

Ajith Abraham
Pavel Krömer
Vaclav Snášel *Editors*

Afro-European Conference for Industrial Advancement

Proceedings of the First International
Afro-European Conference for
Industrial Advancement AECIA 2014

Advances in Intelligent Systems and Computing

Volume 334

Series editor

Janusz Kacprzyk, Polish Academy of Sciences, Warsaw, Poland
e-mail: kacprzyk@ibspan.waw.pl

About this Series

The series “Advances in Intelligent Systems and Computing” contains publications on theory, applications, and design methods of Intelligent Systems and Intelligent Computing. Virtually all disciplines such as engineering, natural sciences, computer and information science, ICT, economics, business, e-commerce, environment, healthcare, life science are covered. The list of topics spans all the areas of modern intelligent systems and computing.

The publications within “Advances in Intelligent Systems and Computing” are primarily textbooks and proceedings of important conferences, symposia and congresses. They cover significant recent developments in the field, both of a foundational and applicable character. An important characteristic feature of the series is the short publication time and world-wide distribution. This permits a rapid and broad dissemination of research results.

Advisory Board

Chairman

Nikhil R. Pal, Indian Statistical Institute, Kolkata, India
e-mail: nikhil@isical.ac.in

Members

Rafael Bello, Universidad Central “Marta Abreu” de Las Villas, Santa Clara, Cuba
e-mail: rbellop@uclv.edu.cu

Emilio S. Corchado, University of Salamanca, Salamanca, Spain
e-mail: escorchado@usal.es

Hani Hagras, University of Essex, Colchester, UK
e-mail: hani@essex.ac.uk

László T. Kóczy, Széchenyi István University, Győr, Hungary
e-mail: koczy@sze.hu

Vladik Kreinovich, University of Texas at El Paso, El Paso, USA
e-mail: vladik@utep.edu

Chin-Teng Lin, National Chiao Tung University, Hsinchu, Taiwan
e-mail: ctlin@mail.nctu.edu.tw

Jie Lu, University of Technology, Sydney, Australia
e-mail: Jie.Lu@uts.edu.au

Patricia Melin, Tijuana Institute of Technology, Tijuana, Mexico
e-mail: epmelin@hafsamx.org

Nadia Nedjah, State University of Rio de Janeiro, Rio de Janeiro, Brazil
e-mail: nadia@eng.uerj.br

Ngoc Thanh Nguyen, Wroclaw University of Technology, Wroclaw, Poland
e-mail: Ngoc-Thanh.Nguyen@pwr.edu.pl

Jun Wang, The Chinese University of Hong Kong, Shatin, Hong Kong
e-mail: jwang@mae.cuhk.edu.hk

More information about this series at <http://www.springer.com/series/11156>

Ajith Abraham · Pavel Krömer
Václav Snášel
Editors

Afro-European Conference for Industrial Advancement

Proceedings of the First International
Afro-European Conference for Industrial
Advancement AECIA 2014

Editors

Ajith Abraham
Machine Intelligence Research Labs
(MIR Labs)
Scientific Network for Innovation and
Research Excellence
Auburn Washington
USA

Václav Snášel
Department of Computer Science
Faculty of Elec. Eng. & Comp. Sci.
VSB-Technical University of Ostrava
Ostrava-Poruba
Czech Republic

Pavel Krömer
Department of Computer Science
Faculty of Elec. Eng. & Comp. Sci.
VSB-Technical University of Ostrava
Ostrava-Poruba
Czech Republic

ISSN 2194-5357 ISSN 2194-5365 (electronic)
Advances in Intelligent Systems and Computing
ISBN 978-3-319-13571-7 ISBN 978-3-319-13572-4 (eBook)
DOI 10.1007/978-3-319-13572-4

Library of Congress Control Number: 2014955369

Springer Cham Heidelberg New York Dordrecht London

© Springer International Publishing Switzerland 2015

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

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

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

Springer International Publishing AG Switzerland is part of Springer Science+Business Media
(www.springer.com)

Preface

This volume of *Advances in Intelligent Systems and Computing* contains accepted papers presented at AECIA 2014, the First International Afro-European Conference for Industrial Advancement. The aim of AECIA is to bring together the foremost experts as well as excellent young researchers from Africa, Europe, and the rest of the world to disseminate latest results from various fields of engineering, information and communication technologies.

