P+tree* [14] combined with pivot filtering to compare our technique with a relatively new high dimensional index. After several trials, the best setting we found for parameter *od* for P+tree was 12. We could not execute the D-index for *TED* with $\rho < 7$ because s_i performed worse than sequential search. This is due to the exponential number of read operations required when $\rho > r$.

By running several trials, we found that the best number of pivots is 16 for dataset *MTD* and 30 for dataset *TED*. Pivot sets for SD-index and P+tree are equal. For D-Index, different sets were generated for each level, and nine levels were created. Table 1 shows the order for each level of D-index. The pivot selection strategy employed is called *incremental selection* and it is described by [15]; the parameters $l = 1000$ and $m = 1000$ were used.

Table 2 shows index creation time without considering pivot selection. For different parameters of D-index, we present an average of the running times. P+tree took less time in dataset *TED* because the current implementation avoids duplicate computations. With some effort SD-index could also be optimized in the same way. On the contrary, the D-index cannot be optimized because most of the time is spent creating the levels, a step that cannot be avoided.

Table 3 shows the distribution of objects on the buckets. Clearly, SD-index is able to distribute the data better among the buckets. In the case of D-index, as ρ becomes smaller, buckets have less objects and the standard deviation becomes smaller.

¹ <http://www.obsearch.net/pr1.tar.gz>

² <http://www.obsearch.net>

Table 2. Index creation time in hours

Index	<i>MTD</i>	<i>TED</i>
<i>D</i>	0.80	97
<i>SD</i>	0.15	13.75
<i>P+</i>	0.11	7

Table 3. Object Distribution per Bucket. Average, σ , and minimum and maximum number of objects per bucket. The “Count” column contains the total number of buckets created for SD-index and different ρ values of D-index.

	Index	Avg.	σ	Min	Max	Count
<i>TED</i>	SD	5	119	1	21 534	46681
	<i>D</i> (7)	61250	122495	1	244993	4
<i>MTD</i>	SD	55	421	1	15 241	6338
	<i>D</i> (1)	85	1607	1	99980	4067
	<i>D</i> (4)	843	13254	1	269502	414
	<i>D</i> (8)	2477	25928	1	307854	141
	<i>D</i> (12)	4365	35875	1	321122	80

4.3 Performance Results

In Figure 8, the results for *MTD* show the performance of D-index for different values of ρ . As expected when $\rho = 1$, D-index reduces the amount of computations because the pivot vectors are larger. This is because a greater number of objects is inserted into the upper levels of the structure. SD-index and P+tree benefit from an optimal pivot vector size, executing less distance calculations. Regarding bucket access counts, $\rho = 12$ and $\rho = 8$ are accessing data in a small number of buckets that contain a large percentage of the data. For example, in Table 1, the exclusion bucket of $\rho = 8$ holds 88% of the data. Although few buckets are accessed, note that data read is considerable because the buckets that are accessed have a large number of pivot vectors. From an IO perspective, the P+tree is lagging behind D-index and SD-index because it performs many disk access operations. For D-index, $\rho = 12$ performs better than $\rho = 1$ because $r \leq \rho$ and this implies that no more than one bucket is being accessed per level. Recall from Table 1 that in this configuration, D-Index is working like sequential search on the pivot vectors.

As the distance computation count reflects, $\rho = 1$ minimizes distance computations. As expected, smaller ρ increases the number of objects that are stored in the upper levels of the structure. In those levels, the lower bound l_∞ has better resolution because the pivot set cardinality is the largest. Time-wise, $\rho = 1$ lags behind other configurations because most of the time is spent accessing an exponential number of buckets. The heuristics introduced in Section 3 help SD-index to obtain the best of the two extremes by both reducing distance computations and execution time at the expense of moderate I/O increase.

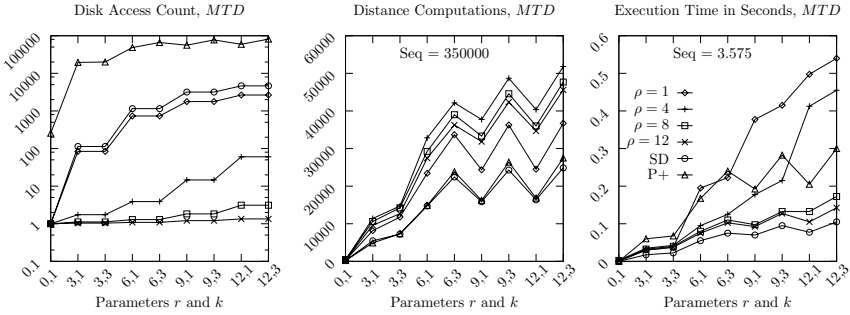


Fig. 8. Results for data-set *MTD*. Performance comparison among SD-index, P+tree and different ρ configurations of the D-index.

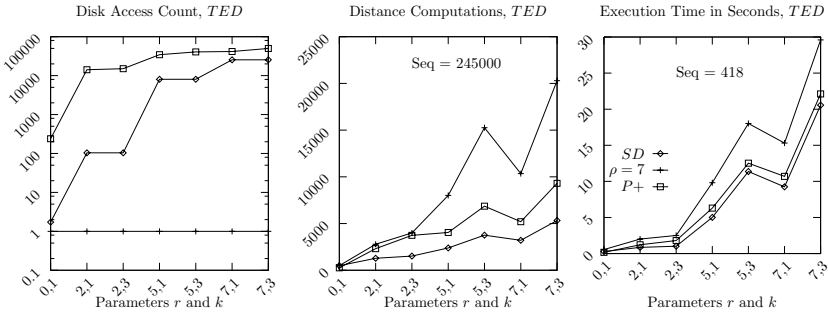


Fig. 9. Results for data-set *TED*. Performance comparison among SD-index, P+tree and a $\rho = 7$ D-index.

Figure 9 shows the results for data-set *TED*. As in data-set *MTD*, SD-index outperforms D-index and P+tree in distance computations and execution time. Although 3x to 4x improvements can be observed for distance computations, the execution performance is only improved 1.5x. This is partially caused by the overhead of accessing secondary storage. In general, distance function cost dominates the execution. D-index performed worse than sequential search for $\rho < 7$. In fact, we could not complete any of the experiments in an acceptable time frame because for $\rho < r$, a large number of buckets was read by function s_i .

4.4 Discussion

When ρ is large, D-index resembles a sequential search on the pivot projection. The fact that the D-index performs well for large ρ is not surprising as it has been experimentally verified that sequential search based on pivot vectors outperforms other hierarchical methods under certain circumstances [6]. The P+tree performs the largest amount of read operations in all the experiments. On the other hand, SD-index and P+tree are very close to each other in high dimensions because

their pivot set cardinality is the largest possible. The advantage obtained by SD-index comes from heuristic one, as we descend first to buckets that are more likely to have more objects.

5 Conclusion

We have studied the D-index in the context of high dimensional objects. Our modification, the SD-index, consistently outperforms the D-index in different configurations. Additionally, it guarantees that fewer distance computations will be executed as the pivot set cardinality is maximized.

The D-index cannot scale well for small ρ and range $r \geq \rho$ and when the number of pivots is large. In this case, the index must access an exponential number of buckets that depends on the number of pivots. In our experiments D-Index with 30 pivots was considerably slower than plain sequential search of a very costly distance function. On the other hand, if we set a large ρ such that $r \leq \rho$, the index degenerates into a sequential search.

References

1. Müller-Molina, A.J., Shinohara, T.: Fast approximate matching of programs for protecting libre/open source software by using spatial indexes. In: SCAM 2007, pp. 111–122. IEEE Computer Society, Washington (2007)
2. Dohnal, V., Gennaro, C., Savino, P., Zezula, P.: D-index: Distance searching index for metric data sets. *Multimedia Tools Appl.* 21(1), 9–33 (2003)
3. Micó, L., Oncina, J., Vidal, E.: An algorithm for finding nearest neighbours in constant averagetime with a linear space complexity. In: *Recognition Methodology and Systems*, pp. 557–560 (1992)
4. Ruiz, E.V.: An algorithm for finding nearest neighbours in (approximately) constant average time. *Pattern Recogn. Lett.* 4(3), 145–157 (1986)
5. Shinohara, T., Ishizaka, H.: On dimension reduction mappings for approximate retrieval of multi-dimensional data. In: *Progress in Discovery Science*, London, UK, pp. 224–231. Springer, Heidelberg (2002)
6. Filho, R.F.S., Traina, A.J.M., Caetano Traina, J., Faloutsos, C.: Similarity search without tears: The omni family of all-purpose access methods. In: *Proceedings of the 17th International Conference on Data Engineering*, pp. 623–630. IEEE Computer Society, Washington (2001)
7. Micó, M.L., Oncina, J., Vidal, E.: A new version of the nearest-neighbour approximating and eliminating search algorithm (aesa) with linear preprocessing time and memory requirements. *Pattern Recogn. Lett.* 15(1), 9–17 (1994)
8. Gennaro, C., Savino, P., Zezula, P.: Similarity search in metric databases through hashing. In: *MIR 2001*, pp. 1–5. ACM, New York (2001)
9. Yianilos, P.N.: Excluded middle vantage point forests for nearest neighbor search. Technical report, NEC Research Institute, Princeton, NJ (1998)
10. Müller-Molina, A.J., Shinohara, T.: On approximate matching of programs for protecting libre software. In: *CASCON 2006*, pp. 275–289. ACM Press, New York (2006)

11. Bloom, B.H.: Space/time trade-offs in hash coding with allowable errors. *ACM Commun.*, 422–426 (1970)
12. Müller-Molina, A.J., Hirata, K., Shinohara, T.: A tree distance function based on multi-sets. In: Chawla, S., Washio, T., Minato, S.-i., Tsumoto, S., Onoda, T., Yamada, S., Inokuchi, A. (eds.) *PAKDD 2008*. LNCS, vol. 5433, pp. 87–98. Springer, Heidelberg (2009)
13. Demaine, E., Mosez, S., Rossman, B., Weimann, O.: An optimal decomposition algorithm for tree edit distance. In: *ALP*, pp. 146–157. Springer, Heidelberg (2007)
14. Zhang, R., Ooi, B.C., Tan, K.L.: Making the pyramid technique robust to query types and workloads. In: *ICDE 2004*, p. 313. IEEE Computer Society, Washington (2004)
15. Zezula, P., Amato, G., Dohnal, V., Batko, M.: *Similarity Search: The Metric Space Approach*. Springer, Secaucus (2005)

Automatic Chinese Text Classification Using N-Gram Model

Show-Jane Yen¹, Yue-Shi Lee¹, Yu-Chieh Wu¹,
Jia-Ching Ying², and Vincent S. Tseng²

¹ Dept. of Computer Science and Information Engineering, Ming Chuan University
5 The-Ming Rd., Gwei Shan District, Taoyuan County 333, Taiwan

{sjyen, leeys}@mail.mcu.edu.tw

² Dept. of Computer Science and Information Engineering, National Cheng Kung University
1 University Rd., Tainan City 701, Taiwan

tsengsm@mail.ncku.edu.tw, jashying@idb.csie.ncku.edu.tw

Abstract. Automatic Chinese text classification is an important and well-known research topic in the field of information retrieval and natural language processing. However, past researches often ignore the problem of word segmentation and the relationship between words. This paper proposes an N -gram-based language model for Chinese text classification which considers the relationship between words. To prevent from the out-of-vocabulary problem, a novel smoothing method based on logistic regression is also proposed to improve the performance. The experimental result shows that our approach outperforms the previous N -gram-based classification model above 11% on micro-average F-measure.

Keywords: Text Classification, N -gram, Feature Selection, Word Segmentation, Logistic Regression.

1 Introduction

In recent years, there are more and more digital text services on the internet with the rapid growth of the World Wide Web. Examples include Google, Yahoo, etc. Search engine does not only provide user with the information, but also manage the information effectively. In order to make it easy, many search engine systems classify the texts in advance to organize their taxonomy. The methods of automatic text classification help people to classify the documents. At present, many techniques have been proposed to deal with the text classification problems. Yang and Liu [21] experiment on the news text of Reuter. They reported that support vector machine (SVM) and k -nearest neighbor (KNN) classifiers achieved the best accuracy in comparison to the other four methods. Sebastian [14] pointed out that SVM had better performance in general case.

However, support vector machine classifiers still have many problems. First, it is designed to solve the binary-class classification problem. Therefore it should be converted to handle multi-class classification problem. Second, the support vector

machine classifiers generally could not generate probability. Most support vector machine classifiers represent the similarity between the class and the text using the cosine similarity.

Unlike probability, the cosine similarity is very abstractive to represent the similarity between the class and the text. These problems are solved by Tipping [17]. He combines the logistic regression in the support vector machine classifiers to generate the probability. Meanwhile, the multi-class classification problem also can be transformed into several binary-class classification problems. Nevertheless, a known problem with the Tipping's approaches is their relative inability to scale with large problems like text classification. Fortunately, Silva and Ribeiro [15] solved the problem by the method of dimension reduction.

Nonetheless, these standard approaches to Chinese text categorization has so far been using a document representation in a word-based "input space" [6], i.e. as a vector in some high (or trimmed) dimensional Euclidean space where each dimension corresponds to a word. The advantages of this method are effective and high feasibility, but it assumes each word is independent with each other. Hence we need to solve word segmentation first for Chinese text. Most text classifiers treat each feature are mutually independent, word segmentation in Chinese is a difficult problem and we could not ensure that every words which is segmented by some word segmentation process are mutually independent. To solve it, character level N -gram models [3][4][12][13][16] could be applied. These approaches generally outperform traditional word-based methods in term of accuracy. Moreover, these approaches can avoid the word segmentation step and model the word-by-word relations [5][11].

Cavnar and Trenkle [3] proposed a classifier for English categorization, named N -gram-based text categorization, which works with letters level feature. They used N -grams as features for traditional feature selection process and then deployed classifiers based on calculating feature-vector similarities. Cavnar and Trenkle's approach has a shortcoming. It was intractable when there is a large feature set, and standard feature selection approaches do not always cope well in such circumstances. For example, given sufficient features, the cumulative effect of uncommon features can still have an important effect on classification accuracy, even though infrequent features contribute less information than common features individually. Teahan and Harper [16] used a PPM (prediction by partial matching) model for text classification where they seek a model that obtains the best compression on a new document. Later, Peng, et al. [12][13] proposed a classifier for Chinese classification without the word segmentation step which works with letters character feature by means of using N -gram model. In their experiment, their approach outperforms Cavnar and Trenkle's approach [3], Support Vector Machine and naive Bayes in most case. Although Peng, et al. proposed a novel classification to improve the traditional classifier for Chinese text classification, the approach also involve in several problems. First, the research of Peng, et al. entirely does not process the feature selection. If the testing text is too large, the classifier maybe uses a lot of unnecessary conditional probability to estimate the probability of a document to assign some class. Because the N -gram model uses the expression (1) to estimate the probability of an event that a sequence of the words $w_1...w_T$ occur in a data set.

$$P(w_1...w_T) = \prod_{i=1}^T P(w_i | w_{i-n+1}...w_{i-1}) \quad (1)$$

Therefore, if there exists some N -gram $w_1...w_n$ whose conditional probability $P(w_n|w_1...w_{n-1})$ are extremely similar in each class, the N -gram $w_1...w_n$ would not affect the result of classification. But standard methods of feature selection cannot be used in this kind of approach.

In this paper, we propose a novel feature selection which fits N -gram model for Chinese text classification. First, we consider the difference of the conditional probability of some N -gram in each class to select significant features. Second, there are many kinds of method for N -gram smoothing to solve sparse data problems, but approaches [12][13] just use the Katz's back-off estimator [2] without further adjusting. Because the original purpose of N -gram model is not to classify a text and the smoothing estimator is the kernel of the model, we believe that the model should be smoother and more suitable real world event. The logistic regression model is a famous classification model, it use a sigmoid function to estimate the probability. Therefore, a novel smoothing estimator using logistic regression is presented in this paper. Third, past research studies [12][13] treat the whole document as a sequence of words. The basis ideal of N -gram model is that the probability happening in some word, w_i , relates to the first $n-1$ continuous words. This means that when we estimate the probability of some word, we perhaps consult it to the last words of last sentence. However, we can regard each sentence as an independent event. In accordance with this point, our approach adjusts the whole model.

2 Logistic-Regression-Based N-Gram Models

The framework of our approach is shown in Figure 1.

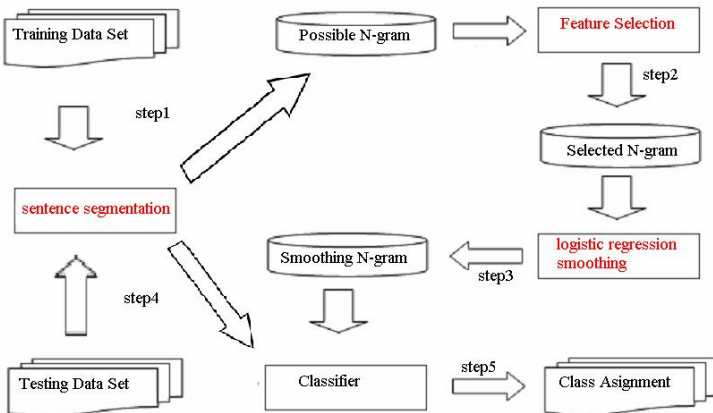


Fig. 1. The framework of our approach

The first step is to discover the frequency of every k -grams where $1 \leq k \leq n$ by scanning the training dataset and estimate each conditional probability of possible n -gram with maximum likelihood estimation. Next, we adopt chi-square statistic to

select important n -grams. In the third step, we build a logistic regression model to estimate each conditional probability $P(w_i|w_{i-n+1}...w_{i-1})$. According to the third step we can make sure each conditional probability $P(w_i|w_{i-n+1}...w_{i-1})$ converge within the interval 0 to 1(i.e. $0 < P(w_i|w_{i-n+1}...w_{i-1}) < 1$).

2.1 Significant N-Grams

According to the previous researches for feature selection [1], the chi-square statistic is better than mutual information and information gain in selecting useful features. However, the chi-square statistic cannot be applied to N -gram-based model without amendment. Therefore, we amend it to suit N -gram-based model. First, we use an example to expound what kind feature is significant in N -gram-based model as follows: Assume we wish to calculate a probability of a word sequence $s=w_1w_2w_3w_4$ occurred in two class, C_1 and C_2 , using bi-gram model and the distribution of w_1w_2 and w_j between C_1 and C_2 is given in Table 1.

Table 1. The distribution of w_1w_2 and w_j between C_1 and C_2

	#(w_1w_2)	#(w_j)
C_1	6	12
C_2	1	2

$$P_{C_1}(w_1w_2w_3w_4) = P_{C_1}(w_1) \times P_{C_1}(w_2 | w_1) \times P_{C_1}(w_3 | w_2) \times P_{C_1}(w_4 | w_3)$$

$$P_{C_2}(w_1w_2w_3w_4) = P_{C_2}(w_1) \times P_{C_2}(w_2 | w_1) \times P_{C_2}(w_3 | w_2) \times P_{C_2}(w_4 | w_3)$$

Because $P_{C_1}(w_2 | w_1) = P_{C_2}(w_2 | w_1)$, we believe $P_{C_1}(w_2 | w_1)$ and $P_{C_2}(w_2 | w_1)$ are unconcerned with the result of classification. Reversely, if there is an n -gram $w_1w_2...w_n$ whose distribution over all classes is more different than the distribution of $(n-1)$ -gram $w_1w_2...w_{n-1}$ over all classes, we can say that the n -gram $w_1w_2...w_n$ is significant. Based on this idea, we not only consider the relation and inverse relation between some n -gram and some class c , but also consider the relation and inverse relation between the n -gram $w_1w_2...w_n$ and the $(n-1)$ -gram $w_2...w_n$.

Because the n -gram $w_1w_2...w_n$ and the $(n-1)$ -gram $w_2...w_n$ have the distribution over all classes, their contingency table should be different.

Table 2. Contingency table with n -gram $w_1w_2...w_n$ and class c

		Assigning the document d to class c		Total
		Assign	Not Assign	
Occurrence of the n -gram $w_1w_2...w_n$ in some document d	Occur	A_I	B_I	A_I+B_I
	Not Occur	P_I	Q_I	P_I+Q_I
Total		A_I+P_I	B_I+Q_I	N_I

Table 3. Contingency table with $(n-1)$ -gram $w_1w_2\dots w_{n-1}$ and class c

		Assigning the document d to class c		Total
		Assign	Not Assign	
Occurrence of the $(n-1)$ -gram $w_1w_2\dots w_{n-1}$ in some document d	Occur	A_2	B_2	A_2+B_2
	Not Occur	P_2	Q_2	P_2+Q_2
Total		A_1+P_1	A_2+P_2	B_2+Q_2

In Tables 2 and 3, Let A_I be the number of times $w_1w_2\dots w_n$ and c co-occur, let B_I be the number of times $w_1w_2\dots w_n$ occurs without c , let P_I be the number of times c occurs without $w_1w_2\dots w_n$, let Q_I be the number of times neither $w_1w_2\dots w_n$ nor c , and let N_I be the total number of n -gram occurred in corpus.

The chi-square statistic of n -gram $w_1w_2\dots w_n$ with class c is represented as follows:

$$\chi^2(w_1w_2\dots w_n, c) = \frac{N_1 \times (A_1Q_1 - B_1P_1)^2}{(A_1 + P_1) \times (B_1 + Q_1) \times (A_1 + B_1) \times (P_1 + Q_1)}$$

Identically, the chi-square statistic of $(n-1)$ -gram $w_1w_2\dots w_{n-1}$ with class c is represented as follows:

$$\chi^2(w_1w_2\dots w_{n-1}, c) = \frac{N_2 \times (A_2Q_2 - B_2P_2)^2}{(A_2 + P_2) \times (B_2 + Q_2) \times (A_2 + B_2) \times (P_2 + Q_2)}$$

It is clear that if an n -gram $w_1w_2\dots w_n$ whose chi-square value is more different than chi-square value of $(n-1)$ -gram $w_1w_2\dots w_{n-1}$, then their distributions over all classes are more different. Therefore, we propose a method as follows:

$$F(w_1\dots w_n, C) = \frac{1}{|C|} \sum_{c \in C} F(w_1\dots w_n, c)$$

where $F(w_1\dots w_n, c) = \begin{cases} \frac{\chi^2(w_1\dots w_n, c)}{\chi^2(w_1\dots w_{n-1}, c)}, & \text{if } \chi^2(w_1\dots w_n, c) \geq \chi^2(w_1\dots w_{n-1}, c) \\ \frac{\chi^2(w_1\dots w_{n-1}, c)}{\chi^2(w_1\dots w_n, c)}, & \text{if } \chi^2(w_1\dots w_n, c) < \chi^2(w_1\dots w_{n-1}, c) \end{cases}$

The $\chi^2(w_1\dots w_n, c)$ is the chi-square statistic of n -gram $w_1w_2\dots w_n$ with class c , and the $\chi^2(w_1\dots w_{n-1}, c)$ is the chi-square statistic of $(n-1)$ -gram $w_1w_2\dots w_{n-1}$ with class c . We can determine them by contingency table as follows: If the significance probability (i.e. p-value) of the chi-square statistic of some n -gram $w_1w_2\dots w_n$ is less than the level of significance, the n -gram is called significant n -gram. In experiment, we set the level of significance equal to 5%, 10%, 15%, 20%, and 25%. Therefore, we can calculate the significance probability of each n -gram. If the significance probability of the n -gram is greater than the level of significance we set, we will not consider the n -gram into the classifier.

2.2 Smoothing Estimator for N-Gram Model

N -gram model is a well-known technique in natural language processing domain. As we know, the goal of language modeling is to predict the probability of natural word sequences; or more simply, to put high probability on word sequences those actually occur (and low probability on word sequences that never occur). Because of the heavy tailed nature of language one is likely to encounter novel n -grams that were never witnessed during training. Therefore, some mechanism for assigning non-zero probability to novel n -grams is a central and unavoidable issue [10]. In recent year, the combining estimation has the best performance in most case. One standard approach of combining estimation to smoothing probability estimates to cope with sparse data problems is to use some sort of linear interpolation estimator.

$$P_{\hat{\theta}}(w_i | w_{i-n+1} \dots w_{i-1}) = \lambda_1 P(w_i) + \lambda_2 P(w_i | w_{i-1}) + \dots + \lambda_n P(w_i | w_{i-n+1} \dots w_{i-1})$$

Where $0 < \lambda_i < 1$ and $\sum_i \lambda_i = 1$. This parameter can be set automatically using the Expectation-Maximization (EM) algorithm [10].

However, the estimation of some probability by some linear way is not perfect enough. Especially, the occurrence of some n -gram is a kind of response of a human been. Most statistician believe that logistic regression can solve the problem which the estimation is a kind of response of a human been. The logistic regression assumes that there is a sigmoid function, called response function, to represent the probability of occurrence of a response as follows:

$$P(r = 1) = \frac{e^{\alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n}}{1 + e^{\alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n}}$$

In the following, we take an example to describe the logistic regression modeling to estimate the conditional probability of bi-grams in some class c . Table 4 represents the response variable and independent variables.

Table 4. The response variable and independent variables

Response Variable		Independent Variables	
$\#(w_{i-1}w_i)$	$\#(w_{i-1})$	$P_{ML}(w_i w_{i-1})$	$P_{ML}(w_i)$
n_1	r_1	x_{11}	x_{21}
n_2	r_2	x_{12}	x_{22}
...
n_j	r_j	x_{1j}	x_{2j}

We can treat $\#(w_{i-1})$ as the frequency we observe and $\#(w_{i-1}w_i)$ as the frequency of the response which bi-gram is occurred. Assume there is a probability p_i to represent the probability of the response which bi-gram is occurred, the joint probability density function is

$$\prod_{i=1}^j \binom{n_i}{r_i} p_i^{r_i} (1 - p_i)^{n_i - r_i}$$

According to the assumption of the logistic regression $p_i = \frac{1}{1 + e^{\alpha_0 + \alpha_1 x_{i1} + \alpha_2 x_{i2}}}$, the joint probability density function can be written as

$$\begin{aligned} & \prod_{i=1}^j \binom{n_i}{r_i} \left(\frac{e^{\alpha_0 + \alpha_1 x_{i1} + \alpha_2 x_{i2}}}{1 + e^{\alpha_0 + \alpha_1 x_{i1} + \alpha_2 x_{i2}}} \right)^{r_i} \left(\frac{1}{1 + e^{\alpha_0 + \alpha_1 x_{i1} + \alpha_2 x_{i2}}} \right)^{n_i - r_i} \\ &= \prod_{i=1}^j \binom{n_i}{r_i} (e^{\alpha_0 + \alpha_1 x_{i1} + \alpha_2 x_{i2}})^{r_i} (1 + e^{\alpha_0 + \alpha_1 x_{i1} + \alpha_2 x_{i2}})^{-n_i} \end{aligned}$$

Then we estimate the parameters (i.e. $\alpha_0, \alpha_1, \alpha_2$) using maximum likelihood estimation and produce the estimator of the conditional probability as follows:

$$P_{lo}(w_i | w_{i-1}) = \frac{e^{\alpha_0 + \alpha_1 P_{ML}(w_i | w_{i-1}) + \alpha_2 P_{ML}(w_i)}}{1 + e^{\alpha_0 + \alpha_1 P_{ML}(w_i | w_{i-1}) + \alpha_2 P_{ML}(w_i)}}$$

Unlike linear interpolation, logistic regression is a sigmoid function. Therefore the parameters (i.e. $\alpha_0, \alpha_1, \alpha_2$) need not to be constrained. Besides, logistic regression has another advantage that it will not over-evaluate conditional probability which is estimated to be 0 by ML estimation as follows Figure 2.

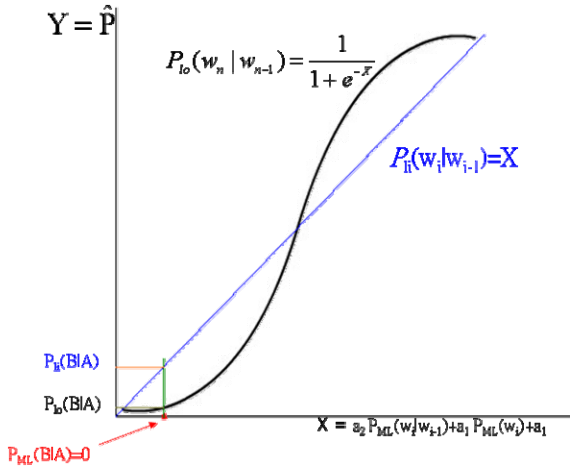


Fig. 2. Comparison of logistic regression and linear interpolation

As we see, if the frequency of AB is equal to zero but the frequency of B is very high, the conditional probability P(B|A) estimated by the linear interpolation is more over-evaluate than the conditional probability P(B|A) estimated by the logistic regression.

2.3 Language Models as Text Classifiers

Unlike the approach of Peng, et al. [12], our approach considers that a document is a set of sentences, and a sentence is a word sequence. Because we believe that calculating the probability of a word considering the end of last sentence is unreasonable. There is no proof the end of last sentence would affect the probability of a word in the start of a sentence.

Assume we wish to classify a document $d = \{s_1, s_2, \dots, s_{|d|}\}$ into a category $c \in C = \{c_1, \dots, c_{|C|}\}$, where $s_i = w_1 w_2 \dots w_{|s_i|}$. A natural choice is to pick the category c that has the largest posterior probability given the text. That is,

$$c^* = \arg \max_{c \in C} \Pr(c | d)$$

Using Bayes rule, this can be rewritten as

$$c^* = \arg \max_{c \in C} P(c) \Pr(d | c)$$

Considering each sentence is independent of other sentence in a document, this can be rewritten as

$$c^* = \arg \max_{c \in C} P(c) \prod_{i=1}^{|d|} \Pr(s_i | c)$$

Using N -gram model applying to classification as

$$c^* = \arg \max_{c \in C} P(c) \prod_{i=1}^{|d|} \left[\prod_{j=1}^{|s_i|} \Pr(w_{ij} | w_{i(j-n+1)} \dots w_{i(j-1)} \wedge c) \right]$$

As the above equation, our model considers sentence level.

3 Empirical Evaluation

The Chinese data set we used has been previously investigated in [19]. The corpus is a subset of the TREC-5 People's Daily news corpus published by the Linguistic Data Consortium (LDC) in 2001. The entire data set contains 164,789 documents on 59 categories, including IT, economy, etc. However, we could not use the entire data for our experiments, because there many categories include too few documents. Therefore, if we use the entire data for our experiments, there will be the imbalanced distribution data problem [22] in it. To void the imbalanced distribution data problem, we just choose particular of the entire data. The entire data set which we choose contains 105,733 documents on 22 categories.

3.1 Experimental Paradigm

Our goal is to classify text into a set of categories, and each document is assigned a single category. Each n -gram-based classification, such as Peng et al's approach [12],

is generally set the parameter n . For the experiments on Chinese data, we follow Peng et al’s approach and experiment on uni-gram model, bi-gram model, and tri-gram model. In each case, we randomly select 90% of documents evenly from among the remaining categories to form the training data. The testing set contains the other 10% documents of original data set. The training set and testing set do no overlap and do not contain replicated documents. Then we divide the training data into 9 parts uniformly. We pick 5 parts as training data 1, and pick all parts as training data 2.

3.2 Measuring Classification Performance

In the experiments, we measured classification performance by *micro-averaged* F-measure scores. To calculate the micro-averaged score, we formed an aggregate confusion matrix by adding up the individual confusion matrices from each category. The micro-averaged precision, recall, and F-measure can then be computed based on the aggregated confusion matrix.

For example, if there are $|C|$ categories included in training data set. For each category $c \in C$, we could determine a confusion matrix as follows.

Table 5. Confusion Matrix

Actual \ Predicted	Text belong to C_i	Text NOT belong to C_i
Text belong to C_i	a_i	p_i
Text NOT belong to C_i	b_i	q_i

The precision, recall, and F-measure of the category c are shown as follows.

$$precision = \frac{a_i}{a_i + b_i}, recall = \frac{a_i}{a_i + p_i}, F - measure = \frac{2 \times precision \times recall}{precision + recall}$$

The aggregated confusion matrix and the micro-averaged precision, recall, and F-measure of the category c are shown as follows.

Table 6. Aggregated Confusion Matrix

Actual \ Predicted	Text belong to $c \in C$	Text NOT belong to $c \in C$
Text belong to $c \in C$	$A = \sum_{i=1}^{ C } a_i$	$P = \sum_{i=1}^{ C } p_i$
Text NOT belong to $c \in C$	$B = \sum_{i=1}^{ C } b_i$	$Q = \sum_{i=1}^{ C } q_i$

$$\text{micro-averaged precision} = \frac{A}{A+B}, \text{micro-averaged recall} = \frac{A}{A+P}$$

$$\text{micro-averaged } F\text{-measure} = \frac{2 \times \text{micro-averaged precision} \times \text{micro-averaged recall}}{\text{micro-averaged precision} + \text{micro-averaged recall}}$$

3.3 Experimental Result and Analysis

Figures 3 and 4 give the results of comparisons of our approaches and Peng et al.'s approach with feature selection on training data set 2 in tri-gram and bi-gram models. In Figures 3 and 4, we can see that our approach outperform Peng et al.'s approach in most case. Tri-gram model outperform bi-gram model in most case. Remarkably, tri-gram model smoothing by logistic regression is more stable than bi-gram model smoothing by logistic regression, because the difference in F-measure of tri-gram model smoothing by logistic regression between best case and worst case is less than 25%, but the difference in F-measure of bi-gram model smoothing by logistic regression between best case and worst case is greater than 43%.

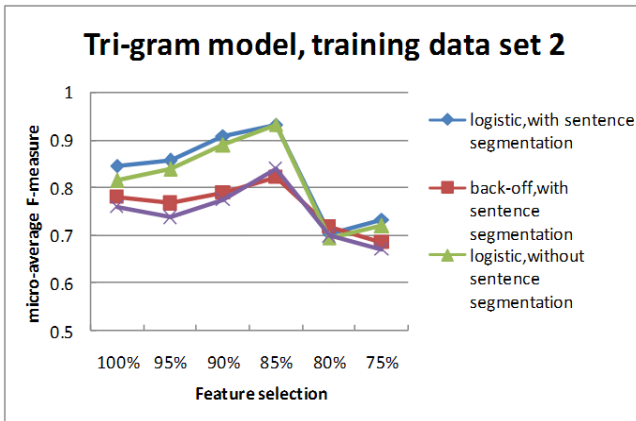


Fig. 3. Comparisons of our approaches and Peng et al.'s approach with feature selection on training data set 2 in tri-gram model

Figures 5 and 6 give the results of comparisons of our approaches and Peng et al.'s approach with feature selection on training data set 1 in tri-gram and bi-gram models. In Figures 5 and 6, we can see that our approach outperform Peng et al.'s approach in most case. There are evidences that the F-measure of back-off smoothing method will not be improved by feature selection. Therefore, the F-measure of back-off smoothing method we be deteriorated following the decreasing of the number of features.

However, if we just select 80% number of features to train the classifiers, the F-measure of back-off smoothing method looks few better then the F-measure of logistic regression smoothing method, because there are too many “significant” features are filtered out in this case. We also can find that the effect of sentence segmentation is marginally in bi-gram model, because bi-gram model considers that each word is observed is just related to the last word. In other words, there are rare bi-grams in head and tail of sentences.