The first edition of AECIA was organized jointly by Addis Ababa Institute of Technology, Addis Ababa University, and VŠB - Technical University of Ostrava, Czech Republic. It received 70 submissions from 19 countries from Africa, Europe, North and South America, and Asia. After a thorough review process performed by foremost experts from academic institutions from 16 countries, a total of 31 papers was selected for presentation at the conference, setting the acceptance rate at less than 45 percent.

The organization of the AECIA 2014 conference was entirely voluntary. The review process required an enormous effort from the members of the International Technical Program Committee, and we would therefore like to thank all its members for their contribution to the success of this conference. We would like to express our sincere thanks to the host of AECIA 2014, Addis Ababa University, and to the publisher, Springer, for their hard work and support in organizing the conference. Finally, we would like to thank all the authors for their high quality contributions. The friendly and welcoming attitude of conference supporters and contributors made this event a success!

September 2014

Ajith Abraham
Pavel Krömer
Václav Snášel

Organization

Conference Chairs

Ajith Abraham
Václav Snášel
Daniel Kitaw
Richard Chbeir

MIR Labs, USA
VŠB - TU Ostrava, Czech Republic
Addis Ababa Institute of Technology, Ethiopia
Laboratoire LIUPPA, Université de Pau et des Pays
de l'Adour (UPPA) IUT de Bayonne, France

Program Chairs

Youakim Badr
Ivan Zelinka
Aboul Ella Hassanien

INSA-Lyon, France
VŠB - TU Ostrava, Czech Republic
Cairo University, Egypt

Publicity Co-chairs

Jan Platoš
Pavel Krömer
Beteley Tekola

VŠB - TU Ostrava, Czech Republic
University of Alberta, Canada
Addis Ababa Institute of Technology, Ethiopia

Conference Organizers

Hussein Soori
Eshetie Berhan
Yalemzewd Negash
Birhanu Beshah Abebe
Jakub Stolfa
Svatopluk Stolfa
Katerina Kasparova

VŠB - TU Ostrava, Czech Republic
Addis Ababa Institute of Technology, Ethiopia
Addis Ababa Institute of Technology, Ethiopia
Addis Ababa Institute of Technology, Ethiopia
VŠB - TU Ostrava, Czech Republic
VŠB - TU Ostrava, Czech Republic
VŠB - TU Ostrava, Czech Republic

International Program Committee

Janos Abonyi	University of Pannonia, Hungary
Ajith Abraham	MIR Labs, USA
Muhanned Alfarras	Gulf university, Bahrain
Ahmed Ali	Suez Canal University, Egypt
Youakim Badr	National Institute of Applied Sciences (INSA-Lyon), France
Addisu Bekele	Adama Science and Technology University, Ethiopia
Eshetie Berhan	Addis Ababa Institute of Technology, Ethiopia
Birhanu Beshah	Addis Ababa Institute of Technology, Ethiopia
Joo P. S. Catalo	University of Beira Interior, Portugal
Hilmi Berk Celikoglu	Technical University of Istanbul, Turkey
Ahmed El Oualkadi	ENSA Tangier, Morocco
Eid Emary	Cairo university, Egypt
Jiri Dvorsky	VŠB - TU Ostrava, Czech Republic
Belachew Gessesse	Bahir Dar Institute of Technology, Ethiopia
Mengist Hailemariam	Adama Science and Technology University, Ethiopia
Aboul Ella Hassanien	Cairo University, Egypt
Carlos Henggeler Antunes	University of Coimbra, Portugal
Dusan Husek	Academy of Sciences of the Czech Republic
Richard Chbeir	Laboratoire LIUPPA, Université de Pau et des Pays de l'Adour (UPPA) IUT de Bayonne, France
Konrad Jackowski	Wroclaw University of Technology, Poland
David Jezek	VŠB - TU Ostrava, Czech Republic
Habib M. Kammoun	REGIM-Lab., Research Groups on Intelligent Machines, University of Sfax, Egypt
Daniel Kitaw	Addis Ababa Institute of Technology, Ethiopia
Jan Kozusznik	VŠB - TU Ostrava, Czech Republic
Bartosz Krawczyk	Wroclaw University Of Technology, Poland
Pavel Kromer	University of Alberta, Canada
Jan Martinovic	VŠB - TU Ostrava, Czech Republic
Santosh Nanda	Eastern Academy of Science and Technology, Bhubneswar, Odisha, India
Eliska Ochodkova	VŠB - TU Ostrava, Czech Republic
Cyril Onwubiko	United Kingdom
Rasha Osman	Imperial College London, United Kingdom
Marcin Paprzycki	IBS PAN and WSM, Poland
Jan Platos	VŠB - TU Ostrava, Czech Republic
Sg Ponnambalam	Monash University Sunway Campus, Malaysia