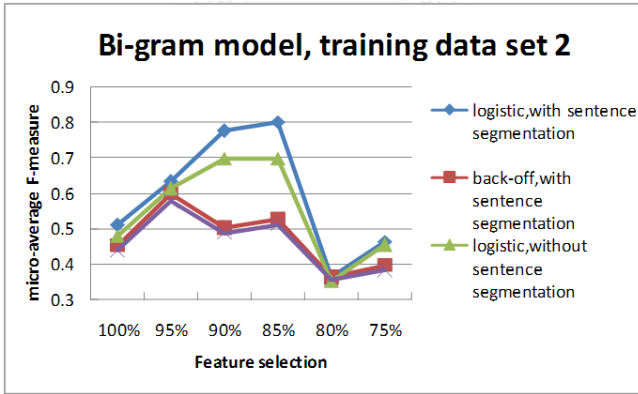


Fig. 4. Comparisons of our approaches and Peng et al.’s approach with feature selection on training data set 2 in bi-gram model

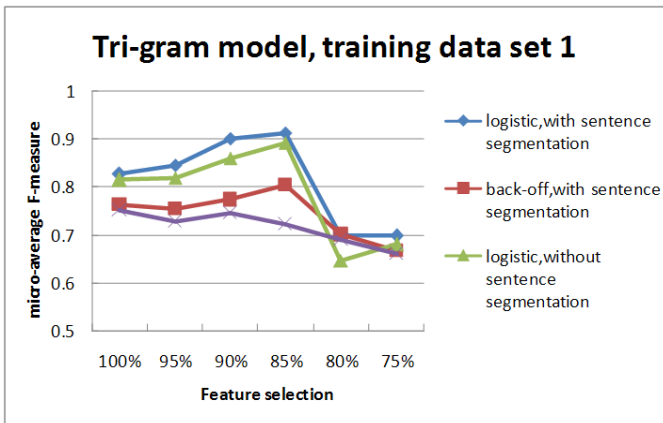


Fig. 5. Comparisons of our approaches and Peng et al.’s approach with feature selection on training data set 1 in tri-gram model

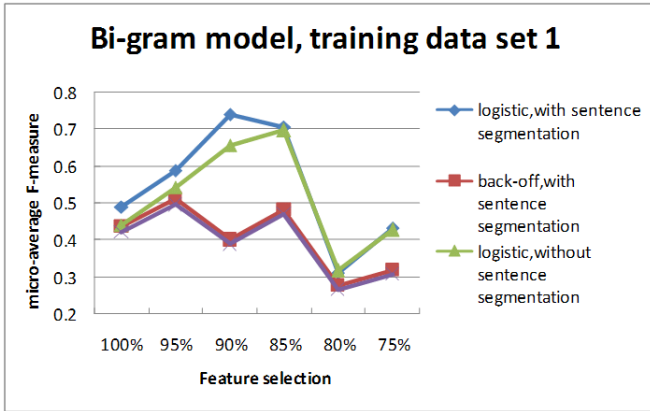


Fig. 6. Comparisons of our approaches and Peng et al.'s approach with feature selection on training data set 1 in bi-gram model

In the following, the results of the best cases of logistic regression smoothing method and back-off smoothing method on each training data set in bi-gram model and tri-gram model is represented in Table 7.

Table 7. The best cases of logistic regression smoothing method and back-off smoothing method

	Logistic		Back-off	
	bi-gram	tri-gram	bi-gram	tri-gram
training data set 2	0.800323	0.931512	0.596471	0.839820
training data set 1	0.739898	0.912540	0.513090	0.803358

In summary, first, generally our approach is better than Peng et al.'s approach in most case. Second, feature selection has no help to classifier smoothing by back-off, but very helpful to classifier smoothing by logistic regression. Third, sentence segmentation impact on F-measure is marginally in classifier smoothing by back-off. Fourth, logistic regression smoothing method generally outperform than back-off smoothing method in most case. Fifth, tri-gram model is much better than bi-gram model.

4 Conclusions and Future Work

In this paper, we have presented an n -gram-based model for Chinese text classification. We use the logistic regression to smooth the probability of n -gram. In our experiment, logistic regression smoothing outperform traditional back-off smoothing, because logistic regression has the ability to process unknown terms and it will not over-evaluate the conditional probability which originally is zero. Besides, we proposed a novel feature selection method which is suitable to N -gram-based model. In

our experiment, we prove that it could improve the F-measure in most case. Third, we consider a text as a set of sentences. According to our experimental result, this could improve system performance especially in the case which use tri-gram model or whose training data set is large enough.

In the future, we plane to find out a method to evaluate the relationship between sentences. Since we believe that considering a text as a set of sentences is not an optimal assumption and there exist some fact could respect the relationship between sentences. We also try to reduce the consumption of memory and calculative time of CPU on feature selection.

Acknowledgments. Research on this paper was partially supported by National Science Council grant NSC98-2221-E-130-019, NSC98-2622-E-130-001-CC3, and NSC98-2221-E-130-022.

References

1. Aizawa, A.: Linguistic Techniques to Improve the Performance of Automatic Text Categorization. In: 6th Natural Language Processing Pacific Rim Symposium, pp. 307–314 (2001)
2. Chen, S., Goodman, J.: An Empirical Study of Smoothing Techniques for Language Modeling. In: 34th Annual Meeting of the Association for Computational Linguistics, pp. 310–318 (1998)
3. Cavnar, W., Trenkle, J.: N-Gram-Based Text Categorization. In: 3rd Annual Symposium on Document Analysis and Information Retrieval, pp. 161–175 (1994)
4. Damashek, M.: Gauging Similarity with N-Grams: Language-Independent Categorization of Text. *Science* 267, 843–848 (1995)
5. Dumais, S., Platt, J., Heckerman, D., Sahami, M.: Inductive Learning Algorithms and Representations for Text Categorization. In: 7th International Conference on Information and Knowledge Management, pp. 148–155 (1998)
6. He, J., Tan, A., Tan, C.: On Machine Learning Methods for Chinese Document Categorization. *Applied Intelligence* 18, 311–322 (2003)
7. Jiang, E.: Learning to Semantically Classify Email Messages. In: 2nd International Conference on Intelligent Computing, pp. 664–675 (2006)
8. Joachims, T.: Text Categorization with Support Vector Machines: Learning with Many Relevant Features. In: European Conference on Machine Learning, pp. 137–142 (1998)
9. Lam, W., Ruiz, M., Srinivasan, P.: Automatic Text Categorization and Its Application to Text Retrieval. *IEEE Transactions on Knowledge and Data Engineering* 11, 865–879 (1999)
10. Manning, C.D., Schuetze, H.: *Foundations of Statistical Natural Language Processing*, pp. 191–227. MIT Press, Cambridge (2004)
11. Peng, F., Huang, X., Schuurmans, D., Cercone, N.: Investigating the Relationship of Word Segmentation Performance and Retrieval Performance in Chinese IR. In: 15th International Conference on Computational Linguistics, pp. 72–78 (2002)
12. Peng, F., Huang, X., Schuurmans, D., Wang, S.: Text Classification in Asian Languages without Word Segmentation. In: 6th International Workshop on Information Retrieval with Asian Languages, pp. 41–48 (2003)

13. Peng, F., Schuurmans, D.: Combining Naive Bayes and N-Gram Language Models for Text Classification. In: 25th European Conference on Information Retrieval Research, pp. 335–350 (2003)
14. Sebastian, F.: Machine Learning in Automated Text Categorization. *ACM Computing Surveys* 34, 1–47 (2002)
15. Silva, C., Ribeiro, B.: Scaling Text Classification with Relevance Vector Machines. In: IEEE International Conference on Systems, Man, and Cybernetics, pp. 4186–4191 (2006)
16. Teahan, W., Harper, D.: Using Compression-Based Language Models for Text Categorization. In: Workshop on Language Models for Information Retrieval, pp. 83–88 (2001)
17. Tipping, M.: Sparse Bayesian Learning and the Relevance Vector Machine. *Journal of Machine Learning Research* 1, 211–214 (2001)
18. Vapnik, V.: *The Nature of Statistical Learning Theory*. Springer, Heidelberg (1995)
19. Wu, Y.C.: Chinese Text Categorization with Term Clustering. M.S. Thesis, Mining Chuan University (2003)
20. Yang, Y.: An Evaluation of Statistical Approaches to Text Categorization. *Information Retrieval Journal* 1, 69–90 (1999)
21. Yang, Y., Liu, X.: A Re-examination of Text Categorization Methods. In: 22nd International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 42–49 (1999)
22. Yen, S., Lee, Y., Lin, C., Ying, J.: Investigating the Effect of Sampling Methods for Imbalanced Data Distributions. In: IEEE International Conference on Systems, Man, and Cybernetics, pp. 4163–4168 (2006)

Genetic Algorithms Evolving Quasigroups with Good Pseudorandom Properties

Václav Snášel¹, Jiří Dvorský¹, Eliška Ochodková¹, Pavel Krömer¹,
Jan Platoš¹, and Ajith Abraham²

¹ Department of Computer Science
Faculty of Electrical Engineering and Computer Science
VŠB – Technical University of Ostrava
17. listopadu 15, 708 33 Ostrava – Poruba, Czech Republic
{pavel.kromer, vaclav.snasel, jan.platos}@vsb.cz

² Center of Excellence for Quantifiable
Quality of Service, Norwegian
University of Science and Technology
O.S. Bragstads plass 2E,
N-7491 Trondheim, Norway
ajith.abraham@ieee.org

Abstract. Quasigroups are a well-known combinatorial design equivalent to more familiar Latin squares. Because all possible elements of a quasigroup occur with equal probability, it makes it an interesting tool for the application in computer security and for production of pseudorandom sequences. Prior implementations of quasigroups were based on look-up table of the quasigroup, on system of distinct representatives etc. Such representations are infeasible for large quasigroups. In contrast, presented analytic quasigroup can be implemented easily. It allows the generation of pseudorandom sequences without storing large amount of data (look-up table). The concept of isotopy enables consideration of many quasigroups and genetic algorithms allow efficient search for good ones.

1 Introduction

The need for random and pseudorandom sequences arises in many applications, e.g. in modeling, simulations, and of course in cryptography. Pseudorandom sequences are the core of stream ciphers. They are popular due to their high encryption/decryption speed. Their simple and cheap hardware design is often preferred in real-world applications. The design goal in stream ciphers is to efficiently produce pseudorandom sequences - keystreams (i.e. sequences that possess properties common to truly random sequences and in some sense are "indistinguishable" from these sequences).

The use of quasigroups and quasigroup string transformations is a recent but successful tendency in cryptography and coding [17]. With quasigroups in the hearth of advanced cryptosystems and hash functions, a need to find good quasigroups becomes hot topic.

Quasigroups and its applications in computer security were studied e. g. in [3]. A design of pseudorandom sequence generator (PRSG) based on quasigroup operation was presented in [4]. The authors have performed an extensive analysis of 2^{16} randomly chosen quasigroups of the orders 5, 6, 7, 8, 9 and 10 and concluded that different quasigroups produce pseudorandom sequences with different period (i.e. the number of elements after which the pseudorandom sequence starts to repeat). They have show that only a small number of quasigroups feature very large value of coefficient of period growth, a property that significantly affects the period of generated pseudorandom sequence [4]. This results encourage research of efficient methods for search for good quasigroups in the field of pseudorandom generators and cryptography.

Genetic algorithms are probably the most popular and wide spread member of the class of evolutionary algorithms (EA). EAs build a class of iterative stochastic search and optimization methods based on mimicking successful optimization strategies observed in nature [7,8,16,21]. The essence of EAs lies in the emulation of Darwinian evolution, utilizing the concepts of Mendelian inheritance for practical applications in computer science [7].

EAs operate with a population of artificial individuals (chromosomes) encoding potential problem solutions. Encoded individuals are evaluated using a carefully selected objective function which assigns a fitness value to each individual. The fitness value represents the quality (relative ranking) of each individual as a solution to given problem. Competing individuals explore in a highly parallel manner problem domain towards an optimal solution [16].

2 Quasigroups

Definition 1. A quasigroup is a pair (Q, \circ) , where \circ is a binary operation on (finite) set Q such that for all not necessarily distinct $a, b \in Q$, the equations

$$a \circ x = b \text{ and } y \circ a = b.$$

have unique solutions.

The fact that the solutions are unique guarantees that no element occurs twice in any row or column of the table for (\circ) . However, in general, the operation (\circ) is neither a commutative nor an associative operation.

Quasigroups are equivalent to more familiar Latin squares. The multiplication table of a quasigroup of order q is a Latin square of order q , and conversely, as it was indicated in [6,9,23], every Latin square of order q is the multiplication table of a quasigroup of order q .

Definition 2. Let $A = \{a_1, a_2, \dots, a_n\}$ be a finite alphabet, a $n \times n$ Latin square L is a matrix with entries $l_{ij} \in A$, $i, j = 1, 2, \dots, n$, such that each row and each column consists of different elements of A .

For $i, j, k \in A$ the ordered triple $(i, j; k)$ is used to represent the occurrence of element k in cell (i, j) of the Latin square. So a Latin square may be represented by the set $\{(i, j; k) \mid \text{entry } k \text{ occurs in cell } (i, j) \text{ of the Latin square } L.\}$

All reduced Latin squares of order n are enumerated for $n \leq 11$ [19]. Let L_n be the number of Latin squares of order n , and let R_n be the number of reduced Latin squares of order n . It can be easily seen that

$$L_n = n!(n - 1)!R_n.$$

Number of distinct Latin squares of a given order grows exceedingly quickly with the order and there is no known easily-computable formula for the number of distinct Latin squares. The problem of classification and exact enumeration of Latin squares of order greater than 11 probably still remains unsolved. Thus, there are more than 10^{90} quasigroups of order 16 and if we take an alphabet $L = \{0 \dots 255\}$ (i.e. data are represented by 8 bits) there are at least $256!255! \dots 2! > 10^{58000}$ quasigroups.

Multiplication in quasigroups has an important property; it is proven that each element occurs exactly q times among the products of two elements of Q , q^2 times among the products of three elements of Q and, generally q^{t-1} among the products of t elements of Q . Since there are q^t possible ordered products of t elements of Q , this shows that each element occurs equally often among these q^t products (see [10]).

2.1 Pseudorandom Sequence Generator Based on Quasigroups

The construction of pseudorandom sequence generator based on quasigroup operations was described in [4]. In this paper we use simplified design of PRSG to verify the proposed quasigroup optimization method.

Definition 3. Let (Q, \circ) be a quasigroup and Q^+ be a set of all nonempty words formed by the elements $q_i \in Q, 1 \leq i \leq n$.

For a fixed $a \in Q$ let the pseudorandom sequence y_1, y_2, \dots, y_n be defined as

$$\begin{aligned} y_1 &= a \circ a \\ y_2 &= a \circ y_1 \\ &\vdots \\ y_i &= a \circ y_{i-1} \end{aligned}$$

2.2 Isotopism of Quasigroups

Definition 4. Let $(G, \cdot), (H, \circ)$ be two quasigroups. An ordered triple (π, ρ, ω) of bijections π, ρ, ω of the set G onto set H is called an isotopism of (G, \cdot) upon (H, \circ) if $\forall u, v \in G, \pi(u) \circ \rho(v) = \omega(u \cdot v)$. Quasigroups $(G, \cdot), (H, \circ)$ are said to be isotopic.

We can imagine an isotopism of quasigroups as a permutation of rows and columns of quasigroup’s multiplication table.

Example 1. Consider a multiplication table for a quasigroup isotopic to the quasigroup of modular subtraction, with operation \circ defined as $a \circ b = (a + n - b) \bmod n$:

o	0	1	2	3
0	0	3	2	1
1	2	1	0	3
2	1	0	3	2
3	3	2	1	0

The table was created from table of modular subtraction. The second and the third rows were exchanged. Permutations π, ρ were identities and $\omega = [0213]$. A multiplication in this quasigroup can be illustrated by e.g. $1 \circ 0 = \omega(1) \circ 0 = 2 \circ 0 = 2$.

Starting with the quasigroup of modular subtraction, we can explore a large class of quasigroups isotopic to the quasigroup of modular subtraction [11,22]. This allows us to utilize quasigroups with very large number of elements without the necessity of their storage in program memory. The multiplication in such isotopic quasigroup is defined as follows:

$$a \circ b = \pi^{-1}((\omega(a) + n - \rho(b)) \bmod n). \quad (1)$$

We call the quasigroup defined by its multiplication formula and three selected permutations an *analytic quasigroup* [15,27].

The notion of analytic quasigroup enables efficient work with large quasigroups. Previous studies in this field used mostly quasigroups of small order [14], or just a small parts of certain quasigroup were utilized mainly as a key for Message Authentication Code. Such small quasigroups are represented as look-up tables in main memory. Larger quasigroup of order 2^{256} is used by NIST's SHA-3 competition [1] candidate, hash function EdonR [13].

The properties of one analytic quasigroup isotopic to the quasigroup of modular subtraction were studied in [15]. The quasigroup was created using three static functions that divided the sequence of n elements of the quasigroup into several parts. The parts were rotated in various directions and exchanged among themselves. It was shown that the investigated quasigroup has some faults in its properties.

2.3 Constructing Quasigroups Isotopic to the Quasigroup of Modular Subtraction

Consider a quasigroup on the length n defined by multiplication $a \circ b = (a + n - b) \bmod n$. Then three permutations π, ρ, ω must be chosen in order to implement isotopic quasigroup, whose multiplication will be defined by (1).

Obviously, there is $n!$ different permutations of n elements. Because three independent permutations are used to define any isotopic quasigroup, there are $n!n!n!$ possible choices of π, ρ and ω .

Permutations of elements cannot be sought for an analytic quasigroup directly, because its elements are not stored in memory. Instead, the permutation needs to be implemented as a function of an element of Q . One way to achieve this goal is the use of bit permutation.

¹ <http://csrc.nist.gov/groups/ST/hash/sha-3/index.html>

A quasigroup over a set of n elements requires $\log_2(n)$ bits to express each element. Each permutation of bits in the element representation represents also a permutation of all elements of the quasigroup (if n is a power of 2). Bit permutation can be implemented easily as a function of $q \in Q$.

The bit permutation is an elegant way of implementing permutations over n elements of Q . Although it enables us to explore only a fragment ($\log_2(n)!\log_2(n)!\log_2(n)!$) of all possible permutation triples over the quasigroup of n elements, it is useful because it does not require all n elements in main memory and therefore fits into the framework of analytic quasigroups.

Bit permutations are computationally more expensive than the static functions used to implement permutation in [15]. However, there are ongoing efforts to implement bit permutation instructions in hardware, which would improve the performance of the proposed algorithm significantly [12].

3 Genetic Algorithms

Genetic algorithms are generic and reusable population-based metaheuristic soft optimization method [5,16,21]. GAs operate with a population of chromosomes encoding potential problem solutions. Encoded individuals are evaluated using a carefully selected domain specific objective function which assigns a fitness value to each individual. The fitness value represents the quality of each candidate solution in context of the given problem. Competing individuals explore the problem domain towards an optimal solution [16]. The solutions might be encoded as binary strings, real vectors or more complex, often tree-like, hierarchical structures (subject of genetic programming [18]). The encoding selection is based on the needs of particular application area.

The emulated evolution is driven by iterative application of genetic operators. Genetic operators algorithmize principles observed in natural evolution. The crossover operator defines a strategy for the exchange of genetic information between parents (sexual reproduction of haploid organisms) while the mutation operator introduces the effect of environment and randomness (random perturbation of genetic information). Other genetic operators define e.g. parent selection strategy or the strategy to form new population from the current one. Genetic operators and algorithm termination criteria are the most influential parameters of every evolutionary algorithm. The operators are subject to domain specific modifications and tuning [21]. The basic workflow of the standard generational GA is shown in Fig. 1.

Many variants of the standard generational GA have been proposed. The differences are mostly in particular selection, crossover, mutation and replacement strategy [16].

In the next section, we present genetic algorithm for the search for good analytic quasigroups. It is an extended version of the initial GA for quasigroup evolution introduced in [27]. In this study, we present modified GA for quasigroup optimization with reengineered fitness function to find quasigroups that produce good pseudorandom sequences.

```

1 Define objective (fitness) function and problem encoding;
2 Encode initial population  $P$  of possible solutions as fixed length strings;
3 Evaluate chromosomes in initial population using objective function;
4 while Termination criteria not satisfied do
5   Apply selection operator to select parent chromosomes for reproduction:
    $sel(P_i) \rightarrow parent1, sel(P_i) \rightarrow parent2$ ;
6   Apply crossover operator on parents with respect to crossover probability to
   produce new chromosomes:
    $cross(pC, parent1, parent2) \rightarrow \{offspring1, offspring2\}$ ;
7   Apply mutation operator on offspring chromosomes with respect to
   mutation probability:  $mut(pM, offspring1) \rightarrow offspring1$ ,
    $mut(pM, offspring2) \rightarrow offspring2$ ;
8   Create new population from current population and offspring chromosomes:
    $migrate(offspring1, offspring2, P_i) \rightarrow P_{i+1}$ ;
9 end

```

Fig. 1. A summary of genetic algorithm

4 Genetic Search for Analytic Quasigroups

The genetic algorithm for the search for analytic quasigroup is defined by encoding of the candidate solutions and fitness function to evaluate chromosomes.

4.1 Encoding

As noted in section 2.3, any analytic quasigroup isotopic to quasigroup of modular subtraction is defined by three permutations. Such permutation triple represents a problem solution and should be mapped to one GA chromosome. Permutations can be for the purpose of genetic algorithms encoded using several strategies. In this study, we use random key encoding.

Random key (RK) encoding is an encoding strategy available for problems involving permutation optimization [24]. In random key encoding, the permutation is represented as a string of real numbers (random keys), whose relative position changes after sorting corresponds to the permutation index. An example of random key encoding is shown in (2).

$$\Pi_5 = \begin{pmatrix} 0.2 & 0.3 & 0.1 & 0.5 & 0.4 \\ 2 & 3 & 1 & 5 & 4 \end{pmatrix} \quad (2)$$

To encode a quasigroup (isotopic to the quasigroup of modular subtraction) of the length $n = 2^l$, we use a vector of $3l$ real numbers $v = (v_1, \dots, v_{l-1}, v_l, \dots, v_{2l-1}, v_{2l}, \dots, v_{3l})$. The vector is interpreted as three concatenated RK encoded permutations of the length l .

This encoding allows us to use traditional implementations of genetic operators, such as n-point crossover and mutation. Crossover was implemented as mutual exchange of genes between selected parents and mutation was implemented as a replacement of gene with a uniform random number from the interval $[0, 1]$.

4.2 Fitness Function

Fitness function is used to rank candidate solutions among themselves. There is a number of methods to measure randomness and evaluate pseudorandom sequences. The DIEHARD battery of tests^[2] and NIST test battery^[3] are de-facto standard tools to detect non-randomness in pseudorandom sequences. They can be used to find sound statistical evidence that a pseudorandom sequence feature non-random patterns. Many randomness tests in both mentioned test suites are based on χ^2 test defined by^[1]:

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i} \quad (3)$$

where O_i is the observed count of occurrences of i -th number and E_i is expected number of occurrences of i -th number.

In this work, we have used the χ^2 statistics of the frequency of numbers produced by the generator described in def. ^[3] as the goodness of fit measure of the generated pseudorandom sequence. The χ^2 statistics was computed after 10000 numbers were produced by the generator. We expect the same uniform probability of every symbol $E_i = \frac{n}{no.of_symbols}$ in the pseudorandom sequence.

The fitness function f was defined as:

$$f = \frac{1}{\chi^2} \quad (4)$$

so that better (i.e. producing more uniform sequence of symbols) pseudorandom generator was characterised by higher fitness value.

5 Experimental Optimization

This section summarizes experimental genetic search for quasigroups isotopic to the quasigroups of modular subtraction with the dimensions 128, 512 and 2048 respectively. We have implemented genetic algorithm with permutation encoding and fitness function as defined above. The parameters of the algorithm (probability of mutation, probability of crossover etc.) were selected after initial tuning of the algorithm. The parameters are summarized in Table ⁽¹⁾.

Table 1. The settings of genetic algorithm for quasigroup search

Parameter	value
Population size	20
Probability of mutation (P_M)	0.02
Probability of recombination (P_C)	0.8
Selection operator	elitist
Max number of generations	1000

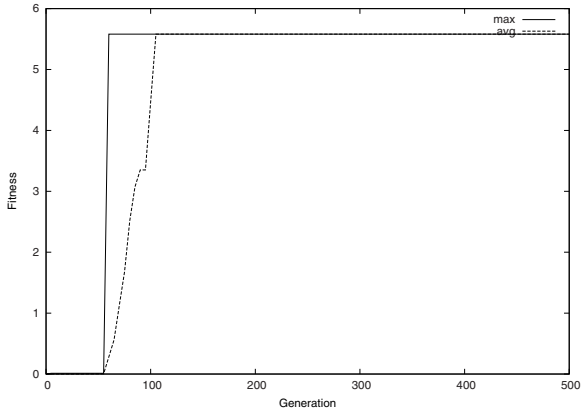


Fig. 2. Maximum and average fitness during an optimization of the quasigroup of dimension 128

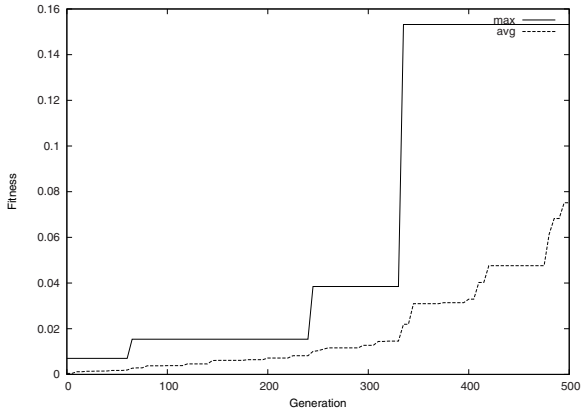


Fig. 3. Maximum and average fitness during an optimization of the quasigroup of dimension 512

The progress of the genetic optimization of analytic quasigroups of the dimensions 128, 512 and 2048 is shown in Fig. 2, Fig. 3 and Fig. 4 respectively.

The evolutionary search has found better quasigroups (in terms of used fitness function) in all experiments. For the quasigroup of length 128, it improved fitness from average 0.00011 in the randomly generated initial population to 5.58. For quasigroup of the length 512, the average initial fitness of 0.00044 was im-

² <http://www.stat.fsu.edu/pub/diehard/>

³ csrc.nist.gov/rng/

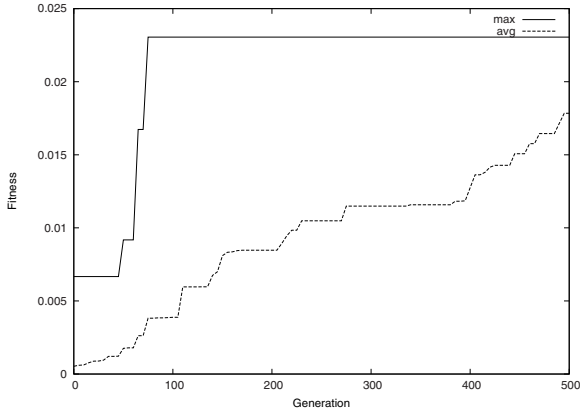


Fig. 4. Maximum and average fitness during an optimization of the quasigroup of dimension 2048

proved to 0.153 and for quasigroup of the length 2048 improved the optimization procedure fitness from 0.00053 to 0.023. In all cases was the final improvement achieved before generation 500.

The results demonstrate that genetic algorithms can improve randomly selected quasigroups (isotopic to the quasigroup of modular subtraction) towards properties defined by used fitness function. The resulting quasigroups are better PRSG generators in terms of normality of symbol frequency in the generated pseudo random sequence. Therefore, the optimized quasigroups are more suitable for applications in computer security.

6 Conclusions

In this paper was described a genetic algorithm for optimization of an analytic quasigroup. The genetic algorithm looks for good bit permutations that are used to construct analytic quasigroups with desired properties. Both, the analytic quasigroup and bit permutation, do not rely on the lookup table of the quasigroup stored in memory. Therefore, large quasigroups can be used and optimized efficiently.

A χ^2 based goodness of fit statistics evaluating the randomness of pseudo-random sequence produced by the quasigroup operation was used as a basis of the fitness function. When used as a pseudorandom number generator, the optimized quasigroups generate better pseudorandom sequences than randomly chosen quasigroups of the same length.

Conducted numerical experiments suggest that the genetic algorithm is a good tool to find optimized quasigroups. In our future work, we aim to find more appropriate fitness function that will allow us to find quasigroups suitable for cryptography and also other application areas.

Acknowledgment

This work was supported by the Czech Science Foundation under the grant no. 102/09/1494.

References

1. Knuth, D.E.: The art of computer programming. In: *Seminumerical Algorithms*, 3rd edn., vol. 2, Addison-Wesley/Longman Publishing Co., Inc. (1997)
2. Marsaglia, G., Tsang, W.W.: Some Difficult-to-pass Tests of Randomness. *Journal of Statistical Software* 7(i03)
3. Markovski, S.: Quasigroup String Processing and Applications in Cryptography. In: *Proceedings 1st Conference of Mathematics and Informatics for Industry*, Thessaloniki, Greece, pp. 278–290 (2003)
4. Dimitrova, V., Markovski, J.: On quasigroup pseudo random sequence generator. In: *Manolopoulos, Y., Spirakis, P. (eds.) Proc. of the 1-st Balkan Conference in Informatics*, Thessaloniki, November 2004, pp. 393–401 (2004)
5. Bäck, T., Hammel, U., Schwefel, H.-P.: Evolutionary computation: comments on the history and current state. *IEEE Transactions on Evolutionary Computation* 1, 3–17 (1997)
6. Belousov, V.D.: *Osnovi teorii kvazigrup i lup*, Nauka, Moscow (1967) (in Russian)
7. Bodenhofer, U.: *Genetic Algorithms: Theory and Applications*. Lecture Notes, Fuzzy Logic Laboratorium Linz-Hagenberg (Winter 2003/2004)
8. Dianati, M., Song, I., Treiber, M.: An introduction to genetic algorithms and evolution strategies, technical report, University of Waterloo, Ontario, N2L 3G1, Canada (July 2002)
9. Dénes, J., Keedwell, A.: *Latin Squares and their Applications*. In: *Akadémiai Kiadó*, Budapest. Academic Press, New York (1974)
10. Dénes, J., Keedwell, A.: A new authentication scheme based on Latin squares. *Discrete Mathematics* (106/107), 157–161 (1992)
11. Dvorský, J., Ochodková, E., Snášel, V.: Hash Functions Based on Large Quasigroups. In: *Proceedings of Velikonoční kryptologie*, Brno, pp. 1–8 (2002)
12. Hilewitz, Y., Shi, Z.J., Lee, R.B.: Comparing fast implementations of bit permutation instructions. In: *Conference Record of the Thirty-Eighth Asilomar Conference on Signals, Systems and Computers*, Pacific Grove, California, USA, November 2004, pp. 1856–1863 (2004)
13. Gligoroski, D., et al.: EdonR cryptographic hash function. Submission to NIST's SHA-3 hash function competition (2008), <http://csrc.nist.gov/groups/ST/hash/sha-3/index.html>
14. Gligoroski, D., Markovski, S., Kocarev, L., Svein, J.: The Stream Cipher Edon80. In: *Robshaw, M.J.B., Billet, O. (eds.) New Stream Cipher Designs*. LNCS, vol. 4986, pp. 152–169. Springer, Heidelberg (2008)
15. Snášel, V., Abraham, A., Dvorský, J., Krömer, P., Platoš, J.: Hash functions based on large quasigroups. In: *Allen, G., Nabrzycki, J., Seidel, E., van Albada, G.D., Dongarra, J., Sloot, P.M.A. (eds.) ICCS 2009, Part I*. LNCS, vol. 5544, pp. 521–529. Springer, Heidelberg (2009)
16. Jones, G.: Genetic and evolutionary algorithms. In: *von Rague, P. (ed.) Encyclopedia of Computational Chemistry*. John Wiley and Sons, Chichester (1998)

17. Knapskog, S.J.: New cryptographic primitives. In: CISIM 2008: Proceedings of the 2008 7th Computer Information Systems and Industrial Management Applications, pp. 3–7. IEEE Computer Society, Washington (2008)
18. Koza, J.: Genetic programming: A paradigm for genetically breeding populations of computer programs to solve problems. Technical Report STAN-CS-90-1314, Dept. of Computer Science, Stanford University (1990)
19. McKay, B.D., Wanless, I.M.: On the Number of Latin Squares. *Journal Annals of Combinatorics* 9(3), 335–344 (2005)
20. Merkle, R.C.: Secrecy, authentication, and public key systems. Stanford Ph.D. thesis, pp. 13–15 (1979), <http://www.merkle.com/papers/Thesis1979.pdf>
21. Mitchell, M.: *An Introduction to Genetic Algorithms*. MIT Press, Cambridge (1996)
22. Ochodková, E., Snášel, V.: Using Quasigroups for Secure Encoding of File System. In: Proceedings of the International Scientific NATO PfP/PWP Conference Security and Information Protection 2001, Brno, Czech Republic, May 9–11, pp. 175–181 (2001)
23. Smith, J.D.H.: *An introduction to quasigroups and their representations*. Chapman & Hall/CRC (2007)
24. Snyder, L.V., Daskin, M.S.: A random-key genetic algorithm for the generalized traveling salesman problem. *European Journal of Operational Research* 174(1), 38–53 (2006)
25. Vojvoda, M.: Cryptanalysis of One Hash Function Based on quasigroup. In: Conference Mikulášská kryptobesídka, Praha, pp. 23–28 (2003)
26. TREC Web Corpus: GOV (2009), http://ir.dcs.gla.ac.uk/test_collections/govinfo.html
27. Snášel, V., Abraham, A., Dvorský, J., Ochodková, E., Platoš, J., Krömer, P.: Searching for Quasigroups for Hash Functions with Genetic Algorithms. In: Proceedings of the 2009 World Congress on Nature & Biologically Inspired Computing, pp. 367–372. IEEE Computer Society, Los Alamitos (2009)

Software Openness: Evaluating Parameters of Parametric Modeling Tools to Support Creativity and Multidisciplinary Design Integration

Flora Dilys Salim and Jane Burry

Spatial Information Architecture Laboratory, RMIT University, Melbourne, Australia
{Flora.Salim, Jane.Burry}@rmit.edu.au

Abstract. The ubiquitous computing era has pushed the Architecture, Engineering, and Construction (AEC) industry towards new frontiers of digitally enabled practice. Are these the frontiers originally identified by the pioneers in the field? Architectural design has progressively shifted from two-dimensional paper based pencil sketched models to digital models drawn in various Computer-Aided Design (CAD) tools. The recent adoption of parametric modeling tools from the aerospace industry has been driven by the need for tools that can assist in rapid flexible modeling. The adaptation of parametric modeling has reformed both pedagogy and practice of architectural design. The question remains if parametric design has answered all the requirements specified by Steven Anson Coons in his 1963 proposal for a Computer-Aided Design (CAD) system. Given the growth of computational power and ubiquitous computing, how has CAD met the visions of its pioneers with respect to the flexibility and ease of communication with the computer and support of simultaneous design conversations with many designers working on the same project? This paper will revisit ideas conceived by the early inventors of CAD, explore the opportunities for advancing parametric modeling with the existing ubiquitous computing infrastructure, and introduces the notion of software openness to support creativity and multidisciplinary design integration.

Keywords: CAD, CAE, Parametric modeling, Parametric design, Software Openness.

1 Introduction

The rise of Computer-Aided Design (CAD) and Computer-Aided Manufacturing (CAM) has led to widespread adoption of digital modeling in the field of manufacturing and construction. There is an increasing demand for flexible modeling tools that allow rapid adaptations to design variations to be produced in real time. CAD is a term originally introduced by Steven Coons [1] and Ivan Sutherland [2] as an integrated design approach which leverages both the creativity and imagination of man and the analytical and computational power of the computer. CAD eventually evolved to tools that truly assist designers. Parametric design software, which was hitherto used mainly in the aerospace and automotive industry, had been adapted to

the Architecture, Engineering and Construction (AEC) industry [3], [4]. We, a computer scientist (the first author) and a practicing architect in academia (the second author), are involved in various integrated multidisciplinary architecture-engineering projects which require explorations of existing design tools and techniques.

This paper aims to investigate the current state-of-the-art of digital modeling and reevaluate parametric design software retrospectively in the light of the aims of early CAD inventors. The review also seeks to discover the requirements for a truly flexible and open parametric design tool that can leverage creative design thinking and multidisciplinary design integration.

The paper is organized as follows. Section 2 presents the summary of the requirements for CAD as proposed by the originators. Section 3 reviews the existing parametric CAD software in the light of the requirements stated in the previous section and outlines the research gap that needs to be investigated. Section 4 presents our proposed notion of *software openness* given the need to leverage creative design thinking and multidisciplinary design integration. Section 5 concludes the paper.

2 Reflections on Requirements for CAD

According to Steve Coons, a CAD system needs to support the design process by assisting creative thinking, increasing productivity, and improving intercommunication and collaboration through a well-defined human-computer interaction, or “*design conversation*” [1]. There are four key requirements of CAD system as stated in the paper.

Firstly, it must have the ability to accept, infer, and store shape description introduced graphically using graphical means as well as structural abstractions (with symbolic languages that is capable of representing interconnected ideas within the design) [1].

Secondly, coupled with the visualization, it must also be an analysis tool that is able to perform all the computations required for the design process. This includes structural analyses, mechanical services, electrical analyses, and other analytical processes that are required to validate and optimize the design [1].

Thirdly, it should be a platform that promotes ease of collaboration by enabling the same model to be accessed and modified by multiple designers in different locations simultaneously (Coons, 1963). The notion of the internet and online collaboration was not even close to reality in 1963. ARPANET (Advanced Research Projects Agency Network), the forerunner of the Internet, was only introduced in 1968 [5]. The idea described in the paper [1] depicts a situation that was only realizable in the early 90s when the Internet had reached one-million users worldwide [5].

Finally, a CAD system must be highly generic to accommodate the creative activities that reside in a trans-disciplinary design domain. The general problems of architects and engineers need a system that is flexible and adaptable to multiple specific domains [1].

Apart from Steve Coon’s manifesto, Kasik et al. has reviewed the progress of CAD departing from Ivan Sutherland’s Sketchpad application and presented ten challenges of CAD and classified them into three different categories: “*computational geometry, interactive techniques, and scale*” [6]. A number of those challenges have been

answered in recent developments of parametric modeling tools. There are profusion of CAD software that perform solid modeling but are not parametric. The scope of this review only covers CAD packages that perform flexible 3D parametric modeling and are referred as parametric modeling tools from this point forward.

3 Parametric Modeling, Are We There Yet?

Parametric modeling was first introduced as a means for design reuse [7]. Parametric modeling uses parameterized relationships between components in the design to define forms [8]. A parametric model comprises variable attributes, which are called parameters and fixed attributes, which are called constraints [4]. A parametric model responds to changes of parameter value or definition without erasure of parts of the model or starting the design model from scratch. Designers alter the value of the parameters to explore design variations that can be generated from the same model [4].

Parametric design tools, such as GenerativeComponents by Bentley Systems, Digital Project by Gehry Technologies (an adaptation from the powerful and expensive CATIA software used by the aerospace industry for the AEC industry), ArchiCAD by Graphisoft, and the Revit family of products by Autodesk, have been taken up more widely by the AEC design and design education communities worldwide only during this decade. In the previous decade, its use was known but through more isolated and research-led projects in practice and academia. AutoCAD was the main software utilized in industry but it was not parametric. The rise of the usage of parametric modeling tools is due to its power to manipulate geometry and generate variations by simply changing its parameters. Geometry manipulation is not limited to simple extrusion or Boolean operations. Parametric modeling tools support the construction of a single model that has many varying geometrical instances. Initially, parameterized mathematical associations between objects in the model need to be defined. For some parameters, it may be possible to determine a range, or the maximum and minimum values based on design criteria. Then by iterative refinement of both the mathematical associations and the values of the parameters in the model, various design options are explored. These design options can be compared and analyzed further and criteria determined for selection of the better and best solutions.

The Sagrada Família church in Barcelona (Fig. 1 – Left), which was started in 1883 and left unfinished by the architect Antoni Gaudí who died in 1926, is an example of a building and a design system in which there is a clearly defined parametric geometry that provides the key to the design process. With the help of parametric ‘recipe’ left by Gaudi in his plaster models and drawings, the ongoing design and construction of the Sagrada Família relies on Gaudí’s technique of prescribing variations of complex geometries based on simple forms exploiting the variation of geometric parameters of those forms. One example of the subsequent influential work to interpret and model Gaudí’s intentions digitally and parametrically is the columnnets model, in which variation of the surface parameters, rotation and Boolean operations are applied on trimmed hyperbolic paraboloids to generate the forms [9]. Two examples of the more recent projects using parametric design tools for design exploration are the Dubai Towers (Fig. 1 – Right), designed by TVS using Generative Components by Bentley Systems, and the “Bird’s Nest” Beijing National Stadium (Fig. 2), designed by Herzog and de Meuron using Digital Project by Gehry Technologies.



Fig. 1. Left: The clerestory window in Sagrada Familia, Image courtesy of: Mark Burry. Right: Dubai Towers. The Lagoons Image courtesy of: Thompson, Ventulett, Stainback & Associates (www.tvsa.com).

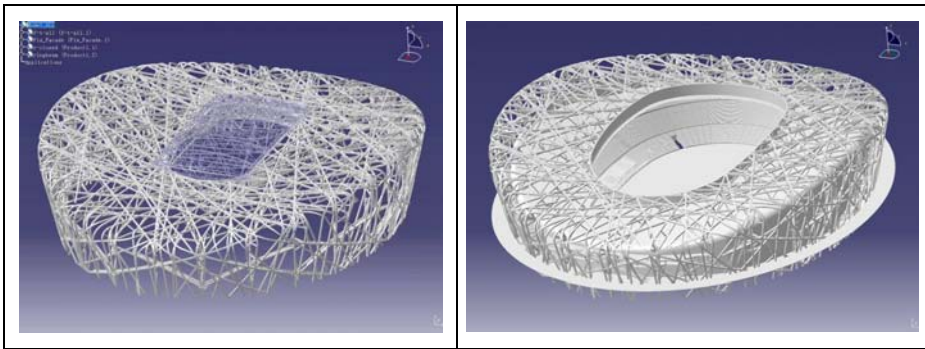


Fig. 2. Left: An early digital model of the Bird's Nest Beijing National Stadium as modeled inside Digital Project. Right: the final model. Image courtesy of: Gehry Technologies LLP (www.gehrytechnologies.com).

We perceive two different groupings of parametric design tools given the varying approach to parametric design. The first group of parametric design tools is based on *associative-geometry*, where parameterized mathematical descriptions and associations between points, curves, surfaces, and solids are possible. The tools belonging to this group include Bentley GenerativeComponents and Rhino Grasshopper. The second group of parametric design tools are focused on *BIM*, where parametric relationships encapsulate parametric descriptions of components of a building design across multiple disciplines [10]. BIM stands for Building Information Modeling, which is regarded as “a data-rich, object-oriented, intelligent and parametric digital representation of the facility” by the American General Contractors (as quoted in [11]). The Autodesk Revit 2010 package (Revit Architecture 2010, Revit MEP 2010, Revit Structure 2010) belongs to this group as well as Gehry Technologies’ Digital Project. Revit 2010, however, is unable to handle NURBS geometry modeling [12] and

therefore Revit is not generally used for modeling complex curved geometries. Digital Project is known to be a very powerful parametric CAD package that handles both complex parametric geometric associations as well as parametric design representations referring to diverse libraries across multiple disciplines. Nevertheless, it is also known to be the most expensive CAD package that is currently around and customizing model via the script interface is not trivial. The first group of tools can also be used for parametric representations that link different disciplines, however, BIM high-level object definitions are not included in the tools by default.

In a reflective mode, we review the existing parametric modeling tools and revisit the CAD requirements by the CAD visionary and retrospectively raise the question – are we there yet? The following evaluation is raised in the light of the use of the current parametric modeling tools and processes.

3.1 Are They Capable of Accepting, Inferring, and Storing Geometries Using Graphical and Symbolical Means?

This section evaluates the existing parametric modeling tools' capabilities to process and manage the geometries in the design using graphical and symbolical means. The abstractions of the design in parametric modeling tools are often represented in a symbolic diagram. These features are now intrinsic to all the existing parametric design tools. GenerativeComponents (GC) encompasses a 3D interactive view, a symbolic view, and an object view [13]. One distinguishing feature of GC is the way it displays both the textual and symbolic representation of the parametric associations between components in the model. This feature is also available in Grasshopper, the new extension of Rhino. The symbolical representation of a parametric model is analogous to an object diagram of a software application. In fact, the notion of components, parameters, constraints, and associations in a parametric model resemble those of the object-oriented programming paradigm. The graphical and symbolical representations undoubtedly assist designers to rethink their design in the light of the components and their associations in the model as well as the parametric variations they would like to generate from the model. Both the graphical and symbolical representations visualize design intent that is captured within the tool, particularly during the early design stage when design options need to be generated rapidly.

Revit 2010 uses the more rigid tree structure to represent the higher level of abstraction of the model. Geometrical associations between components in the model are not visualized in the tree structure. Since each component in a Revit model is an instance of a Revit family type (e.g. an instance of a 36' flush panel type of the door family), the tree structure provides the abstraction and the detailed view of the instances of each type and family that are created in the model. Symbolic representations are used to view the 3D model as engineering plans. However, there is no existing symbolical abstraction to represent the parametric associations of all the parts in a Revit model. Considering that Revit is aimed to be a BIM software, the tree structure is useful for viewing abstract representations of parts of the model to give a non discipline specific overview.

CATIA/Digital Project (DP) primarily uses a hierarchy of linked files. No single file need be very large. It also uses simple graphical files that just show the geometry for representing large extents of the model where it might otherwise take hours or

days to open the Part and Product files. There is a tree structure in CATIA/DP that is an organizational and navigational tool and another that directly maps the model's parametric relationships.

3.2 Are They Capable of Performing Analysis in Different Performance Areas?

To date, CAD tools (i.e. parametric design tools) and analysis tools have been distinct and autonomous. Cross-disciplinary analysis is a challenge since current CAD systems are not integrated well with CAE (Computer Aided Engineering) analysis tools [Kasik05]. Some tools, such as Revit Architecture 2010 has released a plug-in/extension for solar incidence and energy analysis to be performed within the software. Another example is Google Sketchup, which has an extended capability for a simple energy analysis when the IES Virtual Environment plug-in is installed with it. However, there is no existing tool on the market that currently carries the capability to perform all-in-one design and analysis, where the analysis can respond to the design variations made in the same model. Some recent work has come close to realizing this silver bullet [14].

In order for a parametrically-modeled design proposal to be assessed for structural loads, energy efficiency, solar incidence analysis, lighting analysis, and other types of analyses, it must undergo "expert translation" since geometrical representation in various domains are interpreted differently [15]. Translation errors often arise from two key sources, which are floating-point arithmetic and tolerances [6]. Approximations used in numerical calculations reducing floating-point that substitute higher performance for high precision cause errors in models passed to disciplines that require high precision analysis [6]. Tolerances are used to control the level of approximation and need to be managed to produce geometry that can meet different discipline analysis requirements. If the tolerances are set high, analysis can be performed more quickly but with greater likelihood of error in the resulting geometry [6]. If the tolerance is too small, analysis may fail to converge at all [6]. Therefore, selecting a tolerance to be used requires good comprehension of the engineering analyses to be used, the environment of the analyses, and the specific software to be employed [6].

3.3 Are the Models Generic and Adaptable to Various Domains Pertaining to the Design Process?

This section deals with Coons' requirement for a CAD system that can support cross-disciplinary model interoperability. According to a survey by McGraw-Hill construction study [16], 3% of a project cost corresponds to software non-interoperability. In the U.S. construction market worth \$1.3 trillion in 2007, this corresponds to \$36 billion of productivity waste. In the global market worth \$4.6 billion, the figure escalates to \$138 billion [16]. In order to maintain building information efficiently and the vary representations in a broadly reusable way, the building industry has recently attempted to define and encourage the adoption of Building Information Modeling (BIM). With too many industry players and major CAD vendors introducing a number of commercialized BIM tools, however, many have found themselves reinventing the wheel in order to bridge tools created by one vendor to another tool by another vendor.

Industry Foundation Classes (IFC), the accepted industry standard for interoperability and model exchange, are too heavy and complex (IFC consist of 327 data types, 653 entity definitions, and 317 property sets [17]). This poses challenges to performance, scalability, and ease of use and adaptation into specific sub-domains. Green Building XML (gbXML), which was introduced as a lighter option to IFC, only encompassing a subset of IFCs has become a de facto standard for interoperability of green building design. However, semantic translations from those standards to BIM tools and vice versa (import to/export from BIM tools) still generate translation errors. Lee et al argue that parametric modeling needs to be central to BIM and provide mechanisms for translation to various domains [7].

3.4 Do They Promote Collaboration and Communication among the Designers or Users?

Over the top of any provision for collaboration and model sharing and merging within existing parametric design software packages, there is a need for additional communication channels that combine graphical, symbol, and language based information. Currently, separate tools are required Product Lifecycle management tools for manufacturing aim to fulfill this role by integrating information and communication with customers/clients, suppliers, enterprise/planning and systems development information and communication.

However, in practice it is a lightweight tool that allows rapid human interaction as well as capturing the history of informal decision making and rapid change in the design model that is required in combination with parametric modeling. Autodesk Design Review and Octopz are examples of visual annotation environments that support this level of interaction. Versioning software has been integrated with Digital Project parametric software as the means of team model sharing online, allowing visual indicators of ownership of different parts of the model, locked in the versioning software, at any given time within the modeling environment. The concept is familiar to Archicad Teamwork and Bentley Projectwise users.

There are two fundamental truths that Jerry Laiserin has stated of the current state of technology in the industry: that while collaboration across design teams is the most critical factor for successful building information modeling, as model data integration goes up, flexibility of workflow and performance in collaboration go down [18]. This is true even for the most straight forward, explicit geometrical representation of buildings and currently leaves models shared between large teams of a fundamentally more formative and flexible nature largely in the aspirational realm even 40 years on.

3.5 Do They Support the Design Process by Communicating Design Intent?

One of Coons' manifestos is for CAD to promote design conversations between the computer and the designer. In parametric design tools such as Digital Project and GenerativeComponents (GC), the history of geometrical associations and parametric value variations is recorded. In GC, parametric changes are recorded as transactions and users can go back and forth in the transaction list to undo or redo changes in the model. Whenever the value of a parameter or the parametric relationship between one component to another is changed, the value and relationship changes are recorded in

the transaction file. This is one way for CAD tools to carry “design conversations”, since designers constantly perform iterative loops in evaluating their design and mental design conversations that are channeled through their modeling activities can be recorded.

However, difficulties in communicating design intent for multidisciplinary model sharing still persist. When standards of model interoperability exchange such as IFC or gbXML are used for model sharing, parametric representations and associations are not captured either by IFC or gbXML during translation, since the dynamic associations between the components in the model are replaced with static values representing the components themselves. The design intent that needs to be maintained in a parametric model is lost in translation and unrecoverable once the model is imported or exported to other tools or analysis programs.

3.6 Do They Increase Productivity in the Design Process?

Parametric design tools undoubtedly increase productivity. The efficiency of using parametric design approach in modeling depends on the number of the iterations to the design and the extent of the variations in the model [19]. In this respect, a large part of Coons’ future gazing has come to pass. However, there are many measures of productivity. Certainly the speeds of documentation and repetitive reuse of data is massively increased through moving from manual drafting to CAD. More creative activities, combining analytical and synthetic approaches to building shape, structure and articulation have benefited in the ways that were both predicted and prototyped in the 1960s. The size and complexity of building projects has also increased hugely since the 1960s and we can speculate that the technological means have fed this trend. However, the continued use of the term ‘documentation’ hints at the continued lag in the transition of human systems to the new means of production, arguably slowed by the introduction of personal computing and its disaggregating impact on teams and projects [20]. With the growth of computing power and the internet in the last decade plus, the integration of information and reduction in costly design document error and inconsistency has been a major focus in increasing the contribution of IT to productivity measured by building cost.

The usability of parametric design tools has increased with the expansion of computational power. In general, the more complex the software is, the lower the usability. However, in architecture and engineering disciplines, the usability of the software also depends on its capacity to host a model that contains multidisciplinary design representation for multidimensional model editing and real-time visualization. CATIA is capable of performing complex curves manipulation and visualization of the whole airplane model as well as the dependent parts of the model.

3.7 Do They Support the Design Process by Assisting Creative Thinking?

Lars Hesselgren stated “*the advantage of using GenerativeComponents is that it helps me think about what I am doing. The disadvantage is that it forces me so to think*” [13]. Parametric modeling tools have forced designers to quantify their design to components and associations for generating variations that they would like to explore further. However, when a parametric design grows and the number of potential

parametric variations grows exponentially, the flexible model poses the threat of inflexibility given that the proliferation of backward/forward chaining of changes can be unmanageable [21]. This is where customization, scripting, or coding comes in to pull the weight off from visual design. Designers are now becoming more interested in increasing the efficiency of model building through scripting or coding. To designers who code, coding activity is considered as a channel for creativity and means of representing design ideas that goes beyond visual or graphical drawing. Current tools have indeed assisted creative thinking. However, to innovate, designers need to extend the exploration of their design ideas using complementary creative methods to the use of digital/visual synthetic geometric modeling aids.

Lawson [22] stated his desire to see CAD becoming an agent for creativity. As an agent, CAD would act as a subordinate who learns and understands the design intention of the master and is able to assist creative thinking by crawling the internet and finding design ideas that can match the master's interest [22]. Lawson's desire can potentially be realized if parametric modeling is integrated with such computational techniques as machine learning or data mining. The existing work involving machine learning techniques for the design process includes evolutionary design [23], [24] and agent based modeling [25].

3.8 Are We There Yet?

Although it has been a long road to general adoption (and we are still clearly on the uptake curve), the current parametric design tools have to some degree fulfilled the requirements of CAD as prescribed by Coons. Parametric design tools have improved productivity through rapid generation of design variations over simple parametric changes and helped creative thinking in the process of designing. However, there are still gaps and opportunities for advancing the existing parametric design tools.

Firstly, multidisciplinary design model integration is still very challenging [18]. Given the complexity of the process, analysis, and tools required in various sub-domains, an all-in-one engineering analysis tools in a CAD package as proposed by Coons is not feasible in the AEC industry. Cross-CAD and cross-disciplinary model sharing is still a challenge due to software interoperability issues and translation errors. There is a need to escalate the degree of *software openness* in order to carry design intent between designers of one discipline to another. *Software openness* is discussed further in the next section.

Secondly, the need for an online multidimensional design review tool for design collaboration as an additional layer on top of the existing parametric design software is apparent. In the current cyber digital and ubiquitous computing era, a multidisciplinary team working on the same project might be dispersed in different locations across the world. Having to share design variations on the same model requires online collaboration tool and parametric design tool to be integrated.

With the flexibility of parametric modeling, the creative thinking process of the designer can now inform scripts to drive the process of parametric changes. Given the notion of openness of parametric modeling, parameters can be scripted and their values driven from any source, the intent of the designer, the results of performance analysis in the cyber world, or information gathered through sensors, for example, in the physical space of our world. It opens endless possibilities for supporting both creativity and productivity.

The initial need for parametric design tools is the requirement to support flexible modeling. Flexibility should be an inherent attribute of any parametric design tool. In the light of the converging physical and digital spaces, architectural form inside the CAD should keep on transforming and producing emergent results which are reflective on the physical environmental changes as well as the virtual energy simulation results. The current state-of-the-art technology in the cyber/digital world, such as sensing technology, ubiquitous computing, and cloud computing, can potentially be combined with the power of parametric modeling to enhance the flexibility, fluidity, and creativity of designing. Therefore, there is an opportunity to leverage the existing parametric modeling tools by keeping it open to the user, the environment, and other existing tools in the space. In addition, given the complexity of the process, analysis, and tools required in various disciplines, there is a need to escalate the degree of *software openness* in order to carry trans-disciplinary design conversations. This is discussed further in the next section.

4 Software Openness of Parametric Modeling Tools

Hansen [26] quoted Sanford Kwinter that “*the world where everything flows seamlessly together in real time*”, added an argument for the requirement of “*soft systems*” which are driven by their softness and openness to the environment and adaptable to information exchanged with its surroundings. Concurring with Hansen, we propose the notion of *software openness*, which is a state of unimpeded ability of 3D parametric software to be openly customized in order to:

1. exchange information with other software (*open interoperability*) and with the physical environment (*physical-virtual information integration*);
2. modify the exchanged information and the information handling process;
3. mediate, compare, and consolidate the generated design variations for analysis, simulation, and fabrication.

Software openness is required to support design integration across disciplines and creative design uptake of the ubiquitous computing opportunities as the physical realm and virtual worlds become more closely intertwined. This differs from open source software. Software openness does not imply source code access, instead, it requires a methodological approach to defining ways for open input, open process, and open output to parametric modeling tools.

The fundamental of information handling in any software comprises input, process, and output (IPO). Conventionally, it is a linear progression without any feedback loop (Fig. 3). Software openness requires open input, open process, and open output (Fig 4). The input, process, and output parameters of parametric modeling tools need to be opened for multi-devices, multi-tasking, multi-disciplinary and multidimensional interaction effectively to harness creativity, productivity, and collaboration. The information handling in such software should no longer be a linear progression, but a closed loop of IPO that allows real time and open interaction or feedback.

The degree of software openness can be used to assess the suitability of a design tool in the process of choosing the right sets of tools for a multidisciplinary design project. The higher the degree of a tool being receptive to open input, open process,



Fig. 3. Input Process Output

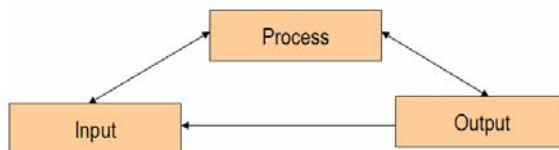


Fig. 4. Open Input, Open Process, Open Output

and open output, the more usable the tool is in the context of a multidisciplinary design project. Since a parametric design tool do not generally include the full range of engineering analysis tools, such as building thermal performance analysis and structural analysis tools, the higher the degree of the software openness, the higher is the possibility of the software being used in design model integration across disciplines and design activities. Since performance analysis and software interoperability are still challenges in multidisciplinary design model integration, the ability to customize design representation, translation, and input/output from one software to another enables engineers to customize the software for the purpose of the project. The need to assess the software openness for interoperability is demonstrated in the PM4D project [27] by Centre for Integrated Facility Engineering (CIFE), Stanford University, where the input and output from one design, engineering analysis, or visualization software to another are compared for the best selection of tools required for the integrated project delivery.

Open Input. Input to a parametric tool can either be models (imported from another CAD) or textual files containing analysis requirements (which could be a result imported from external analysis tools). In a multidisciplinary design project, it is necessary to analyze the input capabilities of the tools being used and the existing model translation mechanisms between one software to another. The interoperability gap between software, where it exists, could then be tackled in two ways. Firstly, the ‘tool-by-tool’ approach which employs a specific tool or plug-in for export and import of a design model between one tool and another, such as demonstrated by Nicholas and Burry [15] in translating a model created in Rhino CAD software to a model understandable by Radiance, a lighting analysis software. Secondly, the ‘bridge tool’ approach which employs an intermediate tool for model translation among various CAD, analysis, and visualization software. This approach was exemplified in the development of DesignLink Software Development Kit (SDK) [28], [29], by Arup and Spatial Information Architecture Laboratory (SIAL), which relies on the translation of various models into DesignLink XML, an XML schema designed for model information exchange in the early design stage.

In order to promote design conversations between human designers and the computer, various sensor and haptic devices need to be explored as an input device to parametric modeling. Handheld devices help designers to converse with the parametric

models better. For example, a tangible view cube was developed to interface with Autodesk tools in order to represent the 3D digital model in a physical environment and promote tangible interaction between the user and the CAD [30]. The ClayTools system from Senseable [31] is a free-form robot hand that allows model creation, crafting, and detailing straight to the virtual model with a force feedback applied to the user. The ClayTools can become input to 3DMax, Maya, and Rhino, which are not parametric modeling software. However, with Grasshopper, Rhino's parametric plug-in, there is a potential in connecting Rhino-Grasshopper virtual parametric models with ClayTools. Another example of a haptic device that can be used for exploring, designing, and generating inputs for parametric models is the Nintendo Wii remote (or Wiimote). A Wiimote plug-in for Autodesk Design Review by Autodesk Labs was developed for designers to navigate through the 3D model and environment using the Wiimote [32]. Another existing work on development of haptic input devices to parametric modeling software is the Wiimote features development for Generative Components (GC) [33], which allows the Wiimote to interact with the GC to generate forms and variations based on real-time data input. The Wiimote can be used to draw in 3D space, modify existing 3D models, or used as a controller of the model (Fig. 5).

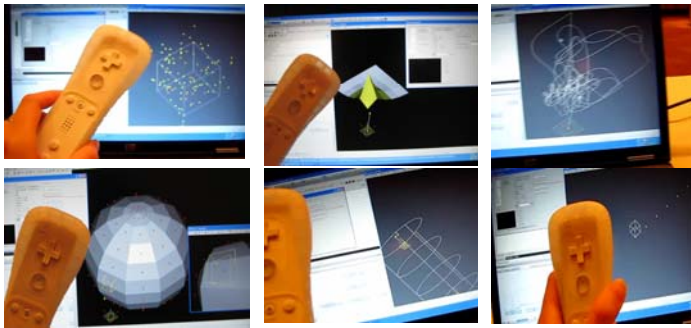


Fig. 5. A Wiimote interacts with GenerativeComponents

Open Process. Software openness can also be used to assess the suitability of the software to be openly customized in order to generate a creative output. Creativity can be achieved not only through physical or digital modeling, but also through creative coding [34]. Parametric design and analysis tools need to accommodate for extensions, customization, and open interoperability, in order that additional plug-ins, class libraries, software components, and scripts can be developed to support creativity and productivity. Interoperability checking and automation through code can be developed when the communication channel between software is open. In Digital Project, customization through codes is available from their internal scripting window, which uses VBA script or by expert use of the CATIA SDK. Given the nature of VBA as procedural language, component and object reuse in the code is a challenge since classes cannot be defined inside the VBA script. The open component based software architecture adopted by GC and Grasshopper (Rhino) provides a better flexibility for customization in comparison to non-parametric modeling tools which largely rely on internal scripting techniques. GC and Grasshopper (Rhino) allow installation of

custom plug-ins or Dynamic Link Libraries (DLLs) to define new features, objects, or parametric relationships or to extend/inherit from the existing type libraries. Autodesk offers an Application Programming Interface (API) to access the classes, methods and operations provided by Revit programmatically on .NET platform. Commercial and non-commercial parametric modeling software could open the process of definition, generation, modification, and deletion of design objects and parametric relationships among the objects in the model to promote collaboration, interoperability, design integration, creativity and interactivity with other software, sensors and haptic devices.

Open Output. Coons' paper opens with reference to the first automatically controlled milling machine at MIT in the 1950s [1]. The CAD/CAM axis so far seems to have proved if anything more vital in the design process than immersive virtual reality. Designing and building in the physical world entails, as it always has, working from and capturing a physical context, representing this virtually, whether through measurement, drawing and physical modeling or through digital capture and representation. Similarly, collective design decision making is still frequently better supported through interaction around physical models. While the cost and speed of traditional physical modeling is in some cases uncompetitive with virtual equivalents during the design process, the opportunity for rapid prototyping allows for the design conversation to flow back and forth between the virtual and physical embodiment.

5 Conclusion and Recommendations

This paper has reviewed the existing CAD software and the current state-of-the-art of parametric modeling software. Although parametric modeling has improved productivity and encouraged creative thinking, the current parametric design tools and processes have not fully answered the requirements of the CAD visionary nor coped with the challenges of the digital age.

The first challenge is to answer the multidisciplinary design model integration, particularly between parametric design and building performance analysis tools. Further, parametric design software are becoming more prevalent, but are still difficult to integrate with analysis tools for informed feedback. In support of the collective use of open parametric design software, there is also a need for multidimensional online collaboration and communication capabilities to be integrated within parametric design tools.

To give their creativity full rein, designers need to think outside the box of the digital screen and digital world in order to hold design conversations that move more effortlessly between the cerebral, virtual and physical worlds. Coding and open interaction with parametric modeling hold the keys to making the transitions between these environments more seamless.

The proliferation of sensor and physical computing devices in the ubiquitous computing era presents immense potential for designing in the mixed reality, where the space of designing is no longer constrained to the physical space. Therefore, software openness can also be measured in the way a parametric design tool is able to cooperate with physical and ubiquitous computing. Designers who model in the digital space

often require real-time interaction and feedback that can be manifest in the physical world. When a parametric design software unlocks the means for software openness, sensors and haptic devices can be linked with parametric design software to achieve tools that are capable of simulating adaptive and responsive forms of architecture.

Finally, the degree of software openness can be used to measure the flexibility of the software to be openly customized for interoperability, model sharing, and supporting creativity.

Acknowledgment

The authors thank Australian Research Council (ARC) for the funding of the research in the ARC Discovery project: Challenging the inflexibility of the flexible model. We also thank the ARC and Queensland State Government's Project Services for the funding of the research in the ARC Linkage project: Assimilation of architectural and services design in early design modeling. We acknowledge John Frazer, Robin Drogemuller, and Bianca Toth from Queensland University of Technology (QUT) and Mark Burry from RMIT University for their contribution to the research team. The authors also acknowledge Daniel Davis from RMIT and Ruwan Fernando from QUT for their assistance.

References

1. Coons, S.A.: An outline of the requirements for a computer-aided design system. In: AFIPS Spring Joint Computer Conference, pp. 299–304. ACM, New York (1963)
2. Sutherland, I.E.: Sketchpad: a man-machine graphical communication system. In: AFIPS Spring Joint Computer Conference, pp. 329–346. ACM, New York (1963)
3. Frazer, J.: Computing without computers. *Architectural Design* 75(2), 34–43 (2005)
4. Hernandez, C.R.B.: Thinking parametric design: introducing parametric Gaudi. *Design Studies* 27(3), 309–324 (2006)
5. Dubendorf, V.A.: *Wireless Data Technologies*. John Wiley & Sons, West Sussex (2003)
6. Kasik, D.J., Buxton, W., Ferguson, D.R.: Ten CAD Challenges. *IEEE Computer Graphics and Applications* 25(2), 81–92 (2005)
7. Lee, G., Sacks, R., Eastman, C.M.: Specifying parametric building object behavior (BOB) for a building information modeling system. *Automation in Construction* 15(6), 758–776 (2006)
8. Monedero, J.: Parametric design: a review and some experiences. *Automation in Construction* 9(4), 369–377 (2000)
9. Burry, M.: Parametric Design and Sagrada Familia. *Architectural Research Quarterly* 1(4), 70–80 (summer 1996)
10. Drogemuller, R.M., Crawford, J., Egan, S.: Linking early design decisions across multiple disciplines. In: *Proceedings of the Conference on Product and Process Modelling in the Building and Construction Industry (ECPPM 2004)*, Istanbul, Turkey (2004)
11. Aranda-Mena, G., Crawford, J., Chevez, A., Froese, T.: Building information modelling demystified: does it make business sense to adopt BIM? *International Journal of Managing Projects in Business* 2(3), 419–434 (2009)
12. Autodesk Inc.: *Autodesk Conceptual Design Curriculum 2010* (2009), <http://usa.autodesk.com/adsk/servlet/item?siteID=123112&id=14038871>

13. Aish, R., Woodbury, R.: Multi-level Interaction in Parametric Design. In: Butz, A., Fisher, B., Krüger, A., Olivier, P. (eds.) *SG 2005*. LNCS, vol. 3638, pp. 151–162. Springer, Heidelberg (2005)
14. Maher, A., Burry, M.: The Parametric Bridge: Connecting Digital Design Techniques in Architecture and Engineering. In: *The 2003 Annual Conference of the Association for Computer Aided Design in Architecture*, Indianapolis, Indiana, pp. 39–47 (2003)
15. Nicholas, P., Burry, C.M.: Import as: Interpretation and Precision Tools. In: *Proc. of the 12th International Conference on Computer Aided Architectural Design Research in Asia, CAADRIA 2007*, Nanjing, China (2007)
16. McGraw-Hill Construction: Interoperability in the Construction Industry (2007), http://construction.ecnext.com/mcgraw_hill/includes/SMRI.pdf
17. Steel, J., Drogemuller, R.: Model Interoperability in Building Information Modelling. In: *Industrialized Software: KISS workshop @ Code Generation 2009*, Cambridge, UK (2009)
18. Laiserin, J.: Next Gen BIM (2009), <http://www.laiserin.com/features/issue25/feature01.pdf>
19. Hoffmann, C.M., Joan-Arinyo, R.: Parametric Modeling. In: Farin, G., Hoschek, J., Kim, M.-S. (eds.) *Handbook of Computer Aided Geometric Design*, ch. 21. Elsevier B. V, Amsterdam (2002)
20. Aish, R.: Migration from an individual to an enterprise computing model and its implications for AEC Research. In: *Berkeley-Stanford CE&M Workshop: Defining A Research Agenda*, Stanford, California (2000)
21. Burry, J., Burry, M.: The Bonds of Spatial Freedom. In: *Proc. of the 26th eCAADe Conference - Architecture in Computro*, Antwerpen, Belgium, pp. 301–308 (2008)
22. Lawson, B.: Oracles, draughtsmen, and agents: the nature of knowledge and creativity in design and the role of IT. *Automation in Construction* 14(3), 383–391 (2005)
23. Frazer, J.: *An Evolutionary Architecture*. Architectural Association Publications, London (1995)
24. Janssen, P.: A generative evolutionary design method. *Digital Creativity* 17(1), 49–63 (2006)
25. Coates, P.S.: Some experiments using agent modelling at CECA. In: *Generative Arts 2004*, Milan, Italy (2004)
26. Hansen, M.: Wearable Space. *Configurations* 10(2), 321–370 (2002)
27. Fischer, M., Kam, C.: *Product Model & 4D CAD - Final Report*. Stanford University (2002)
28. Khemlani, L.: *SmartGeometry 2009 Conference Day* (2009), <http://www.aecbytes.com/feature/2009/SmartGeometry2009.html> (accessed April 27, 2009)
29. Holzer, D., Tengono, Y., Downing, S.: Developing a Framework for Linking Design Intelligence from Multiple Professions in the AEC Industry. In: Dong, A., Moere, A.V., Gero, J.S. (eds.) *Computer-Aided Architectural Design Futures (CAADFutures) 2007*, Sydney, Australia. Springer, Netherlands (2007)
30. Babu, S.: *Tangible View Cube*. CAD Professor (2009), <http://cadprofessor.in/2009/04/tangible-view-cube>
31. Senseable, <http://www.sensable.com>
32. Sheppard, S.: *A Wii Bit of Fun*. In: *Autodesk Labs: It's Alive in the Lab* (2009), http://labs.blogs.com/its_alive_in_the_lab/2007/12/a-wii-bit-of-fu.html
33. Salim, F.D., Mulder, H., Burry, J.: A System for Form Fostering: Parametric Modeling of Responsive Forms in Mixed Reality. In: *Proc. of the 15th International Conference on Computer Aided Architectural Design Research in Asia, New Frontiers*. Hong Kong, China (2010) (accepted for publication)
34. Maeda, J.: *Creative Code*. Thames & Hudson, London (2004)

Dynamic and Cyclic Response Simulation of Shape Memory Alloy Devices

Yutaka Toi and Jie He

Institute of Industrial Science, University of Tokyo,
Komaba 4-6-1, Megruro-ku, Tokyo 153-8505, Japan
toi@iis.u-tokyo.ac.jp

Abstract. Since most of the existing simulation researches on shape memory alloys are focused on their static behavior, the simulation of their dynamic response behavior is the blank field waiting for researchers to explore. In this paper, both one-dimensional and three-dimensional dynamic simulations, as well as an improved simulation model for shape memory alloys' cyclic response behavior will be discussed. The validity of the dynamic and cyclic response simulation model is illustrated by comparing the calculated results with the experiment data. The damping capacity of shape memory alloys is also demonstrated in both one- and three-dimensional simulations.

Keywords: shape memory alloys, dynamic response, finite element method, superelasticity, shape memory effect.

1 Introduction

Since the first large-scale application of shape memory alloys (SMAs) as couplings in the Grumman F-14 aircraft, the implementations of SMA devices have been spread from aerospace to medical field, smart structures, robots etc.. The noninvasive characteristic and excellent bio-compatibility make SMA devices widely applied in biomedical fields. The controllable deformation feature of SMA devices has attracted many attentions in robot researching field. SMAs have been used as actuators and sensors in several experimental robots, such as robots made by Sugiyama and Hirai [1]. Besides, the adjustable stiffness feature of SMAs have been used in passive control devices of vibration controlling, showing a remarkable capability in controlling vibrations with unknown frequency. Even more, the damping capacity of SMAs has become another active research field recently.

However, the practical applications are still limited in small scale such as applications in medical instruments and aerospace devices. One reason is the high price of SMAs, which has made the applications of SMA devices limited in the fields which are not sensitive to cost. Another reason is a lack of understanding of SMAs' properties and simulation tools, which has limited the flexibility when designing SMA devices. The cost problem can be solved by improving the manufacturing process and mass producing, which are jobs for material scientists and manufacturing engineers.

As for the second problem, developing reliable and accurate models for SMA devices is the most important task for researchers in the field of mechanics.

It is clear that, for most of the potential application fields of SMAs, such as SMA actuators, SMA sensors and smart structures with SMAs, dynamic behavior simulation is even more important than static behavior simulation. Since there have been enough static models developed for a variety of SMAs, this paper is focused on modeling the dynamic behavior of NiTi SMA devices. A SMA dynamic model is developed in this paper, whose feasibility is proved in both one-dimensional scale and three-dimensional scales. The model which is used in this paper is based on the static model developed by Brinson [2] and extended by Toi et al. [3]. The dynamic simulation model is inspired by both Seelecke [4], who performed several one-dimensional SMA dynamic simulations using Achenbach's model [5], and by Toi and Choi [6], who have done several impact response simulations of SMA devices by considering a strain rate effect in the previously formulated multi-axial constitutive equation [7]. The present modeling takes into consideration asymmetric tensile and compressive deformation as well as multi-axial stress state which are not considered in the dynamic analysis by Seelecke [4]. However, the strain rate effect, which is included in the constitutive modeling by Toi and Choi [6], is neglected in the present modeling, as the present study discusses relatively slow, vibration problems.