Petrica Pop	North University of Baia Mare, Romania
Radu-Emil Precup	Politehnica University of Timisoara, Romania
Petr Saloun	VŠB - TU Ostrava, Czech Republic
Vaclav Snasel	VŠB - TU Ostrava, Czech Republic
Tinus Stander	University of Pretoria, South Africa
Tammo Steenhuis	Cornell University, USA
Jakub Stolfa	VŠB - TU Ostrava, Czech Republic
Svatopluk Stolfa	VŠB - TU Ostrava, Czech Republic
Thomas Stuetzle	IRIDIA, ULB, Belgium
Theo G. Swart	University of Johannesburg, Republic of South Africa
Fasil Taddese	Mekelle University, Ethiopia
Eiji Uchino	Yamaguchi University, Japan
Juan Velasquez	University of Chile, Chile
Michal Wozniak	Wroclaw University of Technology, Poland
Ivan Zelinka	VŠB - TU Ostrava, Czech Republic
Ahmed Zobaa	Brunel University, UK

Sponsoring Institutions

Addis Ababa Institute of Technology, Addis Ababa University, Ethiopia

Contents

Feature Subset Selection Approach by Gray-Wolf Optimization	1
<i>E. Emary, Hossam M. Zawbaa, Crina Grosan, Abul Ella Hassenian</i>	
Optimization of Wind Direction Distribution Parameters Using Particle Swarm Optimization	15
<i>Jana Heckenbergerova, Petr Musilek, Pavel Krömer</i>	
A Perspective of the Cellular Network of the Future: Cloud-RAN	27
<i>Santhi Kumaran</i>	
Reliable Data Transmission for Multi-source Multi-sink Wireless Sensor Networks	43
<i>Kassahun Tamir, Menore Tekeba</i>	
Extraction of Fetal ECG from Abdominal ECG and Heart Rate Variability Analysis	65
<i>Gizeaddis Lamesgin, Yonas Kassaw, Dawit Assefa</i>	
Addition of Static Aspects to the Intuitive Mapping of UML Activity Diagram to CPN	77
<i>Jan Czopik, Michael Alexander Košinár, Jakub Štolfa, Svatopluk Štolfa</i>	
Flash Assisted Segmented Bloom Filter for Deduplication	87
<i>Girum Dagnaw, Amare Teferi, Eshetie Berhan</i>	
A Train Run Simulation with the Aid of the EMTP – ATP Programme	99
<i>Maroš Ďurica</i>	
Comparative Assessment of Temperature Based ANN and Angstrom Type Models for Predicting Global Solar Radiation	109
<i>Darlington Ihunanyachukwu Egeonu, Howard Okezie Njoku, Patrick Nwosa Okolo, Samuel Ogbonna Enibe</i>	