Afterwards, an improvement of SMA dynamic model considering cyclic effect [8-11] is developed and discussed. The cyclic effect on constitutive behaviors of SMAs in cyclic loading as well as cyclic deformation conditions, which is important in vibration behaviors, was studied experimentally or analytically by Tobushi et al. [8], Dolce and Cardone [9], Gall and Maier [10] and Sun and Rajapakse [11]. As for the constitutive modeling for the cyclic effect, Tobushi et al. [8] and Sun and Rajapakse [11] proposed the use of critical transformation stresses depending on the number of cycles and the loading frequency, respectively. In the present study, the critical transformation stresses and the maximum residual strains are assumed to depend on the accumulated equivalent strain, which can be more conveniently used than the number of cycles in the general, dynamic analysis. The frequency dependence is not considered in the present modeling, as the target of the present study is relatively slow vibration analysis.

2 Description of Shape Memory Alloys

The special features of shape memory alloys are closely related to their internal phase transformation behavior. Generally, there are two stable phases in SMAs: austenite in high temperature and martensite in low temperature or under external loading. As they are shown in Fig. 1, austenite phase is cubic crystal structure, while martensite phases are monoclinic crystal structures. When we consider the SMA body as a whole, if an external load is applied in a horizontal direction, every lattice of the body is subjected to a shear stress. The direction of every lattice changes toward the same direction. When the stress is beyond a critical value, an observable macroscopic change occurs with the phase change from austenite or twinned martensite to detwinned martensite. Detwinned martensite can also be transformed to austenite

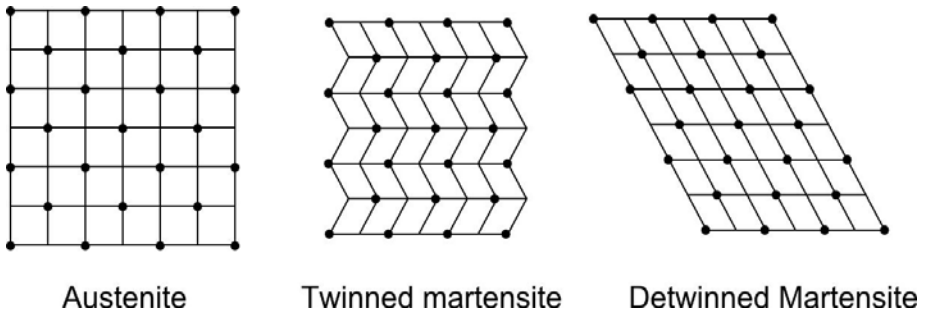


Fig. 1. Phase and crystal structure in SMAs

because of increasing Gibbs free energy when the SMA element is heated in the absence of applied stress.

Figure 2 shows the mechanical characteristics of SMAs. In Fig. 2(left), solid line and dotted line show the stress-strain curve of superelastic effect and shape memory effect respectively. Figure 2(right) is the critical stress-temperature diagram which can be used to determine the phase concentration in different temperature and stress settings. M_s and M_f indicate the starting and finishing temperature of phase transformation from austenite to martensite when stress is 0; while A_s and A_f indicate the starting and finishing temperature of phase transformation from martensite to austenite in the same setting. The sloping solid line indicates the critical phase transformation stress in different temperature. The phase transformation from austenite to martensite occurs when stress is increasing while it is between σ_{M_s} and σ_{M_f} . The phase transformation from martensite to austenite occurs when stress is decreasing

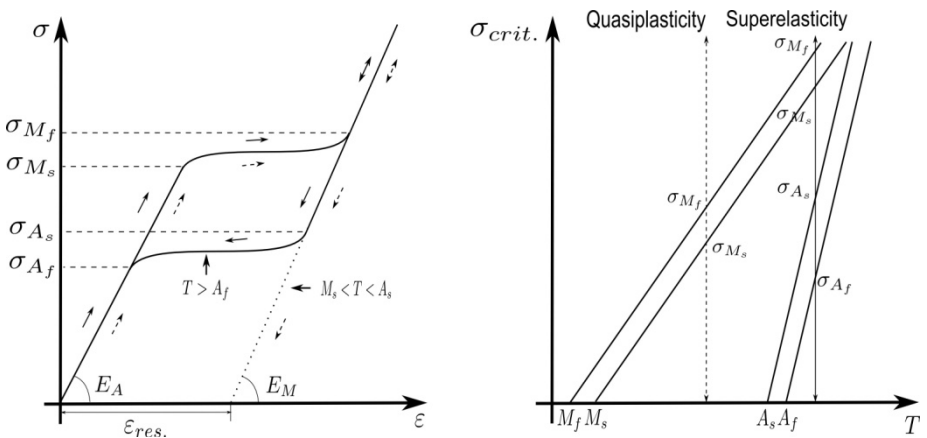


Fig. 2. Equivalent stress vs. strain curve (left); Mechanical property of SMAs (right)

while it is between σ_{A_s} and σ_{A_f} . As examples, the vertical solid line shows a typical superelastic effect route, while the vertical dotted line shows typical shape memory effect route, which correspond to the solid line and the dotted line in Fig. 2(left) respectively.

3 One-Dimensional Dynamic Simulation

3.1 Phase Transformation and Constitutive Equation

The general one-dimensional stress-strain relation of SMAs is written as

$$\sigma - \sigma_0 = E(\varepsilon - \varepsilon_0) + \Omega(\xi_S - \xi_{S0}) + \theta(T - T_0) . \quad (1)$$

where E is Young's modulus; Ω is the transformation coefficient; ξ_S is the stress-induced martensite volume fraction; θ is the thermal elastic coefficient; T is the temperature. The subscript '0' indicates the initial values.

E is expressed as

$$E = E_a + \xi(E_m - E_a) . \quad (2)$$

where E_a is Young's modulus of austenite phase, E_m is Young's modulus of martensite phase, ξ is the total martensite volume fraction.

The total martensite volume fraction ξ is expressed as

$$\xi = (\xi_S + \xi_T) . \quad (3)$$

where ξ_T is the temperature-induced martensite volume fraction.

The transformation coefficient Ω is expressed as

$$\Omega = -\varepsilon_L E . \quad (4)$$

where ε_L is the maximum residual strain, which is shown in Fig. 2 as ε_{res} . To consider the asymmetric behavior under tensile and compressive loading, Drucker-Prager equivalent stress σ^{DP} [7] is applied in phase transformation equations as

$$\sigma^{DP} = \sigma + \beta(\sigma_x + \sigma_y + \sigma_z) . \quad (5)$$

The phase transformation equations are expressed as:

Martensite Transformation, when $T > M_S$ and $C_{M_s}(T - M_s) < \frac{\sigma^{DP}}{(1+\beta)} < C_{M_f}(T - M_f)$

$$\xi_S = \frac{1-\xi_{S0}}{2} \cos \left\{ \frac{\pi}{C_{M_f}(T-M_f)-C_{M_s}(T-M_s)} \left[\frac{\sigma^{DP}}{(1+\beta)} - C_{M_f}(T - M_f) \right] \right\} + \frac{1+\xi_{S0}}{2} . \quad (6)$$

$$\xi_T = \xi_{T0} - \frac{\xi_{T0}}{1-\xi_{S0}} (\xi_S - \xi_{S0}) = \xi_{T0} \frac{1-\xi_S}{1-\xi_{S0}} . \quad (7)$$

$$\xi = (\xi_S + \xi_T) . \quad (8)$$

Reverse Transformation, when $T > A_S$ and $C_{A_f}(T - A_f) < \frac{\sigma^{DP}}{(1+\beta)} < C_{A_s}(T - A_s)$

$$\xi = \frac{\xi_{S0}}{2} \left\{ 1 + \cos \left[\pi \frac{C_{A_s}(T-A_s) - \frac{\sigma^{DP}}{(1+\beta)}}{C_{A_s}(T-A_s) - C_{A_f}(T-A_f)} \right] \right\} . \tag{9}$$

$$\xi_S = \frac{\xi_{S0}}{\xi_0} \xi . \tag{10}$$

$$\xi_T = \frac{\xi_{T0}}{\xi_0} \xi . \tag{11}$$

3.2 Numerical Examples

For simulations of this section, the same SMA bar is used. Material constants of this SMA bar are listed in Table 1.

Table 1. Material constants of SMA bar

Category	Constants(Unit)	Value
Moduli	E_a (MPa)	60000
	E_m (MPa)	20000
	C_{M_c} (MPa/°C)	8.0
	C_{M_f} (MPa/°C)	8.0
	C_{A_c} (MPa/°C)	13.0
	C_{A_f} (MPa/°C)	13.0
Transformation	M_f (°C)	-77.5
Temperature	M_s (°C)	-65.0
	A_f (°C)	-14.5
	A_s (°C)	-21.7
	Maximum residual strains	ϵ_L
	Mass(kg)	0.0087

Simulation of Superelastic Response. For simulations of symmetric superelastic response ($\beta = 0$ in eq. (5)), free vibration case is listed. The temperature of SMA bars is set at 10°C, which is higher than austenite phase transformation finish temperature A_f at -14.5°C. To present simulation results, the simulation time interval is set at 10^{-8} seconds. The horizontal initial velocity of SMA bar is set at 0.17m/s. According to Fig. 3, the upper graph for the displacement time-history shows a clear damping behavior, the amplitude keep decreasing until phase transformation no longer occurs. The lower graph shows a typical stress-strain curve of superelastic effect.

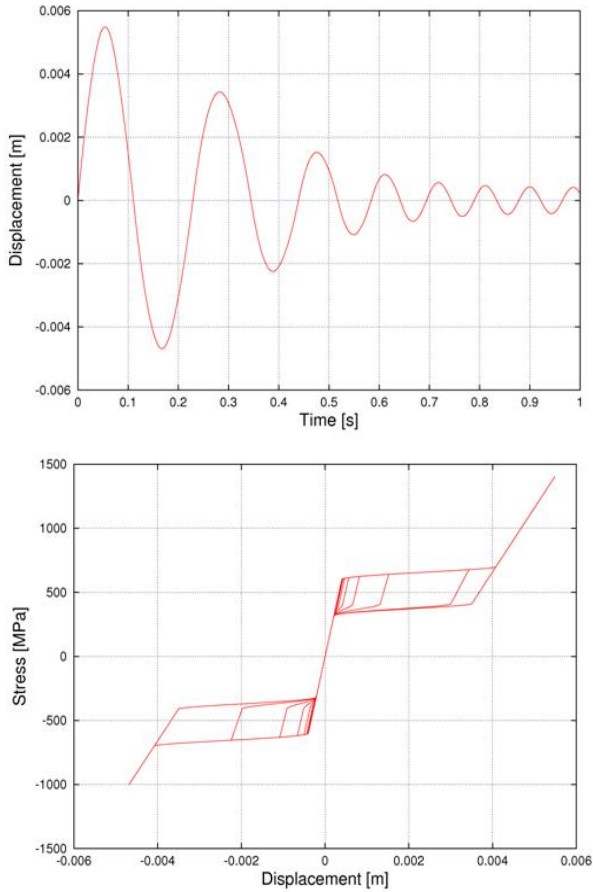


Fig. 3. One-dimensional superelastic free vibration response simulation: Displacement vs. time (upper), Stress vs. displacement (down)

The martensite phase transformation completes only in the first cycle and almost stops after the 5th cycle.

Simulation of Asymmetric Superelastic Response. For simulations of asymmetric superelastic response ($\beta = 0.15$ in eq. (5)), the displacement time-history and the stress-strain graphs are described in Fig. 4. The material constants, boundary condition and initial condition are exactly the same as those in the simulation of superelastic response. Lower level of phase transformations can be observed in compressive

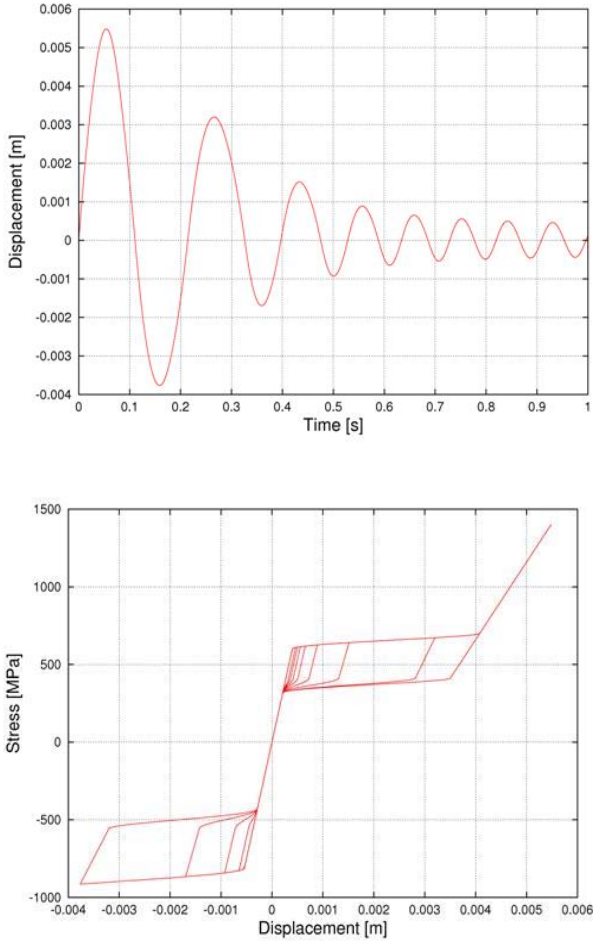


Fig. 4. One-dimensional superelastic free vibration response simulation in asymmetric setting: Displacement vs. time (upper), Stress vs. displacement (down)

side, which is caused by the higher requirement of stress for the phase transformation in asymmetric setting.

4 Three-Dimensional Finite Element Dynamic Simulation

4.1 Constitutive Equation

To describe the three-dimensional dynamic behavior of SMA devices, the one-dimensional stress-strain relation is extended into:

$$\{\sigma\} = [D]\{\varepsilon\} + \xi_S\{\Omega\} + T\{\theta\} . \tag{12}$$

where $\{\sigma\}$ is the stress vector; $[D]$ is the stress-strain matrix; $\{\varepsilon\}$ is the strain vector; ξ_S is the stress-induced martensite volume fraction; $\{\Omega\}$ is the transformation coefficient vector; T is the temperature; $\{\theta\}$ is the thermal elastic coefficient vector.

The transformation coefficient vector $\{\Omega\}$ is expressed as

$$\{\Omega\} = -[D][R_S]\{\varepsilon_L\} . \tag{13}$$

where $[R_S]$ is the residual strain direction matrix; $\{\varepsilon_L\}$ is the maximum residual strain vector. The thermal elastic coefficient vector $\{\theta\}$ is expressed as

$$\{\theta\} = -[D]\{\alpha\} . \tag{14}$$

where $\{\alpha\}$ is the coefficient of thermal expansion. By using Euler's formula for numerical time integration, we can derive the stress-strain relation equation as

$$\begin{aligned} \{\Delta\sigma\} = [D]\{\Delta\varepsilon\} - \xi_S[D][\Delta R_S\varepsilon_L]\{\Delta\sigma\} + \Delta t([\dot{D}]\{\varepsilon\} - \dot{\xi}_S[D][R_S]\{\varepsilon_L\} - \\ \xi_S[\dot{D}][R_S]\{\varepsilon_L\} - T[D]\{\dot{\alpha}\} - T[\dot{D}]) . \end{aligned} \tag{15}$$

or

$$\{\Delta\sigma\} = [X_D]^{-1}[D]\{\Delta\varepsilon\} + [X_D]^{-1}[Y_D] . \tag{16}$$

or

$$\{\Delta\sigma\} = [D_D]\{\Delta\varepsilon\} + [\Theta_D] . \tag{17}$$

where

$$[X_D] = [I] + \xi_S[D][\Delta R_S\varepsilon_L] . \tag{18}$$

$$[Y_D] = \Delta t([\dot{D}]\{\varepsilon\} - \dot{\xi}_S[D][R_S]\{\varepsilon_L\} - \xi_S[\dot{D}][R_S]\{\varepsilon_L\} - T[D]\{\dot{\alpha}\} - T[\dot{D}]) . \tag{19}$$

By using Newmark's β method, the finite element discretized equation of motion governed by the stress-strain relation above can be solved efficiently and accurately.

4.2 Numerical Examples

The forced vibration simulation in superelastic setting is described in this section. The simulation object is listed in Fig. 5. It is a highway-bridge-like concrete structure with

SMA dampers installed in their upper part (in black). The bottom of the structure is fixed in all directions. Therefore not only the dynamic behavior but also the damping capacity of SMAs will be shown in following subsection.

Simulation of Superelastic Forced Vibration Response. To show the three- dimensional SMA device's superelastic dynamic behavior, the temperature is set at 50°C, which is higher than the austenite phase transformation starting temperature $A_s = 42^\circ\text{C}$. The external force is a function of sine, which is expressed as

$$F = F_{\max} \sin(\beta \times \text{time}) . \tag{20}$$

In this simulation, the maximum external force is set at 4×10^7 N. The frequency coefficient β is set at 0.1256; therefore the cycle time is about 50s. The number of eight-node hexahedron elements is 1, 2, and 13 in x-, y- and z-direction, respectively. The calculated result is shown in Fig. 6. Similar to the one-dimensional simulation, a clear stress-strain curve can be observed. But the displacement history is much more complicated than that in the one-dimensional simulation. The concrete elements of this structure show an elastic behavior according to its equivalent stress-strain curve.

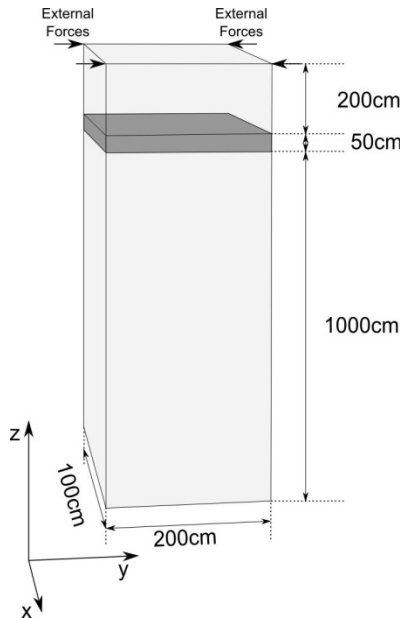


Fig. 5. Forced vibration simulation object

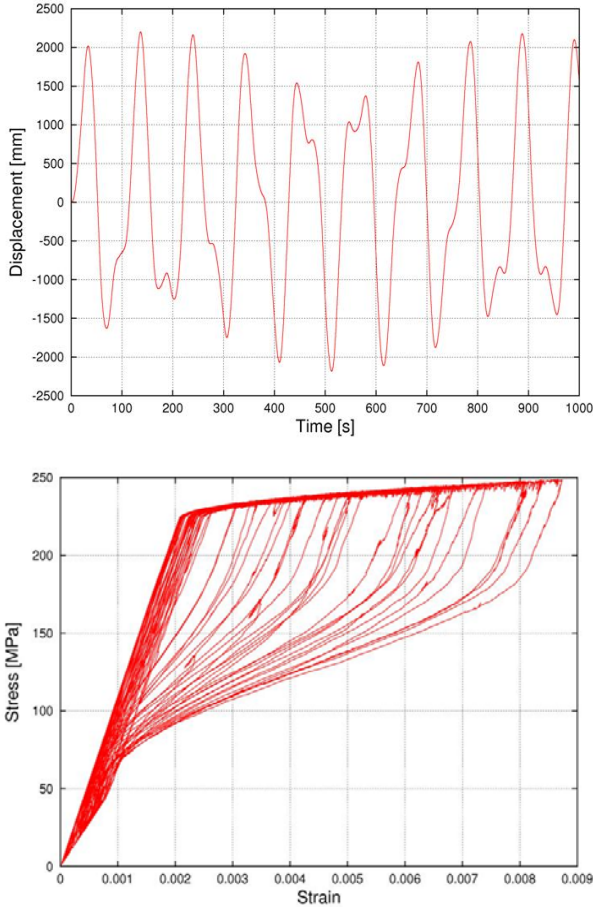


Fig. 6. Three-dimensional superelastic forced vibration response simulation: top displacement vs. time (upper) equivalent stress vs. strain of SMA edge (down)

5 Improved Model for Cyclic Behavior of Shape Memory Alloys

5.1 Constitutive Equation

The constitutive models described in section 3 and section 4 are also applicable for one-dimensional and three-dimensional cyclic dynamic simulation. After careful analysis of experiment data from Tobushi et al. [8], Dolce and Cardone [9], Gall and Maier [10] and Sun and Rajapakse [11], it is found that the cyclic behavior of SMA devices can also be described by making a connection between the accumulative strain and the martensite phase transformation start and finish temperature M_s and

M_f , the reverse phase transformation start and finish temperature A_s and A_f , as well as the maximum residual strain ε_L . The accumulative strain is expressed as follows

$$e = \int \dot{\varepsilon}_{eq} dt . \quad (21)$$

where e is the accumulative strain, $\dot{\varepsilon}_{eq}$ is the equivalent strain rate. The connection between the accumulative strain e and the martensite phase transformation start and finish temperature M_s and M_f , as well as the reverse phase transformation start and finish temperature A_s and A_f are expressed as

$$M_s(e) = M_s(0) + \alpha_{M_s} \left(1 - \exp(-\beta_{M_s} e) \right) . \quad (22)$$

$$M_f(e) = M_f(0) + \alpha_{M_f} \left(1 - \exp(-\beta_{M_f} e) \right) . \quad (23)$$

$$A_s(e) = A_s(0) + \alpha_{A_s} \left(1 - \exp(-\beta_{A_s} e) \right) . \quad (24)$$

$$A_f(e) = A_f(0) + \alpha_{A_f} \left(1 - \exp(-\beta_{A_f} e) \right) . \quad (25)$$

where $M_s(0), M_f(0), A_s(0)$ and $A_f(0)$ are the phase transformation critical temperatures. α_A, α_M and β_A, β_M are the material constants, which have to be determined by experiment data. The maximum residual strain is divided into two coefficient, ε_L^M , which is the maximum residual strain during martensite phase transformation, and ε_L^A , which is the maximum residual strain during reverse phase transformation. ε_L^M and ε_L^A are expressed as

$$\varepsilon_L^A(e) = \varepsilon_L(0) + \alpha_L^A \left(1 - \exp(-\beta_L^A e) \right) . \quad (26)$$

$$\varepsilon_L^M(e) = \varepsilon_L(0) - \alpha_L^M \left(1 - \exp(-\beta_L^M e) \right) . \quad (27)$$

where $\varepsilon_L(0)$ is called the initial maximum residual strain; $\alpha_L^A, \alpha_L^M, \beta_L^A$ and β_L^M are the coefficients which have to be determined by experiment data. By substituting the equations above into the martensite phase transformation equation which was described in section 2, we will be able to obtain the cyclic dynamic behavior of SMA devices.

5.2 Numerical Examples

According to the experimental data provided by Tobushi et al. [8], the simulation object (SMA bar) is similar to the one which is described in the section 2. Unloading starts as soon as the strain reaches 9%. Detailed simulation results are listed in Fig. 7.

The experimental data shows the stress and strain relationship during No. 1, 2, 5, 10, 20, 50, 100 cycles. The initial maximum residual strain is set at 3.5%, the temperature is set at 353K. The experimental data shows that critical stresses of start

and finish martensite phase transformation as well as the reverse transformation are getting lower and lower, while the maximum residual strain is getting larger and larger. The simulation result which is presented in lower graph of Fig. 7 shows a similar behavior. The critical stresses and the maximum residual strain are changing faster at the beginning then slowing down, because of the term $(1-\exp(-X))$ in coefficients revolution function. Only the cycles No. 1, 2, 5, 10, 20, 50 and 100 are presented in Fig. 7 to show a better comparison with the experimental data.

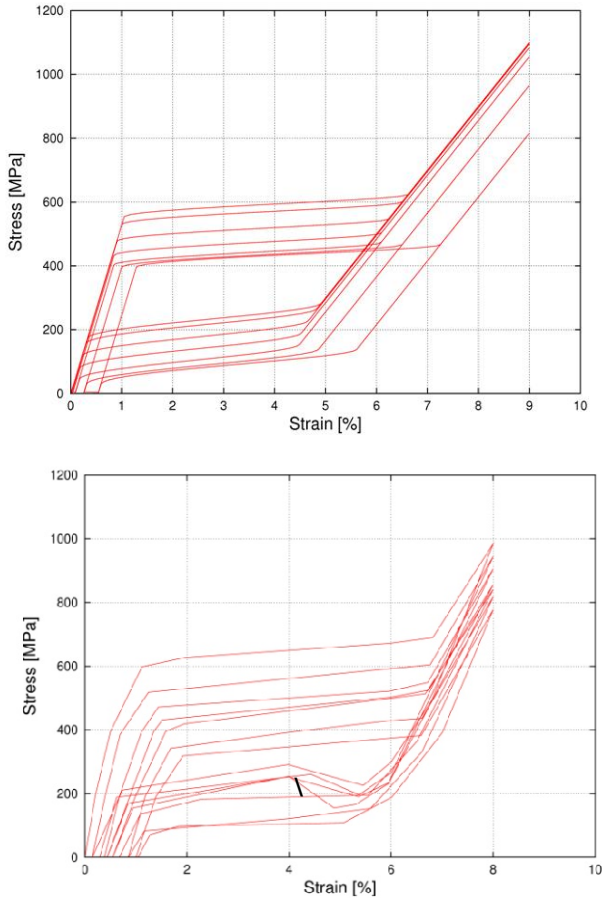


Fig. 7. Experimental data and simulation results under cyclic loading: simulation data (upper), experimental data (down)

6 Concluding Remarks

In this paper, the one-dimensional and three-dimensional dynamic response models for SMA devices have been developed. Besides, model improvement works have

been done by introducing cyclic effect and asymmetric behavior of SMA devices into the phase transformation mechanism. According to the simulations performed in this paper, several conclusions can be obtained as follows:

1. The validity of the dynamic response simulation model as well as the improved phase transformation mechanism to describe SMAs' superelastic and quasi-plastic behavior has been illustrated.
2. An improved SMA model by considering the cyclic effect has been developed and tested.
3. The potential of SMAs being used as damping devices has been demonstrated.

The future works will be focused on enhancing the simulation model's stability, especially for high frequency vibration simulations in three-dimensional scales.

References

1. Sugiyama, Y., Hirai, S.: Crawling and jumping by a deformable robot. *International Journal of Robotics Research* 25(5), 603 (2006)
2. Brinson, L.: One-dimensional constitutive behavior of shape memory alloys: thermomechanical derivation with non-constant material functions and redefined martensite internal variable. *Journal of Intelligent Material Systems and Structures* 4(2), 229 (1993)
3. Toi, Y., Lee, J., Taya, M.: Finite element analysis of superelastic, large deformation behavior of shape memory alloy helical springs. *Computers and Structures* 82(20-21), 1685–1693 (2004)
4. Seelecke, S.: Modeling the dynamic behavior of shape memory alloys. *International Journal of Non-Linear Mechanics* 37(8), 1363–1374 (2002)
5. Achenbach, M.: A model for an alloy with shape memory. *International Journal of Plasticity* 5(4), 371–395 (1989)
6. Toi, Y., Choi, D.: Constitutive Modeling of Porous Shape Memory Alloys Considering Strain Rate Effect. *Journal of Computational Science and Technology* 2(4), 511–522 (2008)
7. Toi, Y., Choi, D.: Computational Modeling of Superelastic Behaviors of Shape Memory Alloy Devices under Combined Stresses. *Journal of Computational Science and Technology* 2(4), 535–546 (2008)
8. Tobushi, H., Iwanaga, H., Tanaka, K., Hori, T., Sawada, T.: Stress-strain-temperature relationships of TiNi shape memory alloy suitable for thermomechanical cycling. *JSME International Journal. Series 1, Solid Mechanics, Strength of Materials* 35(3), 271–277 (1992)
9. Dolce, M., Cardone, D.: Mechanical behaviour of shape memory alloys for seismic applications 2. Austenite NiTi wires subjected to tension. *International Journal of Mechanical Sciences* 43(11), 2657–2677 (2001)
10. Gall, K., Maier, H.: Cyclic deformation mechanisms in precipitated NiTi shape memory alloys. *Acta Materialia* 50(18), 4643–4657 (2002)
11. Sun, S., Rajapakse, R.: Simulation of pseudoelastic behaviour of SMA under cyclic loading. *Computational Materials Science* 28(3-4), 663–674 (2003)

Adaptive Fuzzy Filter for Speech Enhancement

Chih-Chia Yao and Ming-Hsun Tsai

Department of Computer Science and Information Engineering,
Chaoyang University of Technology, Wufong, Taichung 41349, Taiwan
ccyao@cyut.edu.tw

Abstract. In this paper an adaptive fuzzy filter, based on fuzzy system, is proposed for speech signal enhancement and automatic speech recognition accuracy. In the past two decades the basic wavelet thresholding-algorithm has been widely studied and is common applied to filter noise. In the proposed system adaptive wavelet thresholds are generated and controlled by fuzzy rules concerning the presence of speech in noise. First an amplified voice activity detector is designed to improve performance on SNR lower than 5dB. Then an adaptive threshold decision module based on fuzzy inference system is proposed. In this fuzzy inference system overall relations between speech and noise are summarized into seven fuzzy rules and four linguistic variables, which are used to detect the state of signals. The adaptive threshold and membership functions are optimally obtained by particle swarm optimization algorithm so the SNR of the filter output for training signal data can be maximized. Experimental results reveal that our proposed system effectively increases the SNR and the recognition rate of speech.

Keywords: speech enhancement; wavelet thresholding; fuzzy; voice activity detection.

1 Introduction

Speech enhancement remains an important problem within the field of speech and signal processing, such as mobile communication, speech recognition and hearing aids [1]. Speech enhancement algorithms have undergone numerous studies throughout the last two decades yet issues such as distortions to the original speech signal and residual noise, sometimes in the form of musical tones created by the enhancement algorithms, remain unsolved [2, 3, 4].

The main objective of speech enhancement is to reduce corrupting noise while preserving as much original clean speech quality as possible. The objectives of speech enhancement algorithms can be classified into increasing the recognition rate and the signal-to-noise ratio (SNR) of the distorted speech [5, 6]. These can be roughly classified as digital signal processing and statistical analysis. The former usually remove an estimate of the distortion from the noisy features, e.g. spectral subtractive [7]. The latter utilizes a statistical model to describe predictable structures and expected patterns in the signal process [8, 9].

The principle of a spectral subtractive algorithm is to obtain the best possible estimations of the short time spectra of a speech signal from a given contaminated speech. The main attractions of a spectral subtractive algorithm are: (1) Easy and simple implementation. (2) High flexibility against subtraction parameter's variation. On the other hand, spectral subtractive algorithms suffer a number of deficiencies including an imprecise estimation of signal and noise parameters and a mismatch of the probability distribution models of speech or noise [10].

The wavelet thresholding method is often used to shrink the signal and remove the noise. Using an adaptive thresholding algorithm could provide better performance [11, 12]. In this paper we propose an adaptive fuzzy filter, based on wavelet thresholding, to improve performance in speech signal enhancement and automatic speech recognition accuracy. In this developed adaptive fuzzy filter the signals are classified into speech and non-speech segments by an amplified voice activity detector. This amplified voice activity detector employs the full wavelet packet transform (instead of the perceptual wavelet packet transform) to decompose the input speech signal into critical sub-band signals [13]. In each critical sub-band signal the mask construction is obtained by smoothing the Teager energy operator and the entropy of corresponding wavelet coefficients. Thus, the amplified voice activity detector enhances the ability of distinction when the SNR is lower than 5dB.

Although some algorithms have been proposed to generate the adaptive threshold the model between speech and noise is still not clear and the rules judging the existence of noise is usually expressed in an ambiguous style. A new-type of speech enhancement system based on fuzzy rules is proposed here. In the design overall relations between speech and noise are summarized into seven fuzzy rules. Adaptive thresholds are then inferred by those fuzzy rules. Moreover, this filter is optimally designed by particle swarm optimization, so that the SNR of the filter output is maximum; membership functions to represent the rules are not required to be obtained beforehand.

As is well know, using critical parameters is helpful on detecting the presence of speech in noisy environment. Some researchers have focused on energy, zero crossing rate and time duration to find the boundary between the word signal and background noise. Other parameters have been proposed such as linear prediction coefficient, linear prediction error energy, pitch information and time-frequency parameter. However these parameters cannot yet be well adapted to variable-level background noise even if more complex decision strategies are used. In this paper four parameters used as linguistic variables are incorporated in order to detect the presence of speech in noisy environment. The four parameters are energy, zero crossing rate, entropy and standard deviation. Examples of processing actual speech with eight types of noises are shown to demonstrate the high performance of this filter.

The remainder of this paper is organized as follows. In Section 2 we review the basic concept of wavelet transform and wavelet packet decomposition. In Section 3 the scheme of the speech enhancement system is introduced. In Sections 4 we introduce the fuzzy inference model for speech enhancement. In Section 5 a performance evaluation of the proposed system is presented and comparisons with other protocols are made. Our conclusions are made in Section 6.

2 Basic Concepts

2.1 Wavelet Packet Transformation

Wavelet Transformation is widely used in signal processing because processing signals in the frequency domain is often easier to implement [13]. In comparison to short time Fourier transform, wavelet transform allows the use of long-time intervals to obtain more precise low-frequency information and shorter intervals for high-frequency information. The continuous wavelet transform (CWT) of a signal $x(t)$ is given as follows:

$$X_{CWT}(a, \tau) = \frac{1}{\sqrt{a}} \int_{-\infty}^{\infty} x(t) \psi\left(\frac{t-\tau}{a}\right) dt \quad (1)$$

where τ and a represent the time shift and scale variables, respectively, and $\psi(\cdot)$ is the mother wavelet chosen for the transform.

In the discrete version, the wavelet decomposes the signal with variable frames to perform multi-resolution analysis in a dyadic form known as discrete wavelet transform (DWT). In DWT the scale and translation parameters of the discrete wavelet family are given by $a = 2^m$ and $\tau = n2^m$.

Signals transformed by DWT can be regarded as splitting signals by a series of high-pass and low-pass filters. The high-frequency, low-scale components of speech signals through a high-pass filter are preserved as details. In the same way, the low-frequency, high-scale components are preserved as approximations. In discrete wavelet transform only approximations can be decomposed iteratively.

A further generalization of the DWT is the wavelet packet transform (WPT) that offers a richer range of possibilities in signal analysis. In WPT the decomposing process is iterated on both high and low frequency components rather than continuing only on low frequency terms as with a standard DWT.

The principle of wavelet packet transform is that, given a signal, a pair of low pass and high pass filters is used to yield two sequences to capture different frequency sub-band features of the original signal. The depth of the wavelet packet tree can be varied over the available frequency range, resulting in configurable decomposition. The two wavelet orthogonal bases are defined as

$$\psi_{j+1}^{2^p}(k) = \sum_{n=-\infty}^{\infty} h[n] \psi_j^p(k - 2^j n) \quad (2)$$

$$\psi_{j+1}^{2^{p+1}}(k) = \sum_{n=-\infty}^{\infty} g[n] \psi_j^p(k - 2^j n) \quad (3)$$

where $h[n]$ and $g[n]$ denote the low-pass and high-pass filters, respectively. $\psi(n)$ is the wavelet function and parameters j and p are the number of decomposition levels and nodes, respectively.

2.2 Particle Swarm Optimization

In 1995, Kennedy and Eberhart introduced the particle swarm optimization algorithm (PSO) into the field of social and cognitive behavior [14]. Traditionally the main problem in designing a neural fuzzy system is training the parameters. Backpropagation

training is commonly adopted to solve this problem. However the steepest descent approach, commonly used in backpropagation training to minimize the error function, may reach the local minima very quickly and never find the global solution. Accordingly, a new optimization algorithm, called particle swarm optimization (PSO), appears to provide better performance than the backpropagation algorithm.

Like other population-based optimization approaches PSO is initialized with a swarm of random solutions, each swarm consists of many particles. Each particle is characterized by its current position $\vec{x}_i = [x_i^1, x_i^2, \dots, x_i^D]$ and current velocity $\vec{v}_i = [v_i^1, v_i^2, \dots, v_i^D]$, where D stands for the dimensions of the solution space. In the PSO the trajectory of each particle in the search space is adjusted by dynamically altering the velocity of each particle. Then the particles rapidly search the solution space using the moving velocity of each particle. Each of these particle positions is scored to obtain a fitness value based on how to define the solution of the problem. During the evolutionary process the velocity and position of particle i are updated as

$$\vec{v}_i = \omega \times \vec{v}_i + \varphi_1 \times rand_1() \times (NBest_i - \vec{x}_i) + \varphi_2 \times rand_2() \times (GBest_i - \vec{x}_i) \quad (4)$$

$$\vec{x}_i = \vec{x}_i + \vec{v}_i \quad (5)$$

where ω is the inertia weight, φ_1 and φ_2 are the acceleration coefficients, respectively. The second term in Eq. (4), called the cognitive component, reflects the experience of a particle since it is dependent on the best position of the respective particle. The third term is referred to as the social component and contains the information of a social group due to the dependence on the neighborhood best position. The random numbers $rand_1()$ and $rand_2()$ are chosen from the interval $U(0,1)$. In Eq. (4), $GBest$ is the position with the best fitness found so far for the i th particle, and $NBest$ is the best position in the neighborhood. The term \vec{v}_i is limited to the range $\pm \vec{v}_{max}$. If the velocity violates this limit, then it is set to the actual limit. Changing the velocity enables each particle to search around its individual best position and global best position. After initialization of the positions and velocities of the particles update equations are applied to every particle in each iteration until a stopping criterion is fulfilled.

3 Architecture of Adaptive Fuzzy Filter

This section proposes a new framework for speech enhancement called an Adaptive Fuzzy Filter. This utilizes a fuzzy inference system to determine the threshold used to suppress noise. The schematic diagram of an adaptive fuzzy filter is shown in Fig. 1. The structure of an adaptive fuzzy filter comprises wavelet packet transform, a voice activity detection module, a fuzzy inference module, a learning module and a speech filter.

In the adaptive fuzzy filter wavelet packet transform is applied to each input frame and a Level-3 wavelet packet decomposition tree with eight frequency bands is obtained. A voice activity detection module is then used to distinguish the presence or absence of speech. The fuzzy inference module determines the adaptive threshold used in the speech filter in order to suppress noise. After processing the inverse wavelet packet transform is applied to provide the enhanced speech. Particle swarm optimization is used in the learning module to provide the optimal solution.

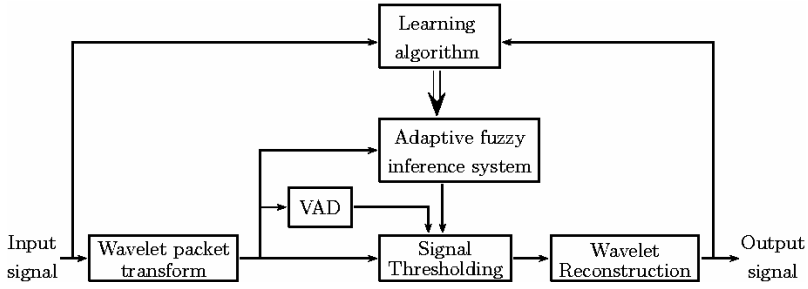


Fig. 1. The schematic diagram of fuzzy speech enhancement model

3.1 Voice Activity Detection (VAD)

Voice activity detection (VAD) is used to distinguish speech from contaminated speech signals and is required in a variety of speech communication systems. In previous studies a Teager energy operator has been adopted to distinguish noise and speech and has been proven to provide excellent performance in both additive noisy and real noisy environments [15]. The discrete form of the TEO is given by

$$T[y(n)] = y^2(n) - y(n+1)y(n-1) \tag{6}$$

where $T[y(n)]$ is called the TEO coefficient of $y(n)$. However a Teager energy operator is insensitive when the signal-to-noise ratio is low. As an example, when the SNR is lower than 0dB the difference between noisy energy and speech energy is not obvious and the performance of TEO is not satisfactory. Fig. 2 shows an example of using TEO to distinguish speech from noise signals contaminated by -5dB Gaussian white noise. To overcome this problem, in our proposed module, a Teager energy operator is cooperated with entropy to improve the ability of distinction. The formula of entropy is shown as follow:

$$E(S) = \sum_i p(s_i) \log p(s_i) \tag{7}$$

Entropy represents the degree of variation. When the noise model is stationary or slight non-stationary the corresponding entropy is kept stable or slight changed. Fig. 3 shows the comparison of using TEO and entropy on VAD to distinguish signals contaminated by -5dB Gaussian white noise.

The proposed VAD algorithm computes signals $w_m(x_i)$ that have been produced by the wavelet package transform on each input frame to produce 2^J sub-band wavelet packets, where J is the number of levels for the wavelet packet decomposition tree and $1 \leq m \leq 2^J$. Then, a set of $T[w_m(x_i)]$ can be derived from Eq. (6). The scheme of the voice activity detection is designed as follows:

$$V_m(x_i) = \begin{cases} 0, & \text{if } \text{var}(T[w_m(x_i)]) < \lambda_j \text{ or } \text{var}(E[w_m(x_i)]) < \zeta \\ x_i, & \text{otherwise} \end{cases} \tag{8}$$

where ζ is user-defined and $\lambda_j = \sigma_j \sqrt{2 \log(N_j)}$, proposed by Johnston and Silverman [16]. $\text{var}(\cdot)$ denotes the variance.

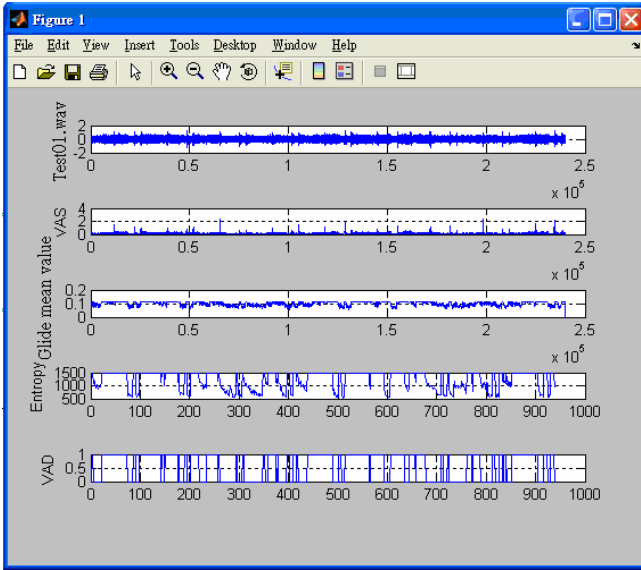


Fig. 2. The property of Gaussian white noise with -5db

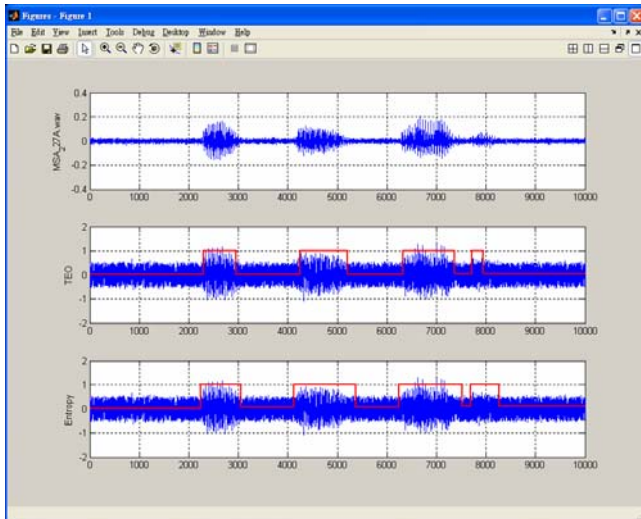


Fig. 3. The comparison of VAD based on TEO and entropy with -5db Gaussian white noise

4 Adaptive Threshold Decision Module Based on Fuzzy Inference System

In this section, an adaptive threshold decision module based on fuzzy inference system is proposed. In the module the adaptive threshold for each sub-band signal is

generated by the fuzzy inference system and the noise components were removed by thresholding the wavelet coefficients.

Let $s(t)$ be clean speech with finite length and $n(t)$ be noise, then contaminated speech $y(t)$ can be expressed as

$$y(t) = s(t) + n(t) \quad (9)$$

If W denotes the wavelet transform matrix, Eq. (9) can be written in the wavelet domain as

$$Y(t) = S(t) + N(t) \quad (10)$$

where $Y(t) = W \cdot y(t)$, $S(t) = W \cdot s(t)$ and $N(t) = W \cdot n(t)$. The estimated speech signal $\hat{S}(t)$ can be obtained by using the thresholding function

$$\hat{S}(t) = F_T(Y, T) \quad (11)$$

where $F_T(Y, T)$ denotes the thresholding function and T is the threshold. The standard thresholding function includes the soft thresholding function which is defined as

$$F_{T_s}(Y, T) = \begin{cases} \text{sign}(Y)(|Y| - T), & |Y| \geq T \\ 0, & |Y| < T \end{cases} \quad (12)$$

In this paper an adaptive threshold decision module based on a fuzzy inference system was proposed. In the fuzzy inference system overall relations between speech and noise are summarized to seven fuzzy rules. Moreover, four linguistic variables are used to detect the state of signals. They are energy, zero crossing rate, standard deviation and average residual. The reasons for adopting those four parameters are listed as follows.

-Energy:

Energy is an effective factor on measuring the degree of noise when the SNR is bigger than 0dB. Studies have shown that energy provides excellent performance on voice activity detection. However, rapid variation of energy in the speech model causes implementation difficulties.

-Zero crossing rate:

Zero crossing rate (ZCR) is a basic acoustic feature that can be computed easily. It is equal to the number of zero crossing of the waveform within a given frame. In general the ZCR of both unvoiced sounds and environment noise are larger than voiced sounds. ZCR is often used in conjunction with volume for end-point detection. It is hard to distinguish unvoiced sounds from environment noise by using ZCR alone since they have similar ZCR values.

-Standard deviation:

Standard deviation is a measure of how wide any given numbers are spread. It is useful in comparing sets of data, which may have the same mean but different range. In a given frame standard deviation is helpful to distinguish the signal's mode. The standard deviation of a stable noise is different to voice sound. The observation of standard deviation is useful in the threshold decision level.

-Average residual:

The statistical property of average residual is similar to standard deviation. In this paper standard deviation is used to detect the variation within a frame but average residual is used to compare the difference between frames. By calculating the average residual across several frames the signal's mode can be distinguished, whether the variation is temporal or not, in a very short time interval. In this paper the average residual is calculated based on the standard deviations of previous frame, current frame and next frame.

In the fuzzy inference system, Eq. (13) is used to obtain the optimal solution.

$$T(k+1) = T(k) + \alpha(k)THS(k) \quad (13)$$

Here, $\alpha(k)$ is a step length for threshold's variation and $THS(k)$ is the k th threshold. $\alpha(k)$ and $THS(k)$ are decided by above four parameters. First the range of threshold is decided by energy and zero crossing rate. Studies have shown that energy is an effective factor in measuring the degree of noise greater than 0dB. However when the threshold is only decided by energy any rapid variation of energy causes difficulty and irrational implementation. On top of which, most mechanisms generally use an 8k sample rate to record the speech signal and, when the frequency of the noise is bigger than the sample rate, energy is ineffective. To overcome above problem zero crossing rate is introduced to cooperate with energy on deciding the threshold.

As described in previous statements the ZCR of unvoiced sounds and environment noise are larger than the ZCR of voiced sounds. In many studies a zero crossing rate is used on voice activity detection. The signal's mode can be easily distinguished by using zero crossing rate. When the noise is high frequency and the speech is low frequency the zero crossing rate can be taken as the auxiliary and be cooperated with energy. When the frequency of noise is low the threshold decision is highly dependent on energy. As the frequency of noise becomes higher zero crossing rate is taken into consideration. In short the importance of zero crossing rate on the threshold decision is proportional to the frequency of noise.

Standard deviation and average residual are then introduced to control the degree of threshold variation. The degree of variation on signal's modes between two adjacent frames is reflected in the difference between standard deviations. As the difference between standard deviations increases so the amount of threshold variation increases. However this does not apply when speech is contaminated by instantaneous impulse signal. An instantaneous impulse signal in a given frame will usually increase the standard deviation and often causes a violent vibration of threshold. To avoid this problem, the average residual of any adjacent three frames is calculated to smooth the threshold's variation. The average residual is calculated from standard deviations of any adjacent three frames and can effectively reduce the impact of instantaneous impulse signals. Should the value of the average residual increase there is a high probability that it is neither temporal nor instantaneous.

4.1 Fuzzy Inference Rules

The threshold can be set by the local characteristics of the input signals but, as previously mentioned, there are many ambiguous cases where constructing the noise model

Table 1. Linguistic variables - Fuzzy set relational table

Linguistic Variable	Fuzzy Sets		
	Low	High	
Energy	Low	High	
Zero crossing rate	Low	High	
Standard deviation	Low	High	
Average residual	Low	High	
Threshold	Low	Median	High
Step length	Low	Median	High

from the contaminated signal is hard to achieve. In the proposed speech enhancement system a fuzzy inference system is adopted to overcome these problems and the adaptive threshold is generated in order to suppress noise. To begin with various models of noises and the relation between noise and speech signals are summarized into several cases. The fuzzy sets of linguistic variables and their corresponding linguistic terms are then defined according to the analysis. The relations between linguistic variables and their fuzzy sets are shown in Table 1.

After eliminating the cases where the noise is impossible to distinguish from speech signal, seven fuzzy rules are concluded, listed as follows.

Rule 1: If Energy is High and Zero Crossing Rate is Low, then Threshold is Medium.

Energy is an effective factor on measuring the degree of noise when the SNR is bigger than 0dB. In most cases speech signal has the property of high energy as opposed to noise. In attempt to retain speech and suppress noise, the threshold can be set lower than the speech signal. This rule is applicable in places like airports or streets with heavy traffic. In this model the zero crossing rate is classified to low when the frequency is between 50Hz and 5 KHz.

Rule 2: If Energy is High and Zero Crossing Rate is High, then Threshold is High.

As the zero crossing rate becomes higher, the probability of noise existence increases. In rule 2 the probability of the existence of noise is higher than in rule 1, so that the threshold can be set higher than rule 1. This rule is applicable when the energy is large and the frequency is high such as noise generated by high frequency electric products or machines. In this model the zero crossing rate is classified to high when the frequency is bigger than 5K Hz.

Rule 3: If Energy is Low and Zero Crossing Rate is Low, then Threshold is Low.

Low zero crossing rate usually means the probability of noise existence is not high. In most case the max energy is detected from the speech signal. Therefore, the threshold is set to low when the energy of contaminated speech signal is low. The type of noises is generated such as by the low frequency electric products (indoor fan, ..., etc).

Rule 4: If Energy is Low and Zero Crossing Rate is High, then Threshold is Medium.

On the basis of rule 4, high zero crossing rate denotes a high probability of noise. The threshold should be bigger than in rule 3 but not exceed the energy of speech signal. This rule is applicable on the environments such as the electromagnetic wave.

Rule 5: If Standard deviation is Low and Average residual is Low, then Step length is Low.

This rule is applicable when the model of noise is stable, whether in a very short time interval or in a longer time interval. The property of this noise is that the standard deviation is small and the average residual calculated from three successive frame blocks' standard deviations is also small. Since the standard deviation is small, the low threshold is necessary. The waveform under this situation is shown in Fig. 4. In Fig. 4, vertical axis represents the amplitude and horizontal axis represents the time. Besides, in Fig. 4, each frame is bounded by the rectangle. The standard deviation of each frame is small as the waveform within each frame is stable. The average residual of the successive three frame's standard deviation is also small since the standard deviations are all small.

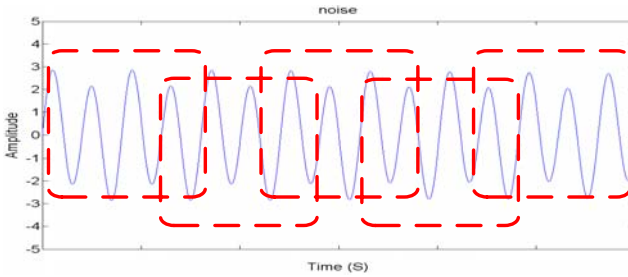


Fig. 4. The property of the noise is low standard deviation and low average residual

Rule 6: If Standard deviation is High and Average residual is Low, then Step length is Medium.

This rule is applicable when noise is stable but variant in a very short time interval. The property of this noise is that the standard deviation is high but the residual that is calculated from current and neighbor frame's standard deviations is small. In rule 6 high standard deviation means the amplitude of noise is higher than the noise in rule 5. Hence, the threshold should be set higher than in rule 5. Under the situation with multiple sources of noise, this rule can bring the function into full play. The waveform under this situation is shown in Fig. 5. In Fig. 5, a big variation of the waveform within the middle frame causes the standard deviation to be big. The average residual of the successive three frame's standard deviation is small since the standard deviations are all big.

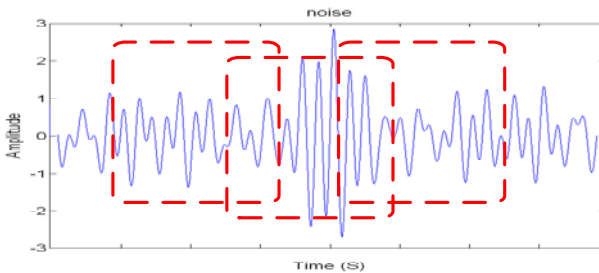


Fig. 5. The property of the noise is high standard deviation and low average residual

Rule 7: If Standard deviation is High and Average residual is High, then Step length is High.

This rule is applicable when noise is unstable, whether in a very short time interval or in a long time interval. In rule 7 the current frame's standard deviation is high and is higher than neighbor frames so that the average residual calculated from current and neighbor frame's standard deviations is high. The waveform under this situation is shown in Fig. 6.

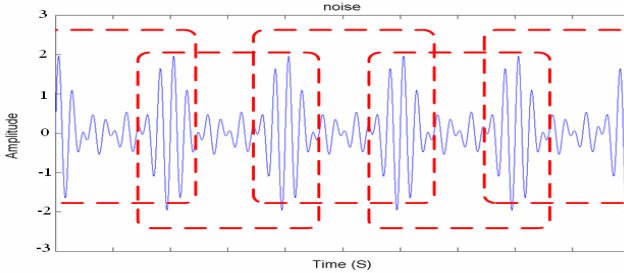


Fig. 6. The property of the noise is high standard deviation and high average residual

The next step is to design this function “*Threshold*” and “*Step length*” for Energy, Zero crossing rate, Standard deviation and Average residual. In most fuzzy systems membership function is obtained by fuzzy approximate reasoning using the membership function of the *Low* and the *High* of Energy, Zero crossing rate, Standard deviation, Average residual and its relationship to *Threshold* and *Step length*. In this paper the Gaussian type and the sigmoid function are used to approximate the membership function. The benefit is efficient calculation. In the field of signal processing the goal is to obtain the optimum filter, in which the output is as close to the original signal as possible; thus some optimization is required in designing this filter. The particle swarm optimization algorithm (PSO) is applied to optimize the parameters in these nonlinear function. As we know, particle swarm optimization algorithm provides better performance on finding the global solution while optimizing the overall structure. Experimental results demonstrate that our proposed system provides excellent performance on speech enhancement.

5 Experimental Results

This section demonstrates the effectiveness of our proposed system. The experimental results that pertain to our proposed speech enhancement system are compared to spectrum subtraction and signal subspace approach.

First, tests for the enhancement and recognition task were carried out on the Aurora 2 databasen [17]. In the database the speech files include recording of the 10 English digits and 26 English letters spoken by males and females. The Aurora 2 database does not only provide the uncontaminated speech of 4004 sentences it also includes the speech of 48048 sentences contaminated with eight types of noise (Subway, Babble, Car, Exhibition, Restaurant, Street, Airport and Station) and are mixed with -5, 0, 5, 10, 15 and 20dB.

In order to achieve a consistent comparison the HTK speech recognition system is adopted as the classifier [18]. The HTK speech recognition system was first developed by the Speech Vision and Robotics Group of Cambridge University. With this tool, the HMM model, language model, training model and etc. can be established. In the Aurora 2 database, HTK has been adopted for the training of the model and for recognition. Therefore, we make use of the same approach to configure and establish the model, and process for recognition.

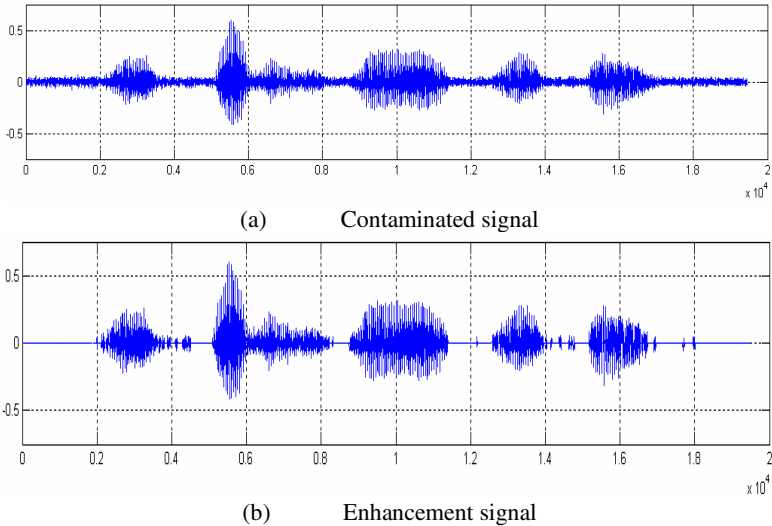


Fig. 7. (a) Speech signal contaminated by 10dB car noise. (b) Enhancement signals denoised by our proposed system.

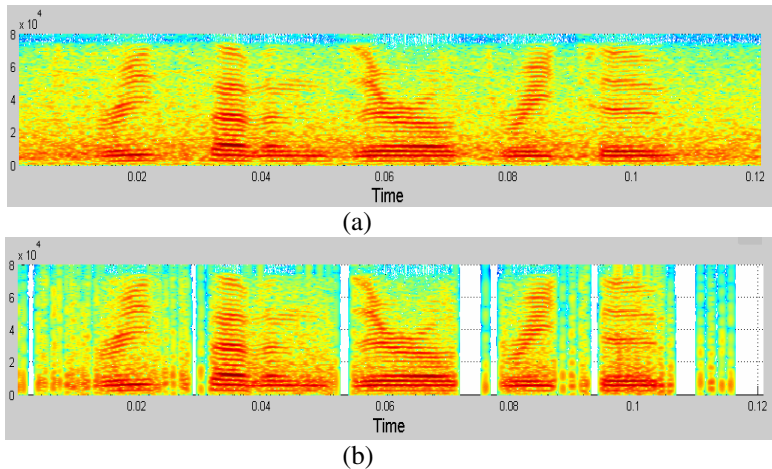


Fig. 8. (a) The frequency of speech signal contaminated by 10dB car noise. (b) The frequency of enhancement signals denoised by our proposed system.

The evaluation of the proposed system includes two parts. First the signal-to-noise ratios on before and after enhanced signals are evaluated. Next applying those enhanced speech signals into the HTK classifier compares the recognition rates of our proposed system with other methods. The experimental results are listed as follows.

Fig. 7 shows the waveform before and after signal denoising by our proposed system. The signal is mixed with 10dB car noise. The corresponding time-frequency diagrams are reveals in Fig. 8.

The comparisons of SNR are included to demonstrate the effectiveness of our proposed system. Tables 2 and 3 show the SNR before and after signal denoising by our proposed system. The experimental results show that our proposed system can effectively remove noise.

Finally, the recognition rates of our proposed system with other methods are listed in Tables 4 - 6. Table 4 shows the rate of recognizing the signal without denoising. Tables 5 and 6 show the rate of recognizing the signal after denoising by our

Table 2. The SNR of the signals in Aurora 2 database

SNR (without enhancing)									
	Subway	Babble	Car	Exhib.	Rest.	Street	Airport	Station	AVE.
20dB	18.5	18.3	18.4	18.6	18.38	18.32	18.45	18.38	18.41
15dB	13.5	13.6	13.4	13.6	13.45	8.91	8.96	13.55	12.37
10dB	8.85	8.9	8.95	8.92	4.97	4.94	4.96	8.93	7.42
5dB	4.95	4.9	4.98	4.96	4.97	4.95	4.92	4.94	4.94
0dB	0.84	0.62	0.08	0.12	0.11	0.3	0.26	0.09	0.30
-5dB	-4.81	-4.26	-4.33	-4.13	-5.01	-4.86	-4.97	-5.06	-4.67

Table 3. The SNR of the signals after enhancing by our proposed algorithm

SNR (after enhancing)									
	Subway	Babble	Car	Exhib.	Rest.	Street	Airport	Station	AVE.
20dB	20.1	18.9	20.3	19.7	18.8	19.2	19.4	18.7	19.38
15dB	15.5	14.8	13.4	14.9	13.9	14.2	13.6	14.4	14.33
10dB	11.4	10.3	11.6	11.3	9.2	10.6	9.8	10.94	10.64
5dB	7.4	6.1	8.1	6.6	5.7	6.7	6.4	6.5	6.68
0dB	4.45	3.68	5.3	4.3	3.16	3.77	3.63	3.53	3.97
-5dB	1.33	1.1	1.64	1.21	0.88	1.16	0.86	1.18	1.17

Table 4. The recognition rate of the signals without enhancing

Recognition Rate % (withouth enhancing)									
	Subway	Babble	Car	Exhib.	Rest.	Street	Airport	Station	AVE.
clean	98.8	98.6	98.8	99.1	99.3	99.25	99.22	99.5	99.07
20dB	97.6	97.8	98.5	97.4	98.4	97.8	98.2	98.95	98.08
15dB	95.1	94.2	95.7	94.9	95.4	94.35	95.3	95.7	95.08
10dB	84.9	81.4	80.5	84.6	85.6	82.15	83.2	82.6	83.11
5dB	63.75	59.4	49.1	56.1	63.3	54.25	60.8	54.9	57.7
0dB	35.13	35.06	22.73	25.15	37.21	28.46	35.8	26.89	30.80
-5dB	15.06	19.17	11.08	12.48	18.33	15.4	18.18	14.57	15.53

Table 5. The comparison of signal subspace with our proposed algorithm

		Recognition rate %								
	Method	Subway	Babble	Car	Exhib.	Rest.	Street	Airport	Station	AVE.
clean	subspace	99.26	99.06	99.14	99.51	99.26	99.24	99.14	99.51	99.26
	Our	99.10	98.90	99.16	99.31	99.36	99.26	99.28	99.57	99.24
20dB	subspace	99.28	97.40	98.09	97.81	97.97	97.67	97.55	97.55	97.79
	Our	98.10	97.96	98.56	97.57	98.6	97.93	98.26	98.84	98.22
15dB	subspace	95.52	92.87	96.69	94.48	94.32	93.23	94.18	94.18	94.43
	Our	95.65	94.20	96.10	95.43	95.24	94.85	95.6	95.91	95.37
10dB	subspace	87.69	81.89	90.3	85.34	84.07	83.92	84.46	84.46	85.26
	Our	87.90	83.60	82.9	87.8	88.3	84.95	84.9	86.7	85.88
5dB	subspace	74.70	64.06	74.59	62.94	65.18	67.02	67.52	67.52	67.94
	Our	71.23	65.40	72.3	67.1	66.8	64.25	69.82	63.9	67.6
0dB	subspace	52.16	39.84	45.45	38.29	41.23	39.02	43.66	43.66	42.91
	Our	53.33	40.21	44.33	39.1	42.3	41.36	45.1	43.89	43.70
-5dB	subspace	28.31	21.49	20.34	16.11	17.93	16.84	21.47	21.47	20.49
	Our	30.20	24.60	22.3	16.9	18.3	17.77	20.6	19.6	21.28

Table 6. The performance of contaminated signals denosing by spectrum subtraction

	Method	Subway	Babble	Car	Exhib.	Rest.	Street	Airport	Station	AVE.
20dB	spectrum	84.81	89.99	89.47	80.61	91.18	91.77	91.43	97.55	93.15
15dB	spectrum	77.22	79.11	81.76	67.78	79.34	85.96	84.81	94.18	90.56
10dB	spectrum	65.71	55.01	66.86	47.64	56.96	67.84	64.9	84.46	77.85
5dB	spectrum	45.52	27.27	38.72	25.09	29.52	43.5	36.08	67.52	47.53

proposed method, the signal subspace approach and spectrum subtraction method, respectively. Experimental results demonstrate that our proposed method outperforms other methods.

6 Discussions

In this paper a novel speech enhancement system controlled by fuzzy rules is proposed for removing additive noise in a contaminated speech signal. In this speech enhancement system the rules about how to set the system's parameters depending on the local characteristics of signals are expressed in an ambiguous style and the system is designed optimally by using partial swarm optimization, so that the mean square error of the system's output can be minimized. In the system four linguistic variables and seven fuzzy rules are used to infer the adaptive threshold. Different kinds of noises can be effectively represented and distinguished by the four linguistic variables. Experimental results demonstrate that our proposed method effectively removes noise and outperform other methods.

Concerning further researches. Searching for more critical parameters to identify the noise and finding more rules to add to this system is under considered. Although some studies have demonstrated algorithms can automatically generate inference rules, inference time increases as fuzzy rules increase. Finding the balance between precision and time complexity is an open issue.

References

- [1] Takahashi, Y., Takatani, T., Osako, K., Saruwatari, H., Shikano, K.: Blind spatial subtraction array for speech enhancement in noisy environment. *IEEE Trans. on Audio, Speech, and Language Processing* 17(4), 650–664 (2009)
- [2] Abramson, A., Cohen, I.: Simultaneous detection and estimation approach for speech enhancement. *IEEE Trans. on Audio, Speech and Language Processing* 15(8), 2348–2359 (2007)
- [3] Wu, S.-N., Wang, J.-S.: An adaptive recurrent neuro-fuzzy filter for noisy speech enhancement. In: *International Joint Conference on Neural Networks*, pp. 3083–3088 (2004)
- [4] Er, M.J., Li, Z.: Adaptive noise cancellation using online self-enhanced fuzzy filters with applications to multimedia processing. In: *Intelligent Multimedia Processing with Soft Computing*. Springer, Berlin (2006)
- [5] Ma, N., Bouchard, M., Goubran, R.A.: Speech enhancement using a masking threshold constrained kalman filter and its heuristic implementations. *IEEE Trans. on Audio, Speech, and Language Processing* 14(1), 19–32 (2006)
- [6] Shao, Y., Chang, C.-H.: A generalized time–frequency subtraction method for robust speech enhancement based on wavelet filter banks modeling of human auditory system. *IEEE Trans. on Systems, Man, and Cybernetic–Part B: Cybernetics* 37(4), 877–889 (2007)
- [7] Furuya, K., Kataoka, A.: Robust speech dereverberation using multichannel blind deconvolution with spectral subtraction. *IEEE Trans. on Audio, Speech and Language Processing* 15(5), 1579–1591 (2007)
- [8] Hao, J., Attias, H., Nagarajan, S., Lee, T.-W., Sejnowski, T.J.: Speech enhancement, gain, and noise spectrum adaptation using approximate bayesian estimation. *IEEE Trans. on Audio, Speech and Language Processing* 17(1), 24–37 (2009)
- [9] Chang, J.-H., Jo, Q.-H., Kim, D.-K., Kim, N.-S.: Global soft decision employing support vector machine for speech enhancement. *IEEE Signal Processing Letters* 16(1), 57–60 (2007)
- [10] Hasan, M.K., Salahuddin, S., Khan, M.R.: A modified a priori SNR for speech enhancement using spectral subtraction rules. *IEEE Signal Processing Letters* 11(4), 450–453 (2004)
- [11] Tahmasbi, R., Rezaei, S.: A soft voice activity detection using GARCH filter and variance gamma distribution. *IEEE Trans. on Audio, Speech, and Language Processing* 15(4), 1129–1134 (2007)
- [12] Davis, A., Nordholm, S., Togneri, R.: Statistical voice activity detection using low-variance spectrum estimation and an adaptive threshold. *IEEE Trans. on Audio, Speech, and Language Processing* 14(2), 412–424 (2006)
- [13] Johnson, M.T., Yuan, X., Ren, Y.: Speech signal enhancement through adaptive wavelet thresholding. *Speech Communication* 49, 123–133 (2007)
- [14] Kennedy, J., Eberhart, R.: Particle swarm optimization. In: *Proc. IEEE Int. Conf. Neural Network*, pp. 1942–1948 (1995)
- [15] Chen, S.-H., Wu, H.-T., Chang, Y., Truong, T.K.: Robust voice activity detection using perceptual wavelet-packet transform and Teager energy operator. *Pattern Recognition Letters* 28, 1327–1332 (2007)
- [16] Johnson, I.M., Silverman, B.W.: Wavelet threshold estimators for data with correlated noise. *J. Roy. Statist. Soc., Ser. B* 59, 319–351 (1997)
- [17] Hirsch, H.G., Pearce, D.: The aurora experimental framework for the performance evaluations of speech recognition systems under noisy conditions. In: *ISCA ITRW ASR 2000: Challenges for the New Millennium*, Paris, France (September 2000)
- [18] Speech Vision and Robotics Group, <http://htk.eng.cam.ac.uk>

Risk Prediction for Postoperative Morbidity of Endovascular Aneurysm Repair Using Ensemble Model

Nan-Chen Hsieh¹, Chien-Hui Chan², and Hsin-Che Tsai¹

¹ Department of Information Management, National Taipei College of Nursing
No. 365, Min-Ten Road 11257, Taipei, Taiwan, ROC
nchsieh@ntcn.edu.tw

² Department of Computer Science and Information Engineering, Tamkang University
No. 151, Ying-Chuan Road, 25137 Tamsui, Taipei County, Taiwan, ROC

Abstract. Endovascular aneurysm repair (EVAR) is an advanced minimally invasive surgical technology that is helpful for reducing patients' recovery time and postoperative morbidity. This study proposes an ensemble model to predict postoperative morbidity after EVAR. The ensemble model was developed using a training set of consecutive patients who underwent EVAR between 2000 and 2008. The research outcomes consisted of an ensemble model to predict postoperative morbidity, the occurrence of postoperative complications prospectively recorded, and the causal-effect decision rules. The probabilities of complication calculated by the model were compared to the actual occurrence of complications and a receiver operating characteristic (ROC) curve was used to evaluate the accuracy of postoperative morbidity prediction. In this series, the ensemble of BN, NN and SVM models offered satisfactory performance in predicting postoperative morbidity after EVAR. Moreover, the Markov blankets of BN allow a natural form of causal-effect feature selection, which provides a basis for screening decision rules generated by granular computing.

Keywords: Endovascular aneurysm repair (EVAR), postoperative morbidity, ensemble model, machine learning, Markov blanket

1 Introduction

Cardiac surgery is a complex surgical operation that is indicated for patients with severe insufficiency in cardiac function. Major cardiac surgical interventions include coronary artery bypass grafting (CABG), repair of congenital heart defects, surgical treatment of atrial fibrillation, heart transplantation, repair or replacement of heart valves, aortic surgery, aneurysm repair or a combination of these surgical procedures. During the operation and the postoperative stay at the ICU and nursing ward, there is considerable morbidity for cardiac surgery patients with postoperative complications, which results in increased hospital mortality and postoperative morbidity. Many prediction models for cardiac surgical outcome apply logistic or multivariable regression to assess preoperative risk [1-3]. EuroSCORE is a predominant model [4]. Most of the risk prediction models in current use were derived for patients undergoing open

abdominal aortic aneurysm (AAA) repair and appear to lack utility when applied to EVAR patients. EVAR is an advanced minimally invasive surgical technology that helps reduce patients' recovery time as well as postoperative mortality and morbidity; it is especially helpful in the treatment of patients judged to be high surgical risk for conventional surgery. EVAR benefits patients with medical co-morbidities, and post-operative complications highly significantly influence longer-term postoperative outcomes in EVAR patients.

The use of machine learning models has become widely accepted in medical applications. Various machine learning models including BNs, NNs, and SVMs have been tested in a wide variety of clinical and medical applications [5]. Soft-computing, including fuzzy set and rough set techniques that work well in descriptive data mining, is also a promising technique. BN is a probability-based inference model that has a wide range of applications and is increasingly used medically as a prediction and knowledge representation modeling technique. Verduijn, M., et al. [6, 7] presented the prognostic BN as a new type of prognostic model that builds on the BN methodology and implements a dynamic, process-oriented view of cardiac surgical prognosis. Lin and Haug [8] proposed BN suitable for exploiting missing clinical data for medical prediction. NNs have featured in a wide range of medical applications, often with promising results. Rowan et al. [9] used NNs to stratify the length of stay of cardiac patients based on preoperative and initial postoperative factors. Eom et al. [10] developed a classifier ensemble-based, including NNs, DTs, and SVMs, clinical decision support system for predicting cardiovascular disease level. SVMs have been successfully used in a wide variety of medical applications. Polat and Güne [11] used a least square support vector machine to assist breast cancer diagnosis. Babaoğlu et al. [12] first used principle component analysis method to reduce data features, and acquired an optimum support vector machine model for the diagnosis of coronary artery disease. Choi [13] proposed the detection of valvular heart disorder (VHD) by wavelet packet decomposition and SVM techniques.