Position Control and Tracking of Ball and Plate System Using Fuzzy Sliding Mode Controller	123
<i>Andinet Negash, Nagendra P. Singh</i>	
A Simulation of Project Completion Probability Using Different Probability Distribution Functions	133
<i>Erimas Tesfaye, Kidist Girma, Eshetie Berhan, Birhanu Beshah</i>	
Dynamic Simulation of T-Track: Under Moving Loads	147
<i>Mequanent Mulugeta</i>	
Utilizing Text Similarity Measurement for Data Compression to Detect Plagiarism in Czech	163
<i>Hussein Soori, Michal Prilepok, Jan Platos, Václav Snášel</i>	
AJAX Speed Up vs. JSP in Portal – Case Study	173
<i>David Ježek, Radek Liebzeit</i>	
Intelligent Decision Support for Real Time Health Care Monitoring System	183
<i>Abdelhamid Salih Mohamed Salih, Ajith Abraham</i>	
An Efficient Segmentation Algorithm for Arabic Handwritten Characters Recognition System	193
<i>Mohamed A. Ali</i>	
Graph Drawing Using Dimension Reduction Methods	205
<i>Tomáš Buriánek, Lukáš Zaorálek, Václav Snášel, Tomáš Peterek</i>	
Design of a Single Stage Thermoelectric Power Generator Module with Specific Application on the Automotive Industry	215
<i>Yidnekachew Messele, Eyerusalem Yilma, Rahma Nasser</i>	
Ethiopian Livestock Husbandry Cluster Identification Using FUZZY-AHP Approach	233
<i>Netsanet Jote, Birhanu Beshah, Daniel Kitaw</i>	
Thermal Analysis, Design and Experimental Investigation of Parabolic Trough Solar Collector	245
<i>Yidnekachew Messele, Ababayehu Assefa</i>	
Performance Improvement by Scheduling Techniques: A Case of Leather Industry Development Institute	261
<i>Abduletif Habib, Kassu Jilcha, Eshetie Berhan</i>	
Response Time Reduction in the Leather Products Manufacturing Industry Using Arena Simulation Method	271
<i>Haftu Hailu, Kassu Jilcha, Eshetie Birhan</i>	

Lead Time Prediction Using Simulation in Leather Shoe Manufacturing . . .	283
<i>Hermela Solomon, Kassu Jilcha, Eshetie Berhan</i>	
Ensemble Neurocomputing Based Oil Price Prediction	293
<i>Lubna A. Gabralla, Hela Mahersia, Ajith Abraham</i>	
Internet of Things Communication Reference Model and Traffic Engineer System (TES)	303
<i>Adel H. Alhamedi, Hamoud M. Aldosari, Václav Snášel, Ajith Abraham</i>	
Modeling Cloud Computing Risk Assessment Using Machine Learning	315
<i>Nada Ahmed, Ajith Abraham</i>	
Adaptation of Turtle Graphics Method for Visualization of the Process Execution	327
<i>Jakub Štolfa, Svatopluk Štolfa, Martin Kopka, Václav Snášel</i>	
NFL Results Predictor as a Smart Mobile Application	335
<i>Petr Kavrda, Ondrej Berger, Ondrej Krejcar</i>	
Security Issues of Mobile Application Using Cloud Computing	347
<i>Richard Cimler, Jan Matyska, Ladislav Balík, Josef Horalek, Vladimir Sobeslav</i>	
SIFT-Based Arabic Sign Language Recognition System	359
<i>Alaa Tharwat, Tarek Gaber, Aboul Ella Hassanien, M.K. Shahin, Basma Refaat</i>	
Stock Market Forecasting Using LASSO Linear Regression Model	371
<i>Sanjiban Sekhar Roy, Dishant Mittal, Avik Basu, Ajith Abraham</i>	
Author Index	383