Traditional medical data modeling tends to use statistical techniques to provide reports. However, most diagnoses and treatments in medical decision making are imprecise and are accompanied by errors. RS is able to circumvent this limitation with its ability to handle imprecise and uncertain data [14]. Theoretically, statistical clinical models use predetermined hypotheses and a priori assumptions in initial hypothesis test design. These models also usually require a perfectly random statistical distribution; furthermore, statistical formulae or equations usually prove poor in characterizing medical data and solving for the characteristics or relationships of that data. The RSs approach takes advantage of a correct proven philosophy to work with medical data without a strong a priori reasoning. As an example, Mitra et al. [15] implemented a rule-based RS decision system for the development of a disease inference engine for ECG classification. Podraza et al. [16] presented a complex data analysis and decision support system for medical staff based on rough set theory. Wakulicz-Deja and Paszek [17] implemented an application of RS theory to decision making for diagnosing mitochondrial encephalomyopathies in children.

This study describes the development of an informative ensemble prediction model consisting of BNs, NNs and SVMs and argument using an RS model for the prediction of postoperative morbidity and the description of causal effects between preoperative variables and complication outcomes in EVAR patients. The development of

this prediction model proceeds as follows. First, an informative discretization method for numerical values, which employs a fuzzy discretization method to guide the categorized numerical features as linguistic fuzzy terms on the basis of entropy/MDL algorithm and laboratory surgeon’s knowledge, is developed. Through fuzzy discretization, numerical values are converted into discrete values with informative fuzzy trapezoid functions. Second, the proposed ensemble-based architecture focuses on fusing three types of models (BNs, NNs, and SVMs) with the RS used to augment the interpretability of the derived knowledge. The model selection scheme, designated a CV-bagging scheme, is a mixture of bagging and cross-validation that is chosen in order to improve overall classification by combining classifiers trained on randomly generated subsets of the entire training set. This strategy increases both the diversity and accuracy of individual ensemble classifiers. Finally, the Markov blanket (MB) concept of BN allows a natural form of feature selection, providing a basis concept for filtering decision rules by RSs. The learned knowledge is presented in multiple forms, including causal-effect diagrams and decision rules.

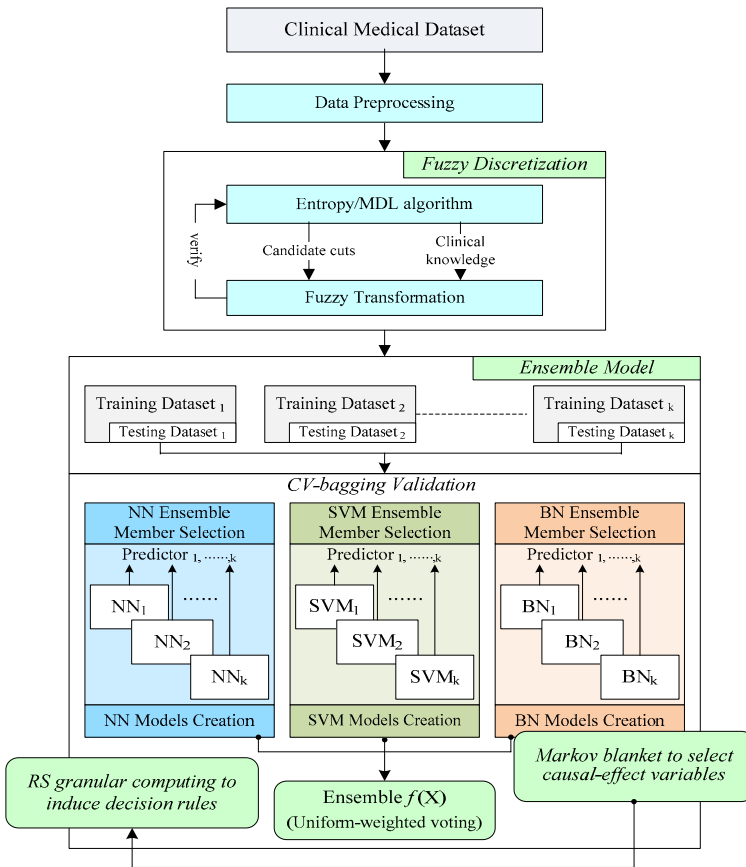


Fig. 1. The proposed architecture

Section 2 of this paper begins with an overview of study background and research materials in general. Section 3 describes the methods and procedures used in this study. Section 3 also empirically tests the proposed method using an EVAR clinical dataset. Section 4 discusses the experimental findings and offers observations about practical applications and directions for future research.

2 Study Background and Materials

Abdominal aneurysm (AAA) is an enlargement that occurs in a weakened area within the largest artery in the abdomen (http://www.vascularweb.org/patients/NorthPoint/Abdominal_Aortic_Aneurysm.html). If an AAA is not treated in due time, the pressure generated by heartbeats causes the aneurysm to continuously grow larger, and the aortic wall continues to weaken. Finally, rupture occurs and massive internal bleeding occurs. The best way to prevent the high mortality associated with AAA is to find the lesion before rupture occurs. However, patients with aortic diseases are often elderly with severe co-morbidities and sometimes devastating morbidity, making them extremely challenging candidates for surgery. For such patients, EVAR represents a lower risk approach than conventional open surgery and is associated with shorter operating times, shorter hospitalizations, more rapid recovery and improved quality of life during the perioperative period and postoperative follow-up. Although long-term data on the clinical outcomes of patients who received EVAR are not yet available, given its importance, building a prediction of postoperative morbidity after EVAR is critical.

We retrospectively examined 140 consecutive patients who underwent EVAR surgery at Taipei Veteran General Hospital, a teaching hospital in Taiwan, between 2000 and 2008. The dataset contains preoperative patient characteristics, details of the operative information, and pathological and laboratory findings from the emergency room, operating room and intensive care unit (ICU). The dataset also included length of ICU stay, variables that describe postoperative complications that frequently occur in EVAR surgery, death during hospitalization, and time of death for patients who expired. Postoperative complication was used as the binary outcome variable of the ensemble model, and types of complications were used as subsidiary outcome variables. The original dataset contained 137 variables, but included many missing values. Preliminary inspection of the dataset showed that many variables contained missing values for at least 50% of the patients; these variables were not included in further analysis. In order to identify significant variables for use in the ensemble model, a number of criteria were employed. Variables that were subjective, ambiguous or inadequately defined were excluded; variables that were frequently incomplete were also excluded from subsequent analysis. Sixty-seven of the 140 patients experienced severe postoperative complications during their stay at hospital. Data collected included preoperative patient characteristics, risk factors, details of the operative information, physical characteristics of the aneurysm, postoperative physiological and laboratory findings, and postoperative complications as the outcome variable.

We employed RS as a preprocessing tool for data reduction. The main advantage of RS is that it does not require any preliminary or additional information about data

such as probability distribution in statistics, basic probability assignment in the Dempster-Shafer theory, or membership degree in fuzzy set theory; RS is especially suitable for dealing with sparse datasets. Based on the indiscernibility relation of RS, redundant variables were identified and eliminated to reduce the number of variables. Originally, RS with Johnson’s search algorithm reported nine variables that showed significant relatedness to the outcome variable. Following suggestions from surgeons and review of relevant research studies [1-3], 10 more interesting variables were included in each patient’s profile. As shown in Table 1, the final dataset contains nine prediction variables, 10 profiling variables and one outcome variable. Besides its usefulness for data reduction, RS is also suitable for finding dependencies among variables, discovering cause-effect relationships, and generating decision rules. Accordingly, we also used RS as the preliminary variable selection mechanism for building MB of BN.

Table 1. Variables used to predict the postoperative morbidity

Variable	Type	Definition	Category	Use as
Heart disease	Categorical	Heart diseases, A ₀ : no, A ₁ : yes	Patc.	Prediction
Smoking	Categorical	A ₀ : no smoking, A ₁ : smoking, A ₂ : uncertain	Patc.	Prediction
Hypertension	Categorical	A ₀ : no, A ₁ : yes	Patc.	Prediction
Hyperlipidemia	Categorical	Hyperlipidemia, A ₀ : no, A ₁ : yes	Patc.	Prediction
CRI	Categorical	Chronic renal insufficiency, A ₀ : no, A ₁ : yes	Patc.	Prediction
COPD	Categorical	Chronic obstructive pulmonary disease, A ₀ : no, A ₁ : yes	Patc.	Prediction
Seg	Numerical	Segmented neutrophil count, 0-100%	Lab.	Prediction
AAA Size	Numerical	Size of abdominal aortic aneurysm, mm	Pato.	Prediction
AAA Shape	Categorical	Shape of abdominal aortic aneurysm, A ₀ : Fusiform, A ₁ : Saccular, A ₂ : Rupture, A ₁₂ : Saccular and Rupture	Pato.	Prediction
BUN	Numerical	Blood urea nitrogen count, mg/dl	Lab.	Profile
Creatine	Numerical	Creatine count, mg/dl	Lab.	Profile
HCT	Numerical	Hematocrit count, 0-100%	Lab.	Profile
Hb	Numerical	Hemoglobin count, g/dl	Lab.	Profile
PLT	Numerical	Platelet count, /CUMM	Lab.	Profile
Age	Numerical	Integer number, 0- 100 years	Patc.	Profile
Gender	Categorical	Male/female	Patc.	Profile
DM	Categorical	Diabetes, A ₀ : no, A ₁ : yes	Patc.	Profile
CVA	Categorical	Cerebral vascular accident, A ₀ : no, A ₁ : yes	Patc.	Profile
AAA site	Categorical	Site of abdominal aortic aneurysm, A ₀ : Infra renal, A ₁ : Supra renal, A ₂ : Juxta renal	Pato.	Profile
Complications	Categorical	Postoperative complications, yes/no		Outcome

i.e. Patc.: patient characteristics; Pato.: pathological; Lab.: laboratory findings.

3 Methods and Procedures

3.1 Entropy/MDL-Based Method for the Discretization of Numerical Values

Most studies dealing with cardiac surgery prediction models have applied logistic or multivariable regression to assess the preoperative risk. Few studies have utilized machine learning algorithms i.e., decision trees, Bayesian networks or neural networks in analyzing clinical data. These state-of-the-art machine learning algorithms are often informative and can represent more knowledge in the clinical data. Generally, clinical datasets involve numerical variables. To satisfy the requirements of prediction models, the employment of a discretization approach is necessary. Discretization is defined as a process that divides numerical values into states of discrete categorical values, leading to informative expressed categorical values. For example, the CART model originally was not designed to handle numerical attributes. During the construction of a CART model, numerical attributes were divided into discrete categorical values. Aside from CART, discretization techniques were frequently adopted in other popular learning paradigms, such as C4.5, BNs, NNs, and genetic algorithms.

This study employed entropy/MDL to discretize numerical domains into intervals. Incorporation investigates fuzzy discretization; the discretized value is characterized by a triple of interval (i,m,a) , where i denotes interval name, m denotes membership value of the value, and a denotes affinity of the value. The steps for automatically finding fuzzy sets from a given dataset are described herein. Assuming that the domain of a numerical attribute ranges from v_1 to v_2 , and that $\{c_1, c_2, \dots, c_k\}$ denote the k cut-points obtained by the entropy/MDL algorithm, using these k midpoints, $k/2 + 1$ linguistic terms or membership functions can be determined for a trapezoidal fuzzy set. The first membership function is computed as:

$$f_{first}(x) = \begin{cases} 1 & \text{if } v_1 \leq x \leq c_1 \\ (c_2 - x)/(c_2 - c_1) & \text{if } c_1 < x < c_2 \\ 0 & \text{if } x \geq c_2 \end{cases}$$

Generally, the middle p th, $p = 2, \dots, k/2$, membership function is computed as:

$$f_p(x) = \begin{cases} 0 & \text{if } x \leq c_{2p-3} \\ (x - c_{2p-3})/(c_{2p-2} - c_{2p-3}) & \text{if } c_{2p-3} < x < c_{2p-2} \\ 1 & \text{if } c_{2p-2} \leq x \leq c_{2p-1} \\ (c_{2p} - x)/(c_{2p} - c_{2p-1}) & \text{if } c_{2p-1} < x < c_{2p} \\ 0 & \text{if } x \geq c_{2p} \end{cases}$$

The final membership function is computed as:

$$f_{final}(x) = \begin{cases} 0 & \text{if } x \leq c_{k-1} \\ (c_k - x)/(c_k - c_{k-1}) & \text{if } c_{k-1} < x < c_k \\ 1 & \text{if } c_k \leq x \leq v_2 \end{cases}$$

For example, AAA_Size is a numerical attribute obtained from the dataset. The domain of AAA_Size ranges from 5.0 to 9.6. Six cut-points, (5.5, 5.7, 5.8, 7.3, 8.1, 9.6), were obtained using the entropy/MDL algorithm. As depicted in Fig. 2, a trapezoidal fuzzy set was obtained with four membership functions corresponding to the six intervals. If AAA_Size={5.65cm}, then it can be transformed to a triple value $[i_2, \max(\text{round}(0.75, 0.8), \text{round}(0.3, 0.4)), i_3] = [i_2, 0.8, i_3]$, meaning that this value belongs to interval i_2 with membership value 0.8, and adheres to interval i_3 . If AAA_Size = {6.55cm}, it can be transformed to a triple value $[i_4, 0.5, i_4]$. Within each analysis, the laboratory system should contain an interval that delimits life-compatible values. Therefore, if cut points fall within the interval, they should be eliminated. For example, one cut point {1.4 mg/dl} of creatine lay within an interval, i.e., [0.75~1.5] for male, [0.5~1.2] for female; therefore the cut point {1.4 mg/dl} was eliminated.

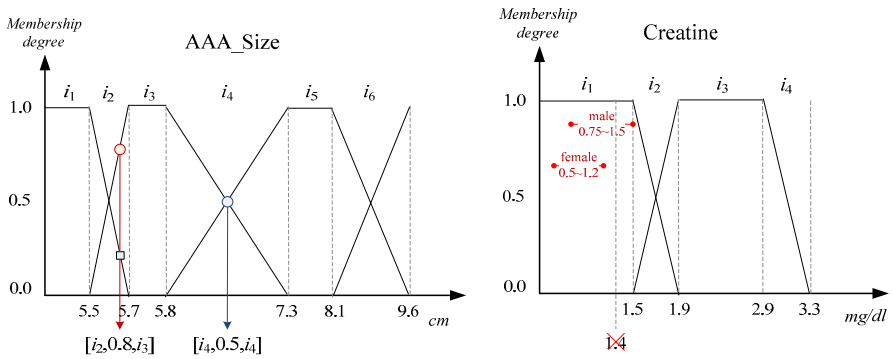


Fig. 2. The trapezoidal fuzzy set of AAA_Size and Creatine

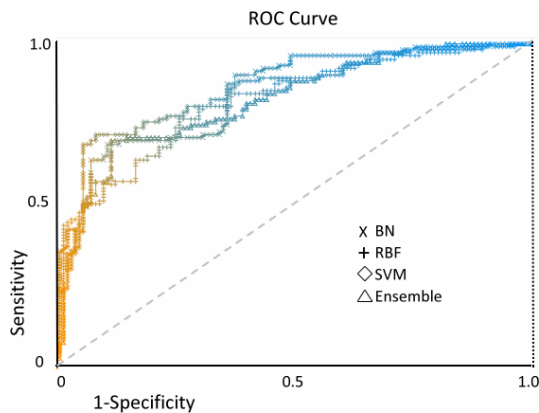
3.2 Ensemble Model for the Prediction of Postoperative Morbidity

BNs, NNs, SVMs and RSs were chosen because they represented different approaches, each of which is simple to implement and has been shown to perform well in medical applications. The rationale of employing these models is that BNs can easily model complex relationships among variables, NNs are generally superior to conventional models, SVMs can be used as a benchmark technique, and RSs can be used to augment the interpretability of the derived knowledge. The detailed configurations of each individual model are as follows: BNs are Bayesian networks with K₂ search algorithm, NNs are RBF neural network with radial basis functions as activation functions, SVMs use John Platt's sequential minimal optimization algorithm with logistic to the outputs for proper probability estimates, and RS is Pawlak's model with Johnson search algorithm for generating the singleton reduce.

In this study, the model selection scheme is a mixture of bagging and cross-validation (CV-bagging) that aims to improve the classification by combining models trained on randomly generated subsets of the entire training set. We first applied a cross validation scheme for model selection on each subset; subsequently, for the sake of simplicity and so as not to run into over-fitting problems, we combined the selected models in a uniform weights approach to reach a final decision. The concept of

uniform weights voting is both appealing and simple to implement; it can generally be applied to any type of classifier without relying on specific interpretations of the output. In order to validate models, we used random subsets as cross-validation folds for reasons of simplicity. In k-fold cross validation, the dataset is partitioned into k subsets. Of these k subsets, a single subset is retained as the validation dataset, and the remaining k-1 subsets are used for training datasets. The cross validation process is then repeated k times with each of the k subsets used only once as the validation data. The k results from the folds can then be averaged to produce a single estimation. In cases where specific inside medical knowledge is not available, such a cross validation method can be used to select a classification method empirically, because it seems to be obvious that no classification method is uniformly superior. For each model, we train several individual classifiers and select the best performing model as final model. If the trained models of each model differently with different initial conditions, we then find proper values for the free parameters of the model. We extend these results by combining models from different classes in order to increase the model diversity. The result, a heterogeneous ensemble, allows classification methods to be used more effectively.

The area under the ROC curve is based on a non-parametric statistical sign test, the purpose of which is to estimate the probability of survival for each of a pair of patients, one of whom did not survive. The surviving patient gains greater probability of survival. In this study, the area under the ROC metric was assessed through stratified 10-fold CV-bagging. The detailed accuracy of individual models and of the ensemble



Models	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area
10 folds CV- bagging weighted average						
BN	0.736	0.270	0.737	0.736	0.734	0.834
NN	0.675	0.330	0.657	0.657	0.674	0.755
SVM	0.757	0.247	0.758	0.757	0.756	0.755
Ensemble	0.757	0.249	0.760	0.757	0.755	0.813

Fig. 3. The results of the postoperative morbidity prediction

model is shown in Figure 3. The results of the experiment show that individual models' performance and improvements in performance were achieved by applying the ensemble of models. This indicates that model combination techniques indeed yield stable performance improvements for ensemble models.

3.3 BN's Markov Blanket for the Causal Interpretability of the Prediction Model

BN models have the ability to reason under conditions of uncertainty on the basis of probabilistic inference mechanisms; therefore, they exhibit more resistance than other machine learning models to noise that commonly occurs in medical datasets. That is, if a dataset contains some samples with missing values, it can still be used to train or test BN models. Because medical datasets tend to have missing values, this property is useful for intelligent data analysis. Unlike most regression techniques, BN models can easily model complex relationships among variables without any requirements for the underlying distributions of variables and are more transparent and intuitive because relationships among variables are explicitly represented by a directed acyclic graph (DAG) for modeling experts' knowledge. The DAG consists of a set of nodes linked by directed arcs, where the nodes represent the variables and the arcs indicate the direct causal relationships among the connected attributes. The strengths of these relationships are expressed by conditional probabilities. By using BNs, we can find a given network's structure and its associated parameters; the resulting model can provide diagnostic, predictive, and inter-causal reasoning. Therefore, the DAG of BN can serve as a causal analysis tool to help surgeons arrive at valuable decision rules generated by RS. Through the use of causal-effect analysis, medical causal knowledge can be exhibited through a causal DAG graph. Such knowledge can help surgeons stimulate the hidden knowledge behind the medical dataset, communicate complex causal concepts, and transform an individual's tacit knowledge into a team's explicit knowledge.

A novel idea of BN for significant feature selection is the Markov blanket (MB) [18], which is defined as the set of input features such that all other features are probabilistically independent of the target features. In this case, based on a general BN classifier, we can obtain a set of features that are present in the MB of the target feature, and features provided by the MB are sufficient for perfectly estimating the distribution of and for classifying the target feature classes. As an example, suppose that we build a BN classifier with a complete input training data set. The MB of the target feature then forms the feature selection mechanism and all features outside the MB are ignored by the BN. The remaining features are then used to induce constrained decision rules. In this study, RS is used to augment the interpretability of ensemble model; all decision rules derived by RS are in the form, $Y \leftarrow X$, where Y is the target variable and X is a subset of variables from the MB. The rationale of this method is that for any feature $x_i \notin \{\text{MB} \cup x_j\}$, $x_i \neq x_j$ and $x_i, x_j \in \text{MB}$, we have $P(x_i | \text{MB}, x_j) = P(x_i | \text{MB})$. That is, the MB of a variable x_i is a minimal set of features MB such that x_i is conditionally independent of any other feature x_j in MB. The MB structure to be used for BN explores not only the relationship between target and predictive variables, but also the relationships among the predictive variables themselves.

Since the conditional independence assumption of BN is not real, BN with a TAN search algorithm can therefore be used to offset this disadvantage. The use of this approach results in significant improvement in terms of classification accuracy, efficiency and model simplicity. Although use of the TAN search algorithm may not always result in the best classification accuracy, this study adopts the TAN search algorithm because it can create a causal effect graph in which the postoperative-complications target variable treated as the supreme parent node is located at the top in the DAG. After a BN is learned, a heuristic is applied to ensure that all nodes in the BN are part of the Markov blanket of the target variable.

Figure 4 depicts an MB for the EVAR dataset composed of 10 variables, 9 predictive variables, and one target variable. The 9 predictive variables contributed significantly to identification of all the variables in the BN that are needed to classify the target variable. We should use these variables as RS's input variables to generate decision rules, connecting with proper profile variables for interpretation. The structure of the MB proposed for analyzing medical decision tasks follow s the logical structure for features defined in a loss causation model.

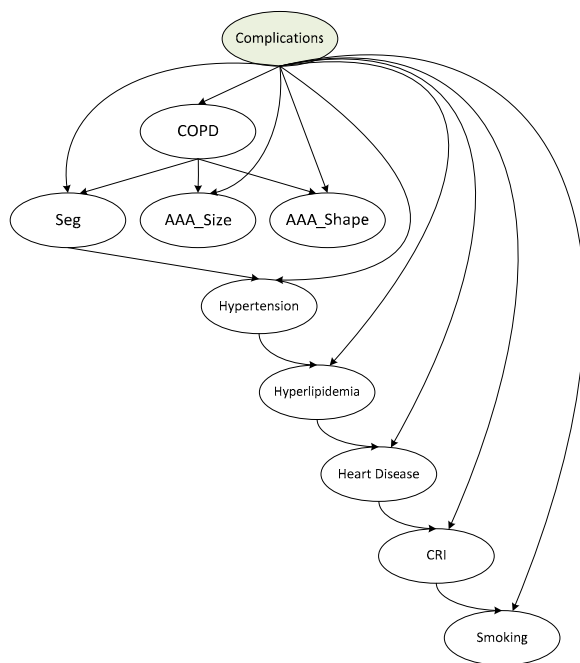


Fig. 4. A Markov blanket for EVAR

3.4 Granular Computing for the Extraction of Causal-Effect Knowledge

In recent years, RSs and granular computing have been proven to be good soft-computing techniques when used in various synergistic combinations with fuzzy logic, neural networks and genetic algorithms. Since RSs do not require prior parameter settings and extract knowledge from the data itself by means of indiscernibility

relations, RSs provide a stronger framework to achieve tractability, robustness, and interpretable decision-making knowledge. RSs require fewer calculations than do other soft-computing techniques, and provide an effective means for analysis of data by synthesizing upper and lower approximations of set concepts from the acquired dataset [19]. The induced decision rules are concise and valuable and can benefit surgeons by revealing hidden medical knowledge in the dataset, even when evidence is sparse. Granular computing refers to computation performed on information granules. Therefore, RS leads to both data compression and gains in computation time, thereby finding wide application.

As an example, rough-fuzzy integration [20] can be considered as a way of emulating the basis of f -granulation in data summarization, where quantifiers have fuzzy boundaries and granular attribute values. Since granular computing is performed on groupings of similar objects, rather than on individual data values, computation time is significantly reduced. Due to its comprehensive logical rules, RS can invariably generate a high performance prediction model, and each logical rule represents an information granule. However, unimportant logical rules obstruct the utility of the induced knowledge. In this study, we proposed to use an RS-BN integration approach to induce necessary decision rules derived from medical knowledge, to generate a prediction model for postoperative morbidity, and to find significant MBs as causal-effect relationships among the variables. MBs from BN are then used as a basis for filtering significant variables for granular computing of RS. The obtained decision rules correspond to different important regions of the feature space or data groups, a feature that is good for interpretation. Due to the suitability of BNs and RSs for handling medical data, this study adopted MBs of BN as the candidate variables selector and used RSs to augment the interpretability of classifiers. The direct information provided by MB on the basis of the probability distribution obtained for each variable, shown in Fig. 4, represents data representing the most typical circumstances arising among the EVAR cases. However, explanation of the meaning of the probability distribution is necessarily rather complex. For easy interpretation, the profiles for postoperative morbidity were established with the help of the RS. We can certainly extract decision rules with whole or partial causality. From this, surgeons and decision makers can create more valuable discrimination rules regarding the possibilities of complications and postoperative morbidity.

Figure 5 depicts a partial causality of BN with probability distributions. We used variables in MB as RS's input variables to generate decision rules, connecting with proper profile variables for interpretation. The structure of the MB proposed for analyzing medical decision tasks follows the logical structure for features defined in a loss causation model. In this network, we can trace how each node will influence and be influenced by other nodes. For example, COPD and Complications will directly influence the state probability of Seg, AAA_size, and AAA_shape through a top-down causal-effect. The state probabilities of the input and profile variables are also included. From the surgeon's perspective, the constrained association rules are reasonable for outcome description. Suitable linguistic quantifiers can be employed to add linguistic confidence.

In most circumstances, medical datasets are uncertain and incomplete and cases are sparse. By conventional rules, judgment standards such as support and confidence are not properly used to select valuable rules. The fuzzy logic-based calculus [21]

provides interpretation and validation of truth statements involving complex linguistic quantifiers such as “always”, “medium” and “often”. The truth statement evaluation is therefore used to validate and select decision rules. The procedure for determining the truth value of a linguistically quantified statement is as follows. Let “Q {t₁,...,t_n} are S” denotes a linguistically quantified statement, and let {t₁,...,t_n} denotes a set of decision rules, DS. If the summary S involves a variable A, and t_i denotes a rule that satisfies the summary S, then the membership value of t_i to S is given by:

$$S(t_i) = \max_{\forall k} (\mu_{EQ}(a_k, b_k)), \text{ for all } a_k \in t_i[A], b_k \in S,$$

where $S(t_i)$ denotes the degree to which t_i satisfies the summary S , and the function $\mu_{EQ}(a_k, b_k) = 0$ if and only if $(a_k \neq b_k) \vee (a_k = b_k \wedge \mu(b_k) = 0)$; $\mu_{EQ}(a_k, b_k) = 1 - |\mu(a_k) - \mu(b_k)|$ if and only if $(a_k = b_k \wedge \mu(b_k) \neq 0)$. Then, the individual truth value of “{t₁,...,t_n} are S” for a variable A over DS is computed as:

$$Truth(\{t_1, \dots, t_n\} \text{ are } S) = \left(\frac{1}{n} \sum_{i=1}^n S(t_i) \right)$$

Moreover, when the linguistic summaries are distributed over m variables with ANDed conditions, that is, $S = S_1 \wedge \dots \wedge S_m$, then the total truth value $Truth(\{t_1, \dots, t_n\} \text{ are } S) = \min_{j=1}^m (Truth(\{t_1, \dots, t_n\} \text{ are } S_j))$. Finally, $T=Q(Truth(\{t_1, \dots, t_n\} \text{ are } S))$ denotes the truth value of the linguistically quantified statement “Q {t₁,...,t_n} are S” to the fuzzy quantifier Q in agreement. We gave an example for the procedure in obtaining truth value using linguistic summaries. For example, let COPD be the observation target; then {COPD, Seg, AAA_Size, AAA_Shape, Complications} is a MB of partial causality. By using granular computing of RS, several decision rules were obtained; we list three relevant decision rules as follows, where “support” means number of cases supporting this rule:

IF	THEN	SUPPORT
COPD(A ₀) and Seg(i ₂ ,0.9,i ₃) and AAA_size(i ₁ ,1.0,i ₂) and AAA_Shape(A ₀)	Complications(Y)	7
COPD(A ₀) and Seg(i ₁ ,0.7,i ₂) and AAA_size(i ₁ ,0.8,i ₂) and AAA_Shape(A ₀)	Complications(Y)	3
COPD(A ₀) and Seg(i ₂ ,0.8,i ₃) and AAA_size(i ₄ ,0.6,i ₃) and AAA_Shape(A ₀)	Complications(Y)	1

Next, we defined the fuzzy quantifiers, Q , as shown in Fig. 6, which are represented as 7 membership functions. Then linguistically quantified statement for validating these decision rules is therefore given by:

“ Q of the patients of postoperative morbidity have

COPD(A₀) and Seg(i₂) and AAA_size(i₁) and AAA_Shape(A₀).”

method of computing the individual truth value for the variable Seg is described as follows. Let the linguistic summary, $S_1(\text{Seg}(i_2))=(i_2,1.0,i_2)$ and a decision rule component $t_i[\text{Seg}]=(i_2,0.9,i_3)$. Then, the membership value of t_i to S_1 of Seg is given by

$S_1(t_1[Seg])=\max(1,0.9)$ if and only if they have overlap interval or $S_1(t_1)=0$. Similarly, $S_1(t_2[Seg])=\max(1,0.7)$ and $S_1(t_3[Seg])=\max(1,0.8)$. For n tuple $\{t_1, \dots, t_n\}$, the individual truth value “ $Truth(\{t_1, \dots, t_n\})$ of Seg to S_1 ” over DS is computed as $\left(\frac{1}{n} \sum_{i=1}^n S(t_i)\right)=1$, where $n=11$. Similarly, for linguistic summary, $S_2(AAA_Size(i_1))$, the individual truth value “ $Truth(\{t_1, \dots, t_n\})$ of AAA_Size to S_2 ” over DS is 0.85. Moreover, COPD and AAA_Shape are categorical values, and data values matched the values in linguistic quantified statement, therefore, their individual truth values are 1. Finally, the linguistic quantified statement is *ANDed* with linguistic summaries, the total truth value for this linguistic quantifier statement is $\min(1,1,0.85,1)=0.85$. A decision rule with truth is obtained in the following form, which denotes a satisfactory truth degree:

“IF COPD(A_0) and Seg(i_2) and AAA_size(i_1) and AAA_Shape(A_0) THEN Patients have postoperative morbidity, with *almost_always* truth 0.95 and support 35%.”

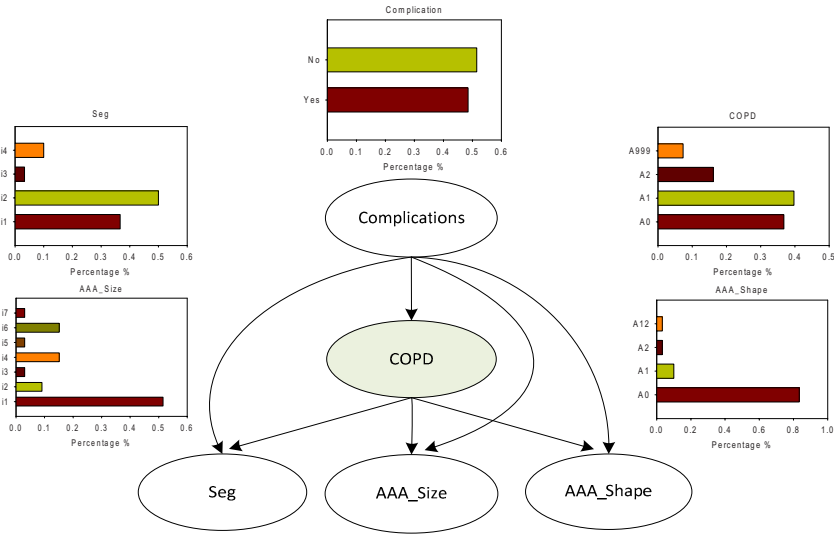


Fig. 5. The partial causality of MB with probability distributions

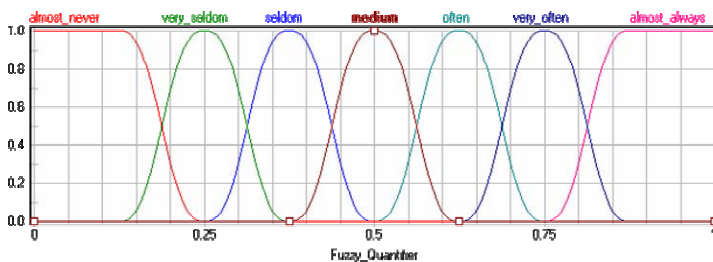


Fig. 6. Fuzzy quantifiers for linguistic summaries

4 Conclusions

Although several risk prediction systems have been proposed for patients undergoing open aneurysm repair (Bohm et al., 2008), they basically rely on traditional statistical methods and provide scant accuracy and utility when applied to EVAR patients. We have proposed an ensemble model to predict postoperative morbidity after EVAR and support clinical decision-making. The proposed ensemble model is constructed by incorporating fuzzy discretization of categorical numerical values; BNs, NNs, and SVMs were used to augment the ensemble model and the dataset was processed by CV-bagging, showing moderate performance. The experimental result shows that the proposed ensemble model predicts postoperative morbidity with relatively satisfactory accuracy, even when data is missing and/or sparse, showing its usefulness in support of clinical decision-making. In particular, the Markov blankets allow for a natural form of causal-effect feature selection, which itself provides a basis for screening decision rules generated by granular computing and generates supplementary information for risk prediction. The learned knowledge is represented in multiple forms, including causal-effect diagrams and causal-effect decision rules. The supplementary nature of multi-models distinguish the proposed model from existing risk scoring systems that are based on conventional statistical methods and from various machine learning models. To summarize, the advantage of using the proposed ensemble model is that it can provide surgeons with practical, relatively accurate aid in their daily diagnostic tasks and enable them to extract meaningful relationships among features of medical datasets through the use of causality graphs and constrained decision rules.

Acknowledgement

This research was partially supported by National Science Council of Taiwan (NSC 97-2410-H-227-002-MY2).

References

1. Barnes, M., Boulton, M., Maddern, G., Fritridge, R.: A model to predict outcomes for Endovascular aneurysm repair using preoperative variables. *European Journal of Vascular and Endovascular Surgery* 35, 571–579 (2008)

2. Bohm, N., Wales, L., Dunckley, M., Morgan, R., Loftus, I., Thompson, M.: Objective risk-scoring systems for repair of abdominal aortic aneurysms: applicability in Endovascular repair. *European Journal of Vascular and Endovascular Surgery* 36, 172–177 (2008)
3. Stijn, C.W., Wouters, L.N., Freek, W.A.V., Rene, M.H.J.B.: Preoperative prediction of early mortality and morbidity in coronary bypass surgery. *Cardiovascular Surgery* 10, 500–505 (2002)
4. Roques, F., Michel, P., Goldstone, A.R., Nashef, S.A.: The logistic EuroSCORE. *European Heart Journal* 24, 1–2 (2003)
5. Bellazzi, R., Zupan, B.: Predictive data mining in clinical medicine: Current issues and guidelines. *International Journal of Medical Informatics* 77, 81–97 (2008)
6. Verduijn, M., Peek, N., Rosseel, P.M., de Jonge, E., de Mol, B.A.J.M.: Prognostic Bayesian networks I: Rationale, learning procedure, and clinical use. *Journal of Biomedical Informatics* 40, 609–618 (2007)
7. Verduijn, M., Peek, N., Rosseel, P.M., de Jonge, E., de Mol, B.A.J.M.: Prognostic Bayesian networks II: An application in the domain of cardiac surgery. *Journal of Biomedical Informatics* 40, 619–630 (2007)
8. Lin, J.-H., Haug, J.: Exploiting missing clinical data in Bayesian network modeling for predicting medical problems. *Journal of Biomedical Informatics* 41, 1–14 (2008)
9. Rowan, M., Ryan, T., Hegarty, F., O'Hare, N.: The use of artificial neural networks to stratify the length of stay of cardiac patients based on preoperative and initial postoperative factors. *Artificial Intelligence in Medicine* 40, 211–221 (2007)
10. Eom, J.-H., Kim, S.-C., Zhang, B.-T.: AptaCDSS-E: A classifier ensemble-based clinical decision support system for cardiovascular disease level prediction. *Expert Systems with Applications* 34, 2465–2479 (2008)
11. Polat, K., Güne, S.: Breast cancer diagnosis using least square support vector machine. *Digital Signal Processing* 17, 694–701 (2007)
12. Babaoğlu, I., Fındık, O., Bayrak, M.: Effects of principle component analysis on assessment of coronary artery diseases using support vector machine. *Expert Systems with Applications* (in Press)
13. Choi, S.: Detection of valvular heart disorders using wavelet packet decomposition and support vector machine. *Expert Systems with Applications* 35, 1679–1687 (2008)
14. Pattaraintakorn, P., Cercone, N.: Integrating rough set theory and medical applications. *Applied Mathematics Letters* 21, 400–403 (2008)
15. Mitra, S., Mitra, M., Chaudhuri, B.B.: A rough-set-based inference engine for ECG classification. *IEEE Transactions on Instrumentation and Measurement* 55, 2198–2206 (2006)
16. Podraza, R., Dominik, A., Walkiewicz, M.: Decision support system for medical applications. In: Hamza, M.H. (ed.) *Applied Simulation and Modelling*, Marbella, Spain (2003)
17. Wakulicz-Deja, A., Paszek, P.: Applying rough set theory to multi stage medical diagnosing. *Fundamenta Informaticae* 54, 387–408 (2003)
18. Cheng, J., Greiner, R., Kelly, J., Bell, D.A., Liu, W.: Learning Bayesian networks from data: an information-theory based approach. *The Artificial Intelligence Journal* 137, 43–90 (2002)
19. Pawlak, Z.: Rough sets. *International Journal of Computer and Information Sciences* 11, 341–356 (1982)
20. Pal, S.K.: Soft data mining, computational theory of perceptions, and rough-fuzzy approach. *Information Sciences* 163, 5–12 (2004)
21. Yager, R.R.: Database discovery using fuzzy sets. *International Journal of Intelligent Systems* 11, 691–712 (1996)

Further Results on Swarms Solving Graph Coloring

Manuel Graña, Blanca Cases, Carmen Hernandez, and Alicia D'Anjou

Grupo de Inteligencia Computacional www.ehu.es/ccwintco UPV/EHU

Abstract. We have proposed the mapping of graph coloring problems into swarm dynamics. Empirical evidence that flock steering behaviors augmented with the notion of hostility (enmity and friendliness) are enough to perform efficiently the task of coloring the nodes of graphs even in the case 3-coloration hard graph topologies. We discuss here what are the minimal cognitive capabilities that allow the emergent behavior of swarms to solve such NP-complete problem without mediating an explicit knowledge representation.

1 Introduction

Swarm Intelligent Systems are models of spatial behavior of populations whose individuals (boids) are steered by decisions taken in a neighboring zone around them. The models of flocking boids developed by Craig Reynolds [9,10] representing the steering behaviors of birds that follow each other following a linear combination of three local rules: separation, coherence and alignment to other flockmates or to static goals situated in the environment. This model has been integrated in a more general models of Self-Organizing Particle Systems (SOPS). Our focus in this work is on the emergent behavior arising when boids have the cognitive ability to perceive the hostility quality of spatially close boids, which are perceived as friend or enemies. This abilities determine the emergence of a spatial distribution at a global level minimizing the number of cohabitating enemies. When the definition of enemies and friends is determined by an underlying graph, this behavior produces a spatial distribution that solves precisely and efficiently the the Graph Coloring Problem (GCP) [2]. The spatial distribution follows a kind of supervised behavior in the sense that we define some regions of the space with singular properties (i.e. representing a color), and the boids are forced to converge to them.

The Supervised Territorial Flock Coloring (SFTC) when taking to its limit expression, that is with flocking parameters of coherence, separation and alignment set to zero, is able to color graphs efficiently. Surprisingly, it improves with hard graphs for 3-coloration. That means that territoriality is a powerful driving force by itself. We identify a key parameter that seems to characterize the problem because it provides a fine prediction of the results that the algorithm will obtain. We have called Hostility this parameter, which is the ratio of the average number of friends to the average number enemies per node in the graph.

There are some other efforts to find other swarm inspired ways to solve the GCP. A k -coloration of a undirected graph with vertex set V and E the set of edges $G = (V, E)$ is a partition of the graph in k classes of equivalence such that the number of connected nodes belonging to the same class is minimized [6]. A recent approach by ACO [5,4] to solve the GCP is [7] identifying ants to colors. In [2] we have shown that our approach improved over some greedy algorithms over a collection of hard coloring problems. Here we are more focused in identifying the minimal swarm dynamics that allow to solve the coloring problem.

Section 2 gives some basic definitions of swarms for the GCP. Section 3 introduces the Supervised Territorial Flock Coloring (SFTC). Section 4 reports the experiments with the extreme version of the SFTC. Finally we gather some conclusions and directions for further work in section 5.

2 Basic Definitions of Swarms for GCP

The aim of this paper is to prove that GCP solutions emerge from the swarm dynamics of flocking boids that have hostility as the basic cognitive ability. They are able to recognize and run away from enemies, and run towards friends. Hostility is defined in terms of the underlying graph to be colored. Boids correspond to nodes in the graph. Enemies of a boid are the boids whose corresponding nodes are connected though an edge in the graph. That is, enemies are nodes at distance 1 over the graph. The friends are boids whose nodes are a distance 2 or greater over the graph. A flocking boid has a spatial position $p = (p_1, p_2)$ and velocity $v = (v_1, v_2)$, which results from the composition of the velocities induced by the different forces pulling the boid, which are combined in the following well known linear model of the boid's velocity, which will be revisited below:

$$v = V_m (\alpha_s v_s + \alpha_c v_c + \alpha_a v_a + \alpha_n v_n) \tag{1}$$

Here s corresponds to separation, c to coherence, a to alignment and n to noise and parameters α are in the range $[0, 1]$. If we introduce v_A^R which denotes the Amity velocity calculated in a neighborhood of radius R and v_E^Z which is the Enmity velocity measured inside radius Z around boid's position, simplifying the notation $v_A^R = v_A$ and $v_E^Z = v_E$, we could introduce these effects into the boid's velocity equation (eq. 1) obtaining the following expression:

$$v = V_m \mathcal{N} (\alpha_{E_s} v_{E_s} + \alpha_{A_c} v_{A_c} + \alpha_{E_a} v_{E_a} + \alpha_n v_n) \tag{2}$$

Here $\|p\|$ denotes the norm of a position or vector p and V_m is a non-negative parameter that limits the norm of a vector and $\mathcal{N}(p) = \frac{p}{\|p\|}$ represents the normalized position or vector and $|S|$ denotes the cardinal of set S . Terms $\alpha_{E_s} v_{E_s}$ and $\alpha_{E_a} v_{E_a}$ are respectively the velocity components of separation and opposite alignment to enemies, $\alpha_{A_c} v_{A_c}$ is the component of coherence to friends and $\alpha_n v_n$ is a noisy component. We take normalized velocities v_{E_s} , v_{A_c} , v_{E_a} , v_n and normalize their linear combination to turn the head of a boid in the direction of the vector, giving after an step forward of length V_m .

2.1 Supervised Territoriality

The defense of territory by birds was modeled mathematically by J. Maynard Smith, who related in 1973 natural selection to the theory of games through the concept of evolutionary stable strategies [2,11]. Territories represent feeding and breeding zones that individuals defend from predators and from competitors of the same specie. In previous works [2] we introduced a territorial concept to model the a k-coloration process. We defined spatial goals associated with node colors, so that boids converging to a goal represent nodes being colored. This idea corresponds to a supervised territoriality modelling of the graph coloring process.

Supervised Territories: for each color $1, \dots, k$ we create a goal region on the swarm space. Usually we place Goals as the vertices of a regular polygon, and flocking boids are attracted by the goals with velocity $\alpha_g v_g$ being equation [3] the one that governs the model:

$$v = V_m \mathcal{N}(\alpha_{E_s} v_{E_s} + \alpha_{A_c} v_{A_c} + \alpha_{E_a} v_{E_a} + \alpha_n v_n + \alpha_g v_g) \quad (3)$$

Satisfaction: In supervised territoriality we attach a level of satisfaction to boids that measures their stress for being out of the desired spatial goals. Satisfaction grows until it reaches a maximum level whenever the boid is safe inside a goal without enemies. Cohabiting with enemies in a goal or being outside a goal decreases the satisfaction until 0 is reached.

Attack: Attack is a behavior triggered by the decrease of their boid satisfaction level. Attackers are individuals that trigger conflict because of their own low level of satisfaction, and the attacks can be of two types:

- External: Unsatisfied boids that are flocking outside the goals select a goal and try to expell enemies from it to replace them there.
- Internal: unsatisfied boids inside a goal try to expell enemies to obtain an homogeneous population inside the goal region.

2.2 The Scope of This Work

Besides, we find that the Unsupervised Flock Coloring (UFC), without spatial goals identified in the swarm evolution space, is able to find the chromatic number of a graph by convergence of the flock to clusters, for a sample of graphs of different topologies. A class of hard graphs for 3-coloration [8] by Bréaz algorithm [1] is explored. We measure the performance of UFC by indirect means, to avoid bias due to the interpretation of the spatial clusters.

3 The Rules of Supervised Territorial Flock Coloring (STFC)

The step from static models to swarms of flocking birds came from the hand of Craig Reynold's works [9,10] on computer animation. Rather than a grid, the

board becomes continuous, and cells are agents that can change their positions in the world. Each individual exhibits a very simple behavior that is specified by a few simple rules that guide them to get along with the collective motion of the flock. The global behavior of the flock emerges from these individual decisions. We will stick to the birds metaphor, so that in the following, we call boids to the agents that compose a flock.

3.1 STFC Scheme

A territorial flock coloring scheme is a tuple $\mathcal{F} = (\mathcal{B}, E, A, \mathcal{G}, \alpha, \beta)$ where:

- \mathcal{B} is a set of nodes that represent a population of boids.
- $E \subseteq \mathcal{B} \times \mathcal{B}$ is a undirected set of edges called Enmity edges. Graph $G = (\mathcal{B}, E)$ is the one we want to color. $E_i = \{b_j \in \mathcal{B} : (b_i, b_j) \in E\}$ is the set of enemies of boid b_i .
- $A \subseteq \mathcal{C}$ is a set of edges called Amity edges, being the friends of node b a subset of nodes at graph distance 1 of b , the enemies of the enemies that are not enemies of b :

$$C = \{(x, y) \in \mathcal{B} \times \mathcal{B} : \exists u ((x, u) \in E \wedge (u, y) \in E) \wedge (x, y) \notin E\}$$

$A_i = \{b_j \in \mathcal{B} : (b_i, b_j) \in A\}$ is the set of friends of boid b_i .

- \mathcal{G} is a set of goals. Each goal $g \in \mathcal{G}$ is in a position p and all them have a goal radius Γ .
- $\alpha = (\alpha_s, \alpha_c, \alpha_a, \alpha_n, \alpha_g) \in [0, 1]^5$ are the flock parameters for equation 3.
- $\beta = (S, R, Z, \Gamma, \rho, V_m, \Sigma, next - goal)$ are the parameters related to the size of the world: the side S , vision radius R , separation zone Z , goal radius Γ , angle of vision ρ , length of step V_m , maximum level of satisfaction Σ and criterion $next - goal$ to calculate the goal velocity.

The angle of vision, denoted ρ , comprises all the positions in the plane $p = p_i + v'$ where $\gamma_{v_i} - \frac{\rho}{2} \leq \gamma_{v'} \leq \gamma_{v_i} + \frac{\rho}{2}$. In this article, the value $\rho = 360^\circ$, meaning that individuals can even observe what happens behind their backs. Let R be the radius of vision: the cone of vision is the disk sector of radius R and angle ρ around the boid. The private radius Z is a zone that separates the boid from others.

Each boid is aware of an spatial region around it, its neighborhood. Let $|\partial_i^X|$ denote the population in the neighborhood of radius X . The set of boids lying inside the neighborhood of the i -th individual of radius X is denoted:

$$\partial_i^X = \partial^X (b_i) = \{b_j : dist(p_i, p_j) < X \wedge p_j \in cone_i\}.$$

where $dist$ is the euclidean distance, R is the ratio of vision, and $cone_i$ represents the population in the cone of vision of the individual i . The population in goal g is $\partial_g^\Gamma = \partial^\Gamma (g) = \{b_j : dist(p_i, p_j) < \Gamma\}$.

3.2 STFC Steering Rules

The steering basic rules, used in our model presented in equation 3, are the classical ones of the Reynolds model: alignment, separation, and cohesion. Combining these rules, the flocking birds, called boids in this model, are able to flight co-coordinately avoiding collisions. The flocking rules for the boid b_i are formalized as follows:

Separation from enemies: steer to avoid crowding local flock enemies inside a private zone of radius Z .

$$v_{s_i}^E = \mathcal{N} \left(- \sum_{b_j \in \partial_i^Z \cap E_i} (p_j - p_i) \right) \tag{4}$$

Cohesion towards friends: steer to move toward the average position c_i of local flockfriends in vision area R :

$$v_{c_i}^A = \mathcal{N} \left(c_i - p_i \text{ where } c_i = \frac{1}{|\partial_i^R \cap A_i|} \sum_{b_j \in \partial_i^R \cap A_i} p_j \right) \tag{5}$$

Alignment in opposite direction of enemies: steer in opposite direction of the average heading of local flockenemies.

$$v_{a_i}^E = \mathcal{N} \left(- \frac{1}{|\partial_i^Z \cap E_i|} \sum_{b_j \in \partial_i^Z \cap E_i} v_j - v_i \right) \tag{6}$$

3.3 Dynamics of STFC Systems

We present here the equations governing the velocity and position of a Boid. Let it be $p_i(0)$ the initial conditions and v_i the model of territorial flock coloring defined in equation 3.

$$p_i(t + 1) = p_i(t) + v_i(t + 1) \tag{7}$$

In a toric geometry, points p are represented as $p \bmod S$. The maximum euclidean distance between two points is hence $\frac{S\sqrt{2}}{2}$.

A boid configuration is a tuple $b_i(t) = (p_i(t), v_i(t), s_i(t), g_i(t))$ where $s_i(t)$ is the level of satisfaction and $g_i(t) \in \mathcal{G} \cup \{nobody\}$ denotes next goal, being *nobody* a default value meaning that no goals are being sought.

$$s_i(t + 1) = \begin{cases} s_i(t) + 1 & \text{in - goal - without - enemies } (b_i(t)) \\ s_i(t) - 1 & \text{in - goal - with - enemies } (b_i(t)) \\ & \vee \text{ enemies - in - all - goals } (b_i(t)) \\ s_i(t) & \text{otherwise} \end{cases}$$

where $x \dot{+} y = \min \{ \Sigma, x + y \}$ and $x \dot{-} y = \max \{ 0, x - y \}$. Predicates *in-goal - without - enemies* and *in-goal - with - enemies* are defines as follows:

$$\textit{in-goal - without - enemies} (b_i) \Leftrightarrow \exists g (b_i \in \partial_g^F \cap \bar{E}_i)$$

$$\textit{in-goal - with - enemies} (b_i) \Leftrightarrow \exists g (b_i \in \partial_g^F \cap E_i)$$

$$\textit{enemies - in - all - goals} (b_i) \Leftrightarrow \forall g (\partial_g^F \cap E_i \neq \emptyset)$$

Next-goal is a criterion that allows the selection of the goal toward which a boid navigates, $g(t + 1) = \textit{next-goal} (b_i(t))$ and can be for example: current goal, nearest goal, nearest free goal,...

A TFC system is a pair (\mathcal{F}, \bar{c}) where \mathcal{F} is a flock scheme and $\bar{c} = \bar{c}(0) = (c_1(0), \dots, c_{|\mathcal{B}|}(0))$ is the vector of initial configurations. A computation of length k of a TFC system is the sequence $\bar{c}(0) \dots \bar{c}(k)$.

4 Supervised Flock Coloring experiments

We implemented in Wolfram Mathematica Mizuno’s [8] generator of hard graphs for 3-coloration by Brélaz heuristic [1], which uses 4 colors for that class of 3-colorable graphs, and producing a sample of 100 hard graphs with population P between 100 and 112 nodes. We also generated a sample of 100 random graphs generated by Kuratowski’s theorem [3], for which $e \leq 3P - 6$ being P the number of nodes a e the number of edges. Complete bipartite graphs of 100 nodes were also analyzed. The side of the world was $S = 400$ and goals are placed forming a regular polygon inscribed in a circle of radius $\frac{S}{2}$.

4.1 Extreme Supervised Territorial Coloring

An Extreme Supervised Territorial Coloring (ESTC) scheme is a STFC $\mathcal{F} = (\mathcal{B}, E, A, \mathcal{G}, \alpha, \beta)$ where the parameters $\alpha = (\alpha_s, \alpha_c, \alpha_a, \alpha_n, \alpha_g)$ that control the flock are zero, $\alpha_s = \alpha_c = \alpha_a = 0$, meaning that the only motors of the system are goal velocity and attack. The ESTC algorithm depends on *next-goal* criterion, that we make equal to nearest goal in the experiments presented here, that extend the work that we began in [2]. Assume that internal attack expells enemies from the goal resetting the maximum satisfaction to both contenders.

As cascade coloration strategy can be used to find the chromatic number of a graph: it is sufficient to start the process of coloring successively the graph with colors $P, P - 1, \dots$. For 3-coloration, the execution of the program has two stages : First, the system attempts to find a 4-coloration (maximum 1500 iterations). Second, eliminate the less populated goal, the individuals newly freed wander to seek a new goal (maximum 3500 iterations). This procedure of cascading coloration is based on known works in reaction-diffusion particle systems [13]. An experiment consists in 25 runs of each one graph, registering the maximum number of well-colored nodes per run. We register the best configuration for each

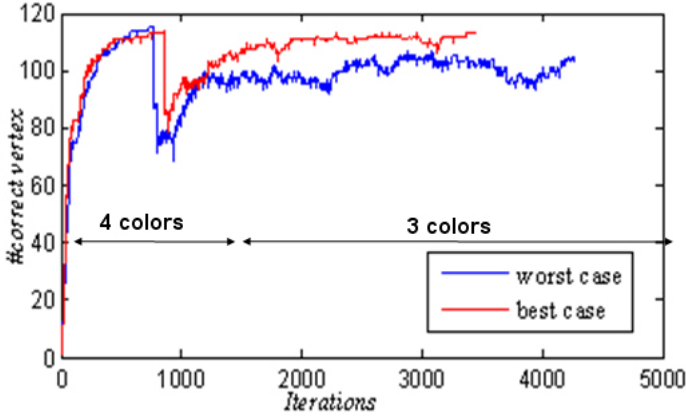


Fig. 1. Evolution of the number of correctly colored nodes in a cascade coloration process

experiment, the one that has minimum errors as well as the step at which the minimum is reached. Each run ends either when a 3-coloring solution is reached (success), or when 5000 iterations are completed in cascade. Figure 1 shows the typical plot of the number of correctly colored nodes along a cascade coloration process of a hard graph.

4.2 Mean Results of Experimentation

Table 1 shows the mean results of the experiments. We calculated the mean number of friends and enemies per node per graph, computing the *Hostility* quotient $H = \frac{\#friends}{\#enemies}$. The final iteration step, the percentage of successful runs per experiment and the percentage of well 3-colored nodes are registered in table 1.

Table 1. Average succesful runs for the main topological graph classes and graph topological measures

Mean values	Trees	Bipartite	Hard	Random
Number of friends	2.92	65.20	6.26	11.76
Number of enemies	2.00	33.83	3.48	3.54
Hostility	1.46	1.93	1.80	3.33
Final iteration	482.74	428.44	2345.70	1467.10
% successful runs	100.00	100.00	52.20	11.04
% well 3-colored nodes	100.00	100.00	99.21	92.16

In order to determine which independent variable (Number of friends, Number of friends or the quotient) is more correlated to dependent variables (Final iteration step, percentage of successful runs per experiment, percentage of well

3-colored nodes per run), we calculated the Pearson coefficient over the rows of **II**. The result is given in table

Table 2. Correlation of topological measures with success measures

Pearson coefficient r	Final iteration	% successful runs	% well 3 colored
Hostility	0,26	-0,86	-0,97
Number of enemies	-0,52	0,50	0,35
Number of friends	-0,50	0,42	0,26

The *Hostility* quotient shows a very high negative correlation to the percentage of successful runs per experiment (runs that find a 3-coloration) and to the percentage of well 3-colored nodes per run. From this result we select Hostility H as the best graph feature that explains tidily the results regardless of the topology of the graph.

4.3 Percentage of Well 3-Colored Nodes

Figure **2** shows (squares) that 3-coloration reaches the 100% of mean well 3-colored nodes in the case of trees and bipartite graphs that are 2-colorable. Hard graphs means are over the 97% of well 3-colored nodes. Random graphs

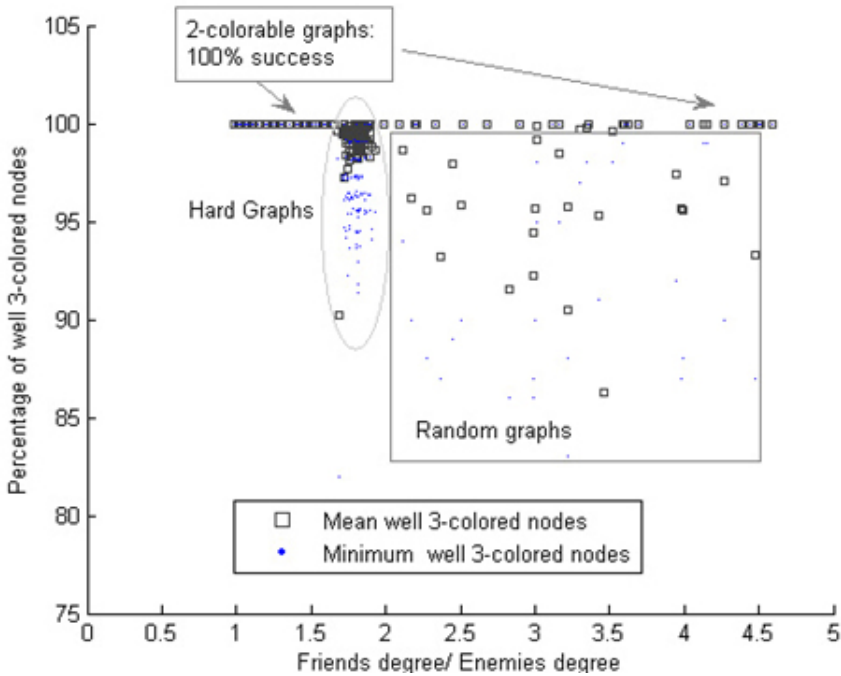


Fig. 2. Distribution of coloring success related to the Hostility quotient

show poor results in the case of graphs that are not 3-colorable since not all planar graphs are 3-colorable. The graphic shows also (dots) the worst cases per experiment: even in this case, hard graphs obtain results above 90%.

4.4 Percentage of Successful Runs Per Experiment

Respecting to the hits, i.e. the number of runs per experiment that end in a 3-coloration, figure 3 shows that the majority of the cases are over the 40% and the mean is about a 75% of runs that end with success. Again, 2-colorable graphs reach the 100% of hits rate. The sample of hard graphs exhibit a wider variability, but 1 of two runs ends with success as can be observed in table 1. Random graphs, remember that not all of them are 3-colorable, obtain quite poor results, but only 4-colorable random graphs reach the 0% of hits success.

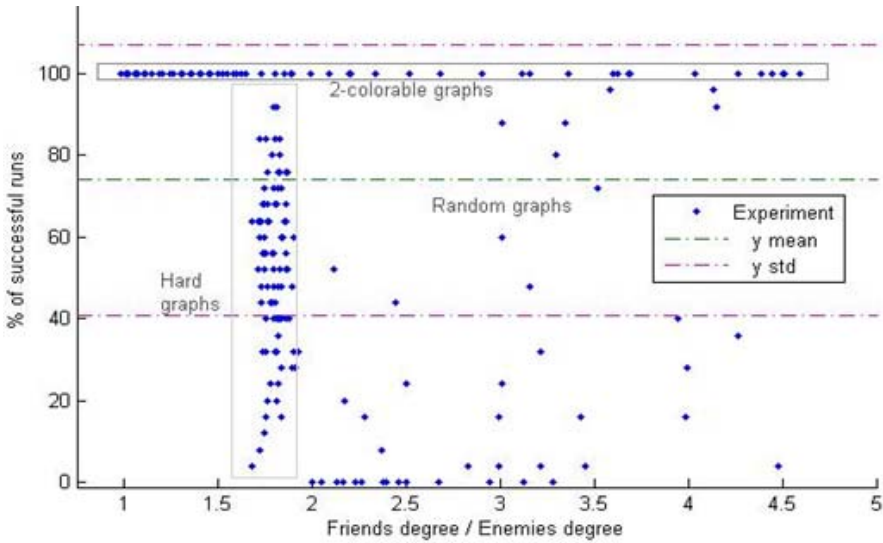


Fig. 3. Number of succesfull runs as a function of the Hostility quotient

4.5 On the Time Needed to Color a Node

To have a measure of computing performance, we estimate the minimal time cost needed to color correctly one node dividing the minimal number of iterations of one run in one experiment by the maximum number of well 3-colored nodes in that experiment. Figure 4 shows the distribution of this value for the experiments performed, identifying the graph topology class. Lowest costs correspond again to 2-colorable graphs, in spite of the proportion of friends over enemies of trees and bipartite graphs are very different. Highest costs correspond to hard graphs; 3-colorable random graphs are divided between those that obtain intermediate time cost values and those that get immediatly a 3-coloration due to previous cascade 4-coloration left the problem almost solved.

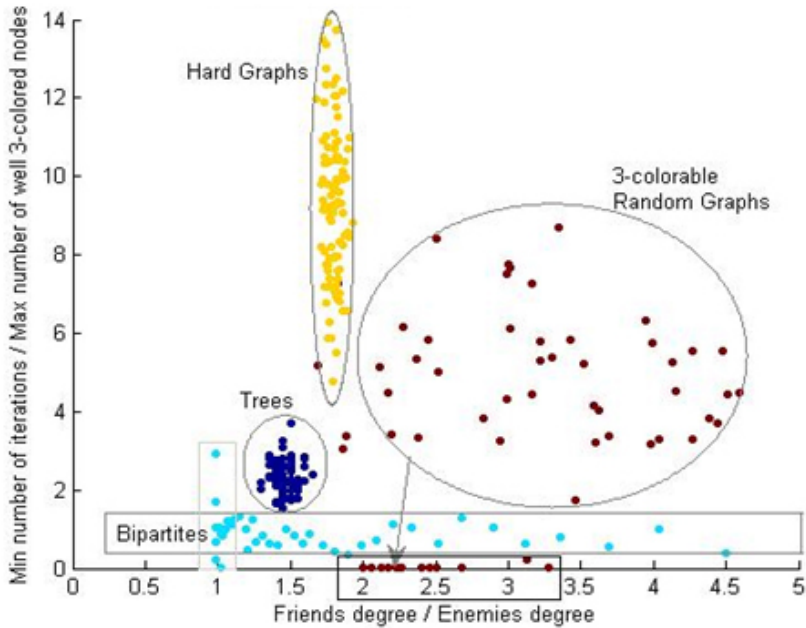


Fig. 4. Number of iterations versus the Hostility quotient

5 Conclusions

We propose in this paper a framework to study GCP in terms of flocking birds that show territorial behavior, calling the algorithm Supervised Territorial Flock Coloring (STFC). We studied separately territoriality and flocking, finding that Extreme Territorial Flock Coloring (ETFC), based in goal velocity and attack is able to solve GCP. ETFC algorithm is able to color Mizuno's hard graphs for 3-coloration while Brélaz algorithm fail for all them.

We found that a global graph parameter, the ratio of nodes at distance one (enemies) to nodes at distance greater than one (friends), which we have called hostility, is a good characterization of the difficulty of solving GCP for each graph. It has negative strong correlation with the success and time to solve the problem.

References

1. Brélaz, D.: New methods to color vertices of a graph. *Communications of ACM* 22, 251–256 (1979)
2. Cases, B., Hernandez, C., Graña, M., D'anjou, A.: On the ability of swarms to compute the 3-coloring of graphs. In: Bullock, S., Noble, J., Watson, R., Bedau, M.A. (eds.) *Artificial Life XI: Proceedings of the Eleventh International Conference on the Simulation and Synthesis of Living Systems*, pp. 102–109. MIT Press, Cambridge (2008)

3. Diestel, R.: Graph Theory, electronic edn. Springer, USA (2005)
4. Dorigo, M., Blum, C.: Ant Colony Optimization theory: a survey. *Theoretical Computer Science* 344, 243–278 (2005)
5. Dorigo, M., Gambardella, L.M.: Ant Colonies for the Traveling Salesman Problem. *BioSystems* 43, 73–81 (1997)
6. Galinier, P., Hertz, A.: A survey of local search methods for graph coloring. *Comput. Oper. Res.* 33, 2547–2562 (2006)
7. Galinier, P., Hertz, A., Zufferey, N.: An adaptive memory algorithm for the k-coloring problem. *Discrete Applied Mathematics* 156(2), 267–279 (2008)
8. Mizuno, K., Nishihara, S.: Constructive generation of very hard 3-colorability instances. *Discrete Applied Mathematics* 156(2), 218–229 (2008)
9. Reynolds, C.W.: Flocks, herds, and schools: A distributed behavioral model. *Computer Graphics* 21, 25–34 (1987)
10. Reynolds, C.W.: Steering Behaviors for Autonomous Characters (1999), <http://www.red3d.com/cwr/papers/1999/gdc99steer.html>
11. Smith, J.M.: Evolution and the theory of games. Cambridge University Press, Cambridge (1992)
12. Smith, J.M., Price, G.R.: The logic of animal conflict. *Nature* 246 (November)
13. Turk, G.: Generating textures on arbitrary surfaces using reaction-diffusion. *SIGGRAPH Computers Graphics* 25(4), 289–298 (1991)

Data Collection System for the Navigation of Wheelchair Users: A Preliminary Report

Yasuaki Sumida, Kazuaki Goshi, and Katsuya Matsunaga

Kyushu Sangyo University, Graduate School of Information Science,
2-3-1 Matsukadai, Hgashi-ku, Fukuoka 813-8503, Japan

K07jk058@ip.kyusan-u.ac.jp, {goshi,matsnaga}@is.kyusan-u.ac.jp

Abstract. In Japan, population of aged people is becoming greater. According to it, number of the wheelchair user would become larger. The government has been improved to lessen the barriers for these people. But depending on the physical power of wheel chair users, it was pointed out that there could be the places where they could not move by their physical power. This is a reason that many wheelchair users do not move by their alone. Therefore we tried to develop the navigation system to eliminate the problem. At the beginning we have developed the system to measure the required force to move the wheelchair, the width of the space of aisle where they could turn back the direction, and the position information. By the data base, we planned to develop the navigation system for wheelchair users who could move alone to working places or other places where they want to go. In this paper, we report the system developed to collect data for navigation of wheelchair users.

Keywords: wheelchair, navigation, GPS, QR-code.

1 Introduction

In Japan, the population of aged people is increasing [1]. Accordingly, wheelchair users are expected to increase. Government regulation has lessened barriers that hinder them. However, depending on the physical power of wheelchair users, there are places where they are unable to move using their own physical power. Consequently, many wheelchair users do not move around by themselves.

For wheelchair users to live independently, they should be supported to move anywhere. For example, the system which helps movement of wheelchair [5] and the system which guides a movable course of the wheelchair [2, 6] would be necessary. However, almost all such studies investigate systems to show routes. No system has been designed to collect information for a navigation database.

This research is undertaken to develop a system that navigates wheelchair users who want to go to the destination and give them the possibility to venture freely outdoors. Regarding the research, we first developed a system to measure the force to move a wheelchair, the size of the space in which a wheelchair might turn back, and position information. These data of forces for wheelchair movement and space sizes can be matched with building plans and roads. As described, we report this system developed to collect data for wheelchair user navigation.

2 Overview of the Developed System

Fig. 1 presents an overview of the system we intend to develop. Using the data collector, information for the navigation of wheelchair users will be collected and transmitted to the database server. The physical power of wheelchair users to move a wheelchair is measured. Using the internet, wheelchair users can find a way to a destination prior to using the database in the server. Moreover, wheelchair users will be navigated in real time. Herein, we report the data collection system.

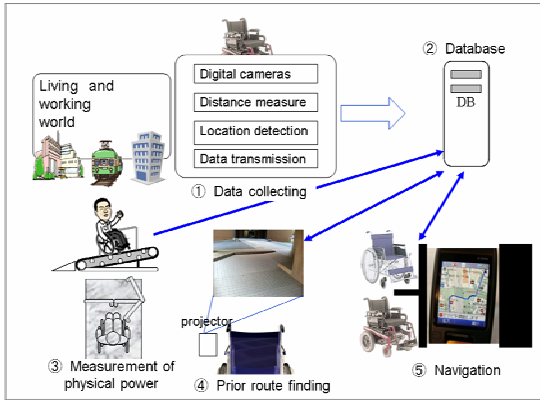


Fig. 1. Overview of the developed system

3 Data Collection Method

The necessary data for the system to navigate wheelchair users is location information of consecutive measuring points, the necessary force to move a wheelchair to go to the destination, information about space around the wheelchair for checking whether it is possible to turn back or not, and the distance and time to the destination.

3.1 Measurement System to Collect Location Data

Force to move the wheelchair from the place of departure to the destination should be corresponding to the map of the building and the road. The wheelchair positions are identifiable using a Global Positioning System (GPS). We use the RMC format of GPS to obtain the position information. The RMC format includes time, latitude, longitude, ground speed, and direction of movement: such location data are sent as serial data. The serial data are input to a built-in PC board via an AVR microcontroller with other measured data. It is thereby possible to record all measured data with location information. However, it is difficult to identify places within the buildings using GPS. We attempted to identify the places in the buildings by the moving distances measured by frequencies of wheelchair rotation from anchor points differentiated using QR codes assigned to floor plans.

The direction in which the wheelchair moves is identified using the QR code. In Fig. 2, sheets on which QR codes are printed are attached on the ceiling at crossing points of the passageways. The QR codes can indicate the wheelchair's location and direction of movement. The wheelchair position in the building is obtained according to the moving direction identified by QR codes and moving distance measured using the frequencies of rotation of the wheel from the position of certain QR codes attached on the ceiling. The QR code includes position information in the building. In Fig. 3, an example of QR code is shown. This QR code contains information of "33.669835,130.444867,11" which means the degrees of latitude, the degree of longitude and the altitude. The QR code includes position information in the building. In Fig. 3, a QR code is shown: it contains information of "33.669835,130.444867,11", denoting the degrees of latitude, the degree of longitude and the altitude. It is thereby possible to obtain information about direction within buildings.

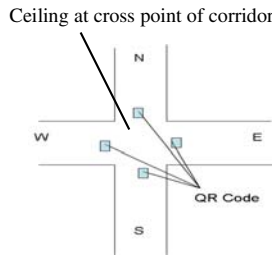


Fig. 2. Places where image charts of QR codes are attached



Fig. 3. A QR code

3.2 System to Measure Forces Moving the Wheelchair

Wheelchair users are not always able to move on any path. For example, it might be difficult to move on steep slopes, roads with different levels, and pathways covered with deep-pile carpets. However, no map shows road surface states in detail. Therefore, the force necessary to move wheelchair must be inferred.

The force necessary to move the wheelchair is determined by the current to the motor of the electric wheelchair. The current to the motor is obtained by conversion of voltage by the resistor connected between the motor and the batteries. The voltage between edges of the resistor is amplified and converted into a digital value by the microcontroller's (AVR, ATMEGA644P: Atmel Corp.) AD converter. The obtained value is calibrated into the value of the drive torque and into the force.

3.3 Space Width Measurement

Wheelchair users want to know the width of spaces on the way to the destination. They cannot turn back in narrow spaces. Therefore, for route navigation, space information must be given that the space is sufficiently wide or not. To turn back, the passage width should be more than 3.6 m on many hand-driven wheelchairs and many electricity-driven wheelchairs [4]. Width is measured using distance measuring devices of ultrasonic type (EZ0 and EZ1, LV-MaxSonar: MaxBotix Inc.). The distance measurement range is 0–6.5 m. The ultrasonic device outputs a pulse width that is proportional to the distance to the object. The distance is obtainable by measuring the pulse width by the counter of the AVR microcontroller. Equation (1) is used for computation to make one count equal to 1 cm.

$$\text{Count cycle} = \{1 \text{ [s]} / (0.01 \text{ [m]} / \text{acoustic velocity [m/s]})\} / 2 \quad (1)$$

The acoustic velocity is altered by temperature, humidity, and air density. For that reason, it is necessary to calibrate it according to these values.