Feature Subset Selection Approach by Gray-Wolf Optimization

E. Emary^{1,4}, Hossam M. Zawbaa^{2,3,4}, Crina Grosan², and Abul Ella Hassenian^{1,4}

¹ Faculty of Computers and Information, Cairo University, Egypt

² Faculty of Mathematics and Computer Science, Babes-Bolyai University, Romania

³ Faculty of Computers and Information, Beni-Suef University, Egypt

⁴ Scientific Research Group in Egypt (SRGE), Egypt

<http://www.egyptscience.net>

Abstract. Feature selection algorithm explores the data to eliminate noisy, irrelevant, redundant data, and simultaneously optimize the classification performance. In this paper, a classification accuracy-based fitness function is proposed by gray-wolf optimizer to find optimal feature subset. Gray-wolf optimizer is a new evolutionary computation technique which mimics the leadership hierarchy and hunting mechanism of gray wolves in nature. The aim of the gray wolf optimization is find optimal regions of the complex search space through the interaction of individuals in the population. Compared with particle swarm optimization (PSP) and Genetic Algorithms (GA) over a set of UCI machine learning data repository, the proposed approach proves better performance in both classification accuracy and feature size reduction. Moreover, the gray wolf optimization approach proves much robustness against initialization in comparison with PSO and GA optimizers.

Keywords: Gray-wolf Optimization, feature selection, evolutionary computation.

1 Introduction

Feature selection algorithm explores the data to eliminate noisy, irrelevant, redundant data, and simultaneously optimize the classification performance. Feature selection is one of the most important stage in data mining, multimedia information retrieval, pattern classification, and machine learning applications, which can influence the classification accuracy rate [1],[2].

The main purpose of feature selection is to determine a minimal feature subset from a problem domain while retaining a suitably high accuracy in representing the original features [3]. In real world problems, feature selection is a must due to the abundance of noisy, misleading or irrelevant features [4]. By removing these factors, learning from data techniques can useful greatly. The motivation of feature selection in data mining, machine learning and pattern recognition is to reduce the dimensionality of feature space, improve the predictive accuracy of a classification algorithm, and develop the visualization and the comprehensibility of the induced concepts [5].

Genetic Algorithm (GA), Ant Colony Optimization (ACO), and Particle Swarm Optimization (PSO) are popular meta-heuristic optimization techniques. The Grey Wolf

Optimizer (GWO) is a new optimization algorithm which simulate the grey wolves leadership and hunting manner in nature. These techniques have been inspired by simple concepts. The inspirations are typically related to physical phenomena, animals behaviors, or evolutionary concepts [6]. In recent years, a lot of feature selection methods have been proposed. There are two key issues in structure a feature selection method: search strategies and evaluating measures. With respect to search strategies, complete , heuristic [7] , random [8] [9] strategies were proposed. And with respect to evaluating measures, these methods can be nearly divided into two classes: classification [10],[11],[12] and classification independent [13],[14],[15]. The previous employs a learning algorithm to evaluate the quality of selected features based on the classification accuracies or contribution to the classification boundary, such as the so-called wrapper method [10] and weight based algorithms [16],[17]. While the latter constructs a classifier independent measure to evaluate the importance of features, such as inter-class distance[13] mutual information[18], dependence measure[14] and consistency measure [15].

In recent years, a lot of feature selection methods have been proposed. There are two key issues in structure a feature selection method: search strategies and evaluating measures. With respect to search strategies, complete , heuristic [7] , random [8] [9] strategies were proposed. And with respect to evaluating measures, these methods can be nearly divided into two classes: classification [10],[11],[12] and classification independent [13],[14],[15]. The previous employs a learning algorithm to evaluate the quality of selected features based on the classification accuracies or contribution to the classification boundary, such as the so-called wrapper method [10] and weight based algorithms [16],[17] . While the latter constructs a classifier independent measure to evaluate the importance of features, such as inter-class distance[13] mutual information[18], dependence measure[14] and consistency measure [15].

In [22], the population of particles split into a set of interacting swarms. These interacting swarms applied the simple competition method. The winner is the swarm which has a best fitness value. The loser is eject and re-initialized in the search space, otherwise the winner remains. In [23], the swarm population divided into sub-populations species based on their similarity. Then, the repeated particles are removed when particles are identified as having the same fitness. After destroying the repeated ones, the new particles are added randomly until its size is resumed to its initial size.