3.4 Movement Distance and Velocity Measurement

As described previously, it is difficult for the GPS to obtain location data in buildings. Therefore, moving distances are measured using the wheelchair wheel's rate of rotation. It is obtained by counting by a device of counter a circular pattern of monochrome alternation inside the wheel. The counter using a photograph interrupter outputs pulses when it detects a white pattern. The rate of wheel rotation was obtained by counting the standing up of the pulse every second using an AVR microcontroller. The movement distance was obtained using eq. (2). The movement velocity is calculable according to the movement distance and movement time.

$$\text{Movement distance each one second [m]} = \{ \text{wheel circumference [m]} / \text{number of white pattern} \} \times \text{pulse count each second} \quad (2)$$

4 Design and Implementation

4.1 Hardware and Software Based on ASSIST

We developed the Assistant System for Safe Driving by Informative Supervision and Training (ASSIST) [3], which records driving parameters including velocity, headway distance, and position: it sends the data to a supervisor at a remote place to educate drivers about safe driving. We decided to develop the data collection system based on ASSIST because the ASSIST hardware and software for recording and transmitting are applicable to collection of routing information for wheelchair users.

4.2 Hardware Design

4.2.1 Camera to Obtain Front View Images

People can easily find their way if landmarks are given. We designed the system to get pictures using a video camera while the probe car was collecting data. The video camera was attached on the front of the wheelchair.

4.2.2 Overall Design of the Hardware (W-ASSIST Hardware)

Fig. 4 shows a block diagram of the hardware with an embedded computer board, an interface, video cameras, distance sensors, GPS, and a rotation sensor. The interface was designed with an AVR microcontroller (ATMEGA644P: Atmel Corp.). The interface hands over the data from the sensors to the embedded PC. The embedded PC board was connected with cameras to obtain front view images and QR codes, and a mobile phone for data transmission to the data server and the web server.

Fig. 5 shows the electric wheelchair with our measurement system.

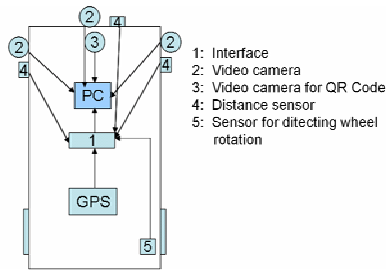


Fig. 4. Overview of the hardware



Fig. 5. The electric wheelchair used for measurements

4.3 Software(W-ASSIST Software)

4.3.1 Design of the Software in the Interface and the Embedded PC Board

Data from all sensors and GPS are input into the AVR microcontroller in the interface board. Then the inputted data are related with the GPS location data and sent to an embedded PC board as serial data. The data transmission is done every second when the data of GPS are complete. The data are transmitted to the data server and the web server. Fig. 6 portrays a block diagram of the system showing the data flow.

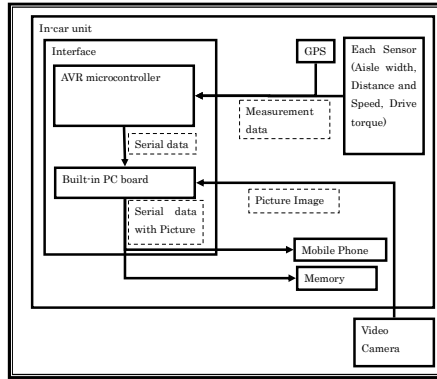


Fig. 6. Block diagram of the system and the flow of data

4.3.2 Software Design

It is necessary to monitor the system to verify its function. We developed software to check the system visually and display graphical images on the screen to confirm the logged data and to monitor the system working in real time. The software was made by revising the ASSIST system software.

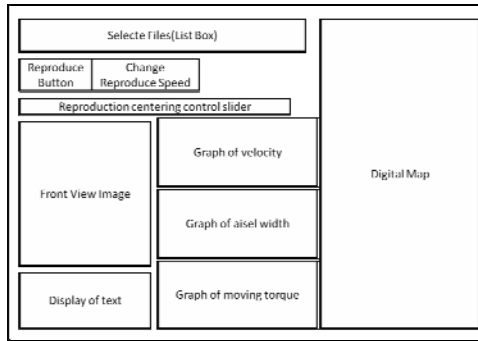


Fig. 7. Design of the log-viewer software

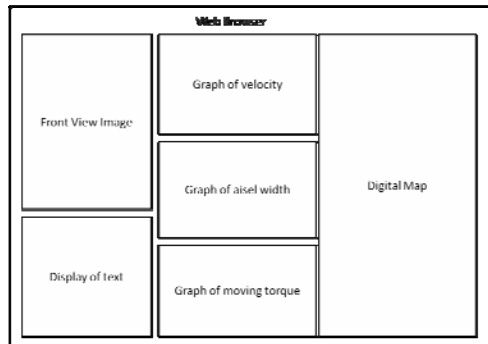


Fig. 8. Design of the live-monitoring software

We designed and developed two software, log-viewer software for checking log data, and the live monitor software for a rear-time check. Fig. 7 shows the design of the log-viewer software. Fig. 8 shows the design of the live-monitoring software.

5 Results

5.1 Results of System Performance

The measured data is stored in a flash memory by the character string by the recorder program. The log-viewer software obtains each data by dividing this character string. The stored data can be checked by the log-viewer software. Fig. 9 shows a method of data processing in log-viewer software. The data is stored in order of the position information from the GPS, measurement data of a microcomputer, data for debugging. The Log-viewer software displays texts and graphs by reading this character string, dividing it, and obtaining each data. Fig. 10 presents a screen display drawn from data logged in the W-ASSIST hardware using the log-viewer software. The hardware and the software of the system worked well to get signals from the GPS, the camera, the current measurement device, the distance measurement devices, and the wheel rotation-counting device. The log-viewer software is developed using Java.

The live-monitoring software consists of a recording program, a transmitting program, a W-ASSIST server program and a web interface. Fig. 11 shows the flow of information of the live-monitoring software. The recording program sends data to summing program through UDP socket. The transmitting program sends live data to the server of ASSIST through TCP socket on the mobile network. The live data consists of times, wheelchair's IDs, measurer's IDs, latitudes, longitudes, measurement data and front view images. The graphs are drawn using measured data for 20seconds

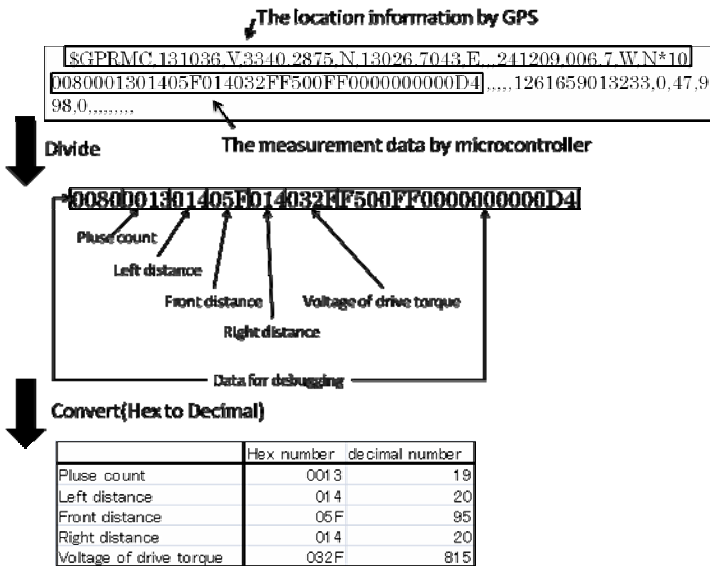


Fig. 9. Method of data processing in log-viewer software

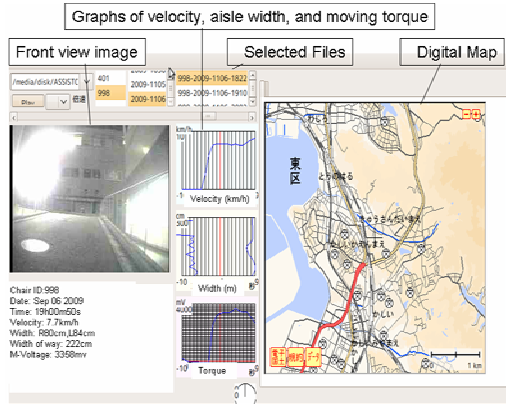


Fig. 10. Screen image displayed by the live-monitoring software through the web

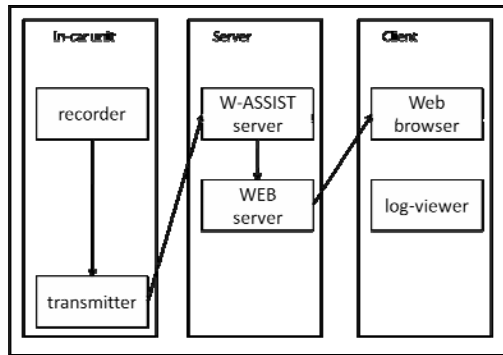


Fig. 11. Block diagram of the live-monitoring software

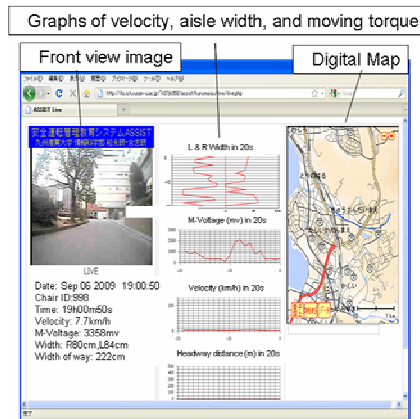


Fig. 12. Screen image displayed by the live-monitoring software through the web

each. Fig. 12 shows a screen display drawn from the data obtained by the W-ASSIST hardware using the live monitor software. Data were transmitted to the servers using the mobile phone system via the internet system. Software for live monitoring of the driving behavior worked well (Fig. 12). The live-monitoring software is developed using PHP and javascript.

5.2 Measurement Accuracy

- The method of experiment

In order to verify about the reproducibility of the measurements, the measurements were done five times by running wheel chair (Fig. 13). In fig. 15, A, B, C and D mean the measurement points.

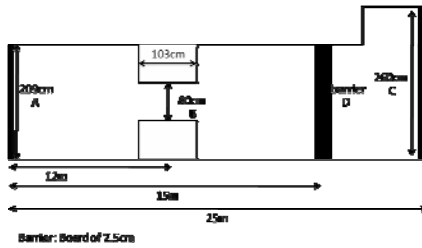


Fig. 13. Experimental environment

- The result of experiment

Table 1 shows the result of moving mileages measured. In Table 1, moving distances were calculated by multiplying the length of every one pulse by the pulse number total. The measured distances have the errors of about 3 meters at most. However, the differences of the totals of pulse numbers are small enough. The cause of the errors of the distance measurements is thought to be derived from the distortion of the tires while wheelchair running.

Table 1. The result of measured mileage

Measurement time	1	2	3	4	5
Total of pulse number	563	562	563	562	563
Moving Distance[cm]	2815	2810	2815	2810	2815
Moving Distance[m]	28.15	28.1	28.15	28.1	28.15

Fig. 14 shows the results of measured aisle width which were gotten by the distance measuring devices. The differences could be thought to be small.

Fig. 15 shows the result of the measured voltage to analyze the currents to the motors. Although there are some differences, the reproducibility could be thought to be high enough.

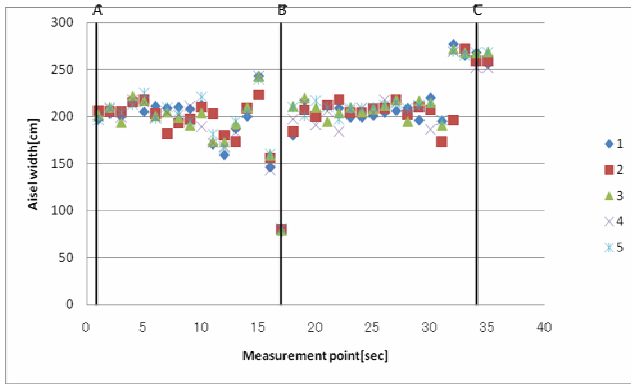


Fig. 14. Result of measured aisle width

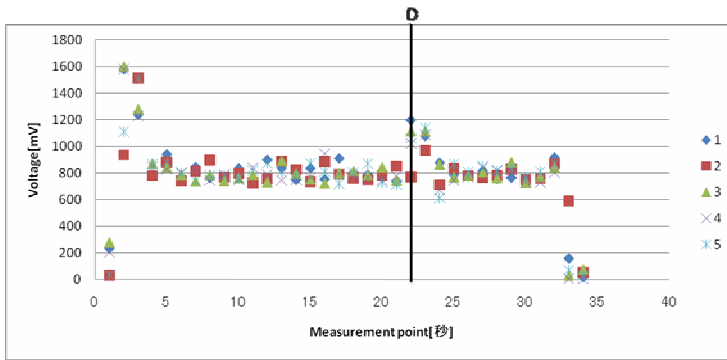


Fig. 15. The result of measured voltage to analyze the driving force

5.2.1 Verification of the Relation between Drive Torque and Actual Power

To confirm the correlation of the voltages and the forces to move wheelchair.

- The method of experiment

The measurements were done five times by running the wheelchair on the slopes of the same asphalt pavement. Those slope's gradient are different. And the forces to move wheelchair were measured by the spring balance.

- The result of experiment

Fig. 16 shows the result of the measured forces to move wheelchair using the spring balance and the voltage while the wheelchair movement for the measurement. Both inclinations could be thought to be proportional to the slope angles. In addition, we conducted multiple linear regression analysis about the starting powers and the running powers, and got the multiple regression equations (eq. (3, 4)).

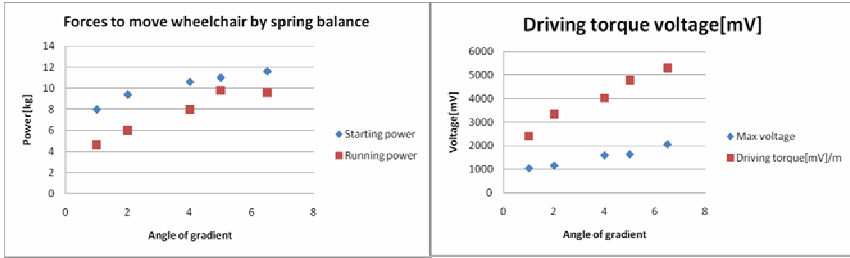


Fig. 16. Result of experiment

$$\text{The needed starting power [kg]} = 5.2054 + 0.0032 * \text{Max voltage[mV]} \tag{3}$$

$$\begin{aligned} \text{The needed running power[kg]} = \\ -0.037 + 0.00192 * \text{Average driving torque voltage[mV]} \end{aligned} \tag{4}$$

The result of the predicted value computed based on these equations is shown in Table 1. It is almost the same value, and if some errors would be permitted, the fore to move wheelchair could be converted into the driving torque by these multiple regression equations.

Table 2. Prediction result

Angle of gradient	Driving torque voltage[mV]/m	Running power[kg]	Predicted running power[kg]	Max Driving torque voltage[mV]	Starting power[kg]	Predicted starting power[kg]
1	2401.4	4.6	4.6	1033.0	8.0	8.6
2	3348.2	6.0	6.4	1146.0	9.4	9.0
4	4025.0	8.0	7.7	1585.0	10.6	10.4
5	4787.5	9.8	9.2	1628.0	11.0	10.6
6.5	5309.5	9.6	10.2	2052.0	11.6	12.0

5.3 Result of a QR Code Reading System

We developed the QRcode reading system using the opencv1.1 library and the libdcodqr library by C language. QR Codes are read in the following procedure. First, the still pictures were gotten using camera. If a QR code in a still picture image was found, it was clipped out. Next, it's image was expanded into double and the QR code was recognized. Finally, the location information in the QR code was analyzed. Fig. 19 shows an example of the screen display where the system was executed actually.

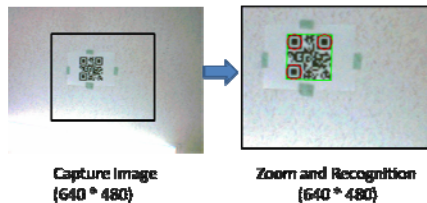


Fig. 17. The examples of execution of the QR code reading system

- The method of experiment

The QR code readings were done using the wheelchair for the measurement under the five QR codes which were posted to the ceiling. We succeeded in the reading of three of the QR cords in five pieces.

5 Conclusion and Future Studies

As described herein, we reported a system developed to collect data to navigate wheelchair users. Results show the possibility of an automatic data collecting system used for the navigation system database. A measurement function will be added to collect more data including route height information and ground surface properties.

References

1. Statistics Bureau in Japan,
<http://www.stat.go.jp/data/jinsui/tsuki/index.htm>
2. Yang, A., Mackworth, A.K.:
<http://people.cs.ubc.ca/~mack/Publications/CCAI07.pdf>
3. Goshi, K., Matsunaga, K., Kuroki, D., Shidoji, K., Matsuki, Y.: Educational Intelligent Transport System ASSIST. In: Proceedings of the Fourth IASTED International Conference Computers and Advanced Technology in Education, Banff, pp. 150–154 (2001)
4. Epoch House, <http://www.epochhouse.com/care05.htm>
5. Wang, H., Salatin, B., Grindle, G., Ding, D., Cooper, R.: Real-time model based electrical powered wheelchair control. *Medical Engineering & Physics* 31(10), 1244–1254
6. Kasemsuppakorn, P., Karimi, H.A.: Personalised routing for wheelchair navigation. *Journal of Location Based Services* 3(1), 24–54 (2009)

Author Index

- Abawajy, Jemal H. III-201
Abraham, Ajith III-472
Ahmadian, Kushan I-574
Aларcon, Vladimir J. I-491, I-501
Amorim, Ronan M. IV-395
Anders, Frauke I-152
Andriamasinoro, Fenintsoa I-476
An, Hong IV-427
Anuar, Mohd Hafiz I-331
Aoki, Takaaki IV-252
Arai, Kohei II-71, II-87, II-336, III-305
Areerachakul, Sirilak III-215
Arikan, Yüksel Deniz II-544
Aritsugi, Masayoshi II-412
Asche, Hartmut I-346, I-515
Atman, Nilüfer II-544
- Baba, Kensuke IV-236, IV-273
Badea, Radu II-215
Bagstad, Kenneth J. I-238
Bai, Songnan IV-520
Barros, Diego Martins Vieira I-430
Basuki, Achmad II-87
Battaglia, Francesco I-1
Benigni, Gladys II-422
Benkner, Siegfried IV-13
Beristain, Andoni I-610
Bernardi, Milena II-206
Bhuruth, Muddun II-570
Bi, Zhongqin IV-482
Blecic, Ivan I-166
Boizumault, Patrice II-432
Boo, Chang-Jin II-99, II-110
Boojhawon, Ravindra II-570
Borruso, Giuseppe I-1
Bramley, Randall II-503
Brankovic, Ljiljana II-586
Brooks, Christopher I-501
Bucur, Razvan II-215
Burry, Jane III-483
- Cacheiro, Javier Lopez IV-41
Cafer, Ferid II-301
Canters, Frank I-89
- Cao, Lu IV-427
Cao, Yanzhao IV-409
Cardell-Oliver, Rachel III-336
Carlini, Maurizio II-177, II-206
Casas, Giuseppe Las I-62
Cases, Blanca III-541
Cassard, Daniel I-476
Castellucci, Sonia II-177
Cattani, Carlo II-155, II-164, II-215,
II-225
Cattrysse, Dirk I-414
Cecchini, Arnaldo I-166
Chai, Kevin II-351
Chan, Chien-Hui III-526
Chang, Maw-Shang II-314
Checiches, Gabriela II-215
Chen, Changqian IV-296
Cheng, Chin-Yi III-431
Cheng, Kai III-395, IV-418
Chen, Gong-liang II-14
Chen, Jiming IV-296
Cheon, Sung Moon I-182
Che, Z.H. II-533
Che, Zhen-Guo II-533
Chiang, C.J. II-533
Chiang, Tzu-An II-533
Chi, Hongmei IV-409
Choi, Bum-Gon III-85
Choi, Hyung-Jin III-118
Choi, Wook III-129
Chong, Dahae III-21, III-31
Choo, Hyunseung III-85, III-96,
III-129, III-158
Chou, Shihchieh III-431
Cho, Yongyun III-258, III-269
Chuang, Chun-Ling III-178
Chu, Hone-Jay I-116
Chung, Min Young III-63, III-85, III-96
Chung, Tai-Myoung III-142, III-352
Chyan, Shyh-Haw I-293
Ciloglugil, Birol II-556
Coluzzi, Rosa I-361
Cong, Ming IV-427
Costantini, Alessandro IV-29, IV-41

- Cracknell, Arthur P. I-545
 Crémilleux, Bruno II-432
 Crisan, Diana II-215
 Crisan, Maria II-215
- Dai, Miao Ru III-368
 Dan, Avishek III-321
 Danese, Maria I-320
 D'Anjou, Alicia III-541
 Daněk, Josef IV-62
 DasBit, Sipra III-321
 Datta, Amitava III-336
 Davoine, Paule-Annick I-445
 de Doncker, Elise II-139
 Delgado-Mohatar, Oscar II-586
 Deris, Mustafa Mat III-201,
 III-405, IV-175
 Desudchit, Tayard III-419
 Dickstein, Flavio II-475, IV-395
 Di Donato, Pasquale I-528
 Dohi, Tadashi IV-441
 Dong, Fei I-131
 Dong, Wei II-463
 Dookhitram, Kumar II-570
 dos Santos Amorim, Elisa Portes II-475,
 IV-395
 dos Santos, Rodrigo Weber II-475,
 IV-395
 Dostálová, Taťjana IV-62
 Dvorský, Jiří III-472
- Edwards Jr., David II-1
 Emanuele, Strano I-32
 Embong, Abdullah IV-83
 Engelen, Guy I-89
 Ervin, Gary I-501
 Esposto, Stefano II-206
- Faisal, Zaman IV-199
 FanJiang, Yong-Yi II-257
 Fidalgo, Robson do Nascimento I-430
 Firoozeh, Nazanin II-370
 Firuzeh, Nazanin II-400
 Fischer, Edward II-285
 Florea, Mira II-215
 Fujimoto, Junpei II-139
 Fujita, Shigeru IV-119, IV-128
 Fuster-Sabater, Amparo II-586
- Gansterer, Wilfried N. IV-13
 Garcia, Ernesto IV-1
- Gavrilova, Marina I-574
 Gensel, Jérôme I-445
 Gervasi, Osvaldo II-422, IV-41
 Ghosal, Amrita III-321
 Ghosh, S.K. I-309
 Gizzi, Fabrizio I-320
 Glorio, Octavio I-461
 Goi, Bok-Min IV-188
 Goldfeld, Paulo II-475, IV-395
 Goshi, Kazuaki III-552, IV-497
 Gotoh, Yusuke II-324
 Graña, Manuel I-610, III-541
 Gutierrez, Eduardo IV-41
- Hadden, John IV-358
 Halder, Subir III-321
 Hamaguchi, Nobuyuki II-139
 Han, Yi IV-263
 Han, Young-Ju III-142, III-352
 Hara, Hideki IV-119
 Harsono, Tri II-71
 Hasegawa, Hidehiko II-60
 Hashim, Mazlan I-331, I-545
 Hatzichristos, Thomas I-140
 Hayashi, Masaki IV-497
 Hayati, Pedram II-351, II-400
 Hedar, Abdel-Rahman IV-457
 He, Jie III-498
 Heng, Swee-Huay IV-188
 Herawan, Tutut III-201, III-405, IV-175
 Hernandez, Carmen III-541
 Hernández, Constanza II-361
 Hirata, Kazuhiro IV-497
 Hirose, Hideo IV-199
 Hlináková, Petra IV-62
 Hong, Youngshin III-52
 Hope, Martin III-228
 Hsieh, Nan-Chen III-526
 Huang, Wong-Rong II-257
 Huang, Zequn IV-520
 Hung, Ruo-Wei II-314
- Ikeda, Daisuke IV-236
 Im, Se-bin III-118
 Inceoglu, Mustafa Murat II-556
 İnceoğlu, Mustafa Murat II-544
 Inenaga, Shunsuke IV-236
 Ishikawa, Tadashi II-139
 Ishiwata, Emiko II-60
 Ismail, Rashad IV-457

- Itokawa, Tsuyoshi II-412
 Ito, Taishi IV-138
 Iwane, Masahiko II-488
 Izumi, Satoru IV-152

 Jamel, Sapiee IV-175
 Janciak, Ivan IV-13
 Jang, Jun-Hee III-118
 Jang, Myungjun I-262
 Jazyah, Yahia Hasan III-228
 Jehng, Jihn-Chang III-431
 Jeong, Seungmyeong III-72
 Jeong, Yeonjune III-158
 Jeung, Jaemin III-72
 Jia, Xiaoqi IV-468
 Jing, Jiwu IV-468
 Ji, Yindong II-463
 Jo, Heasuk IV-510
 Johnson, Gary W. I-238
 Joo, Yongjin I-105
 Jun, Chulmin I-105
 Jung, Jaeil IV-520
 Jung, Soon-Young IV-376
 Ju, Shiguang IV-296

 Kaio, Naoto IV-441
 Kalisch, Dominik I-152
 Kaneko, Kunihiko III-189
 Kang, Ji-Ae III-11
 Kang, Min-Jae II-99, II-110, III-11
 Kang, Seung Goo III-21
 Kawato, Akifumi IV-164
 Khiari, Mehdi II-432
 Kim, Byung-Sung III-158
 Kim, Chang Seup III-85
 Kim, Choel Min III-1
 Kim, Dong In III-42
 Kim, Ho-Chan II-99, II-110
 Kim, Hyeon-Cheol IV-376
 Kim, Hyunduk III-158
 Kim, Jae-Yearn II-119
 Kim, Jingyu III-42
 Kim, Jong-Myoung III-352
 Kim, Junhwan III-31
 Kim, Kyu-Il I-271
 Kim, Kyungill IV-370
 Kim, Sang-Wook III-1
 Kim, Seungjoo IV-510
 Kim, Shin Do I-182
 Kim, Taeyoung III-129

 Kim, Tai-Hoon II-422
 Kinoshita, Tetsuo IV-107, IV-138,
 IV-152, IV-164
 Kitagata, Gen IV-164
 Kitasuka, Teruaki II-412
 Kobayashi, Yusuke IV-152
 Koehler, Martin IV-13
 Köhler, Hermann I-152
 König, Reinhard I-152
 Konno, Susumu IV-107, IV-119
 Krömer, Pavel III-472
 Kudreyko, Aleksey II-155
 Kuo, Jong-Yih II-257
 Kurihara, Yoshimasa II-139
 Kusuda, Tetsuya IV-336
 Kwak, Ho-Young III-11
 Kwong, Kim-hung I-374, I-389
 Kwon, Young Min III-63

 Laganà, Antonio IV-1, IV-41
 Lago, Noelia Faginas IV-29
 Lai, Poh-chin I-374, I-389
 Lanorte, Antonio I-361
 Lasaponara, Rosa I-254, I-361
 Laserra, Ettore II-225
 Le, Thuy Thi Thu I-401
 Lee, Chang H. IV-370
 Lee, Cheng-Chi I-599
 Lee, Eunseok IV-385
 Lee, Im Hack I-182
 Lee, Jin-Kyung I-271
 Lee, Junghoon III-1, III-11, III-52
 Lee, Ju Yong III-63, III-85
 Lee, Kwang Y. II-99
 Lee, Myungsoo III-31
 Lee, Ok Kyung III-63
 Lee, Saebyeok IV-376
 Lee, Sang Joon III-52
 Lee, Seungil I-271
 Lee, Tae-Jin III-85, III-96
 Lee, WonGye IV-376
 Lee, Youngpo III-21, III-31
 Lee, Youngyoon III-21
 Lee, Yue-Shi III-458
 Li, Chun-Ta I-599
 Li, Jian-hua II-14
 Li, Ming II-191
 Lim, HeuiSeok IV-370, IV-376
 Lim, Jaesung III-72
 Lin, Feng-Tyan I-77, I-293

- Lin, Jingqiang IV-468
 Lin, Rong-Ho III-178
 Lin, Yu-Pin I-116, I-224
 Liou, William W. II-25
 Li, Peng IV-427
 Lischka, Hans IV-13
 Li, Tiancheng II-44
 Liu, Dong IV-427
 Liu, Fang II-503
 Liu, Hsiao-Lan I-293
 Liu, Liang I-590
 Liu, Peng IV-468
 Liu, Yuan IV-427
 Li, Yin II-14
 Li, Yu I-590
 Lobosco, Marcelo IV-395
 Lursinsap, Chidchanok III-419
 Lu, Tianbo IV-263
 Lu, Wenjie I-590
- Maggio, Grazia I-210
 Mahmud, Mohd Rizaludin I-331
 Manabe, Yusuke IV-119
 Mancini, Francesco I-210
 Mantelas, Lefteris A. I-140
 Mardiyanto, Ronny II-336
 Marghany, Maged I-331, I-545
 Martel-Jantin, Bruno I-476
 Maruyama, Katsumi II-324
 Ma, Shang-Pin II-257
 Masini, Nicola I-254, I-320, I-361
 Matsunaga, Katsuya III-552, IV-497
 Mazón, Jose-Norberto I-461
 McAnally, William I-501
 Mekhedov, Ivan I-557
 Mestetskiy, Leonid I-557
 Meza, Ricardo Rafael Quintero II-241
 Milani, Alfredo IV-309
 Misra, A.K. II-273
 Misra, Sanjay II-301
 Mitrea, Delia II-215
 Mitrea, Paulina II-215
 Montrone, Silvestro I-17
 Moon, Hyun-joo III-269
 Moon, Jongbae III-258, III-269
 Mukherjee, Indira I-309
 Müller-Molina, Arnoldo José III-443,
 IV-252
 Murgante, Beniamino I-62, I-320
- Nagy, Miroslav IV-62
 Nakamura, Toru IV-236
 Namatame, Akira IV-321
 Nedoma, Jiří IV-62
 Nickerson, Bradford G. I-401
 Nicolas, Lachance-Bernard I-32
 Ninomiya, Daisuke IV-252
 Nishida, Akira II-448
 Nissen, Volker IV-346
 Niyogi, Rajdeep IV-309
 Nomura, Yoshinari II-324
 Nomura, Yusuke II-324
- Ochodková, Eliška III-472
 Ogi, Tetsuro IV-336
 O'Hara, Charles G. I-491
 Oh, Chang-Yeong III-96
 Okamoto, Kouta II-324
 Orshoven, Jos Van I-414
 Osada, Toshiaki IV-164
- Pallottelli, Simonetta IV-29
 Pannacci, Nicola IV-29
 Park, Gyung-Leen III-1, III-11, III-52,
 III-107
 Park, Min-Woo III-142, III-352
 Park, Soohong I-105
 Passeri, Francesco Luca II-422
 Pazand, Babak III-336
 Pecci, Francesco I-46
 Perchinunno, Paola I-17
 Phinitkar, Pattira IV-209
 Pirani, Fernando IV-1
 Platoš, Jan III-472
 Plumejeaud, Christine I-445
 Pontarollo, Nicola I-46
 Porceddu, Andrea I-1
 Potdar, Vidyasagar II-351, II-370,
 II-383, II-400
 Potenza, Maria Rosaria I-320
 Prastacos, Poulicos I-140
 Přečková, Petra IV-62
 Prud'homme, Julie I-445
 Purnami, Santi Wulan IV-83
- Quintero, Ricardo II-361
- Raba, Nikita II-130
 Rahayu, Wenny III-380
 Rampino, Sergio IV-1

- Ridzuan, Farida II-383, II-400
 Robinson, Ian II-44
 Röcker, Carsten IV-93
 Rodriguez, Aurelio IV-41
 Rotondo, Francesco I-283
 Ruckenbauer, Matthias IV-13

 Saft, Danilo IV-346
 Saito, Tsubasa II-60
 Sakatoku, Akira IV-164
 Salim, Flora Dilys III-483
 Sánchez, Leopoldo Z. II-241, II-361
 Sanguansintukul, Siripun III-215,
 III-419
 Sarencheh, Saeed II-370, II-400
 Sari, Anny Kartika III-380
 Scardaccione, Grazia I-62
 Schleupner, Christine I-193
 Scorza, Francesco I-62
 Selicato, Francesco I-210
 Selmane, Schehrazad IV-72
 Sen, Jaydip III-246, III-277
 Sergio, Porta I-32
 Serikawa, Seiichi II-488
 Shan, Pei-Wei II-191
 Shen, Liyong IV-482
 Shibata, Yoshitaka III-168
 Shimizu, Yoshimitsu II-139
 Shin, In-Hye III-1, III-11, III-52, III-107
 Shinohara, Takeshi III-443, IV-252
 Shiratori, Norio IV-138, IV-152, IV-164
 Shon, Minhan III-129
 Shukla, Ruchi II-273
 Şimşek, Ömer II-544
 Singh, Kuldeep IV-309
 Snapp, Robert R. I-238
 Snášel, Václav III-472
 Song, Chonghan III-21, III-31
 Song, MoonBae III-129
 Sophatsathit, Peraphon IV-209
 Sosonkina, Masha II-503
 Stankova, Elena II-130
 Stankute, Silvija I-515
 Steinhöfel, Jens I-152
 Stéphane, Joost I-32
 Suganuma, Takuo IV-138, IV-152
 Sugawara, Kenji IV-119
 Suga, Yuji IV-284
 Sugiyanta, Lipur III-305
 Suh, Soon-Tak I-262

 Suh, Woonsuk IV-385
 Sumida, Yasuaki III-552
 Sur, Sanjib III-321

 Takahashi, Hideyuki IV-138, IV-152
 Takahata, Kazuo III-168
 Takaoka, Tadao II-519
 Takashita, Taiki II-412
 Talevski, Alex II-351, II-383, II-400
 Tang, Cheng-Jen III-368
 Tangkraingkij, Preecha III-419
 Tangman, Desire Yannick II-570
 Taniar, David I-574
 Taniguchi, Hideo II-324
 Tan, Syh-Yuan IV-188
 Tasso, Sergio IV-29
 Tilio, Lucia I-320
 Timothée, Produit I-32
 Tiwari, Ashutosh IV-358
 Toi, Yutaka III-498
 Tokuhisa, Soichiro III-189
 Torre, Carmelo Maria I-17
 Trujillo, Juan I-461
 Trunfio, Giuseppe A. I-166
 Tsai, Hsin-Che III-526
 Tsai, Ming-Hsun III-511
 Tseng, Vincent S. III-458
 Tseng, Wan-Yu I-77
 Tseng, Yuh-Min IV-225
 Tu, Pu III-291
 Turner, Chris IV-358

 Uchida, Noriki III-168
 Uchiya, Takahiro IV-107
 Uddin, Mohammad Mesbah IV-199
 Uemura, Toshikazu IV-441
 Ufuktepe, Ünal IV-53
 Ukey, Nilesh IV-309
 Ukil, Arijit III-277
 Uljee, Inge I-89
 Ushijima, Kazuo IV-418
 Ushioda, Tatsuya IV-128

 van der Kwast, Johannes I-89
 Van de Voorde, Tim I-89
 Vanegas, Pablo I-414
 Villa, Ferdinando I-238
 Villarini, Mauro II-206
 Villecco, Francesco I-590

- Wang, Cheng-Long I-224
 Wang, Hao III-291
 Wang, Lian-Jun I-599
 Wang, Ping III-291
 Wang, Shuai II-463
 Wang, Tao IV-427
 Wang, Ye IV-263
 Wang, Yung-Chieh I-224
 Wolff, Markus I-346
 Wollersheim, Dennis III-380
 Won, Dongho IV-510
 Won, Kyunghoon III-118
 Wu, Bin IV-482
 Wu, Chen-Fa I-116
 Wu, Tsu-Yang IV-225
 Wu, Yu-Chieh III-458

 Xavier, Carolina Ribeiro II-475, IV-395
 Xiang, Limin IV-418

 Yamamoto, Toshihiko IV-321
 Yamawaki, Akira II-488
 Yang, Jian III-291
 Yang, Shiliang I-590
 Yang, Shiyuan II-463

 Yang, Yang II-25
 Yao, Chih-Chia III-511
 Yasuura, Hiroto IV-236
 Yeganeh, Elham Afsari II-370, II-400
 Yen, Show-Jane III-458
 Yilmaz, Buket IV-53
 Ying, Jia-Ching III-458
 Yoe, Hyun III-258
 Yokoyama, Kazutoshi II-324
 Yoon, Seokho III-21, III-31
 Yoo, Yoong-Seok II-119
 Yuasa, Fukuko II-139
 Yu, Hsiao-Hsuan I-224
 Yu, Zhi-hong IV-427

 Zain, Jasni Mohamad IV-83
 Zazueta, Liliana Vega II-241
 Zeng, Zhenbing IV-482
 Zha, Daren IV-468
 Zhai, Yuyi I-590
 Zhang, Hong I-131
 Zhao, Yaolong I-131
 Zotta, Cinzia I-320
 Zou, Zhiwen IV-296
 Zurada, Jacek M. II-110