In [24], the Bat Algorithm (BA) based on type of the sonar, which named echolocation behavior. The micro-bats have the capability of echolocation which attracting these bats can find their prey and discriminate different types of insects even in complete darkness.

In this paper, a classification accuracy-based fitness function is proposed by gray-wolf optimizer to find optimal feature subset. We compare Grey Wolf Optimizer (GWO) algorithm against Particle Swarm Optimization (PSO) and Genetic Algorithm (GA) algorithms for feature selection by applying three different initialization methods and eight different datasets. The results reveal that the GWO resulted in a higher accuracy compared to the other two optimization algorithms.

The rest of this paper is organized as follows: Section 2 presents basics of the gray wolf optimization. Section IV presents the details of the proposed system. In section V, there are experimental results and result analysis. Finally in Section VI, conclusions and future work are presented.

2 Preliminaries

2.1 Gray Wolf Optimization

Gray wolf optimization is presented in the following subsections based on the work in [25].

Inspiration. Grey wolves are considered as apex predators, meaning that they are at the top of the food chain. Grey wolves mostly prefer to live in a pack. The group size is 5-12 on average. They have a very strict social dominant hierarchy. The leaders are a male and a female, called *alpha*. The alpha is mostly responsible for making decisions about hunting, sleeping place, time to wake, and so on. The *alpha* decisions are dictated to the pack. The second level in the hierarchy of grey wolves is beta. The betas are subordinate wolves that help the alpha in decision-making or other pack activities. The beta wolf can be either male or female, and he/she is probably the best candidate to be the alpha in case one of the alpha wolves passes away or becomes very old. The lowest ranking grey wolf is omega. The omega plays the role of scapegoat. Omega wolves always have to submit to all the other dominant wolves. They are the last wolves that are allowed to eat. The fourth class is called subordinate (or delta in some references). Delta wolves have to submit to alphas and betas, but they dominate the omega. *Scouts*, *sentinels*, *elders*, *hunters*, and *caretakers* belong to this category. *Scouts* are responsible for watching the boundaries of the territory and warning the pack in case of any danger. *Sentinels* protect and guarantee the safety of the pack. *Elders* are the experienced wolves who used to be alpha or beta. *Hunters* help the alphas and betas when hunting prey and providing food for the pack. Finally, the *caretakers* are responsible for caring for the weak, ill, and wounded wolves in the pack.

Mathematical Model. In the mathematical model for the GWO the fittest solution is called the alpha (α). The second and third best solutions are named beta (β) and delta (δ) respectively. The rest of the candidate solutions are assumed to be omega (ω). The hunting is guided by α , β , and δ and the ω follow these three candidates. In order for the pack to hunt a prey they first encircling it. In order to mathematically model encircling behavior the following equations are used 1:

$$\vec{X}(t+1) = \vec{X}_p(t) + \vec{A} \cdot \vec{D} \quad (1)$$

where \vec{D} is as defined in 2 and t is the iteration number, \vec{A} , \vec{C} are coefficient vectors, \vec{X}_p is the prey position and \vec{X} is the gray wolf position.

$$\vec{D} = |\vec{C} \cdot \vec{X}_p(t) - \vec{X}(t)| \quad (2)$$

The \vec{A} , \vec{C} vectors are calculated as in equations 3 and 4

$$\vec{A} = 2\vec{A} \cdot r_1 - \vec{a} \quad (3)$$

$$\vec{C} = 2\vec{r}_2 \quad (4)$$

where components of \vec{a} are linearly decreased from 2 to 0 over the course of iterations and r_1, r_2 are random vectors in $[0, 1]$. The hunt is usually guided by the alpha. The beta and delta might also participate in hunting occasionally. In order to mathematically simulate the hunting behavior of grey wolves, the alpha (best candidate solution) beta, and delta are assumed to have better knowledge about the potential location of prey. The first three best solutions obtained so far and oblige the other search agents (including the omegas) to update their positions according to the position of the best search agents. So, the updating for the wolves positions is as in equations 5,6,7.

$$\vec{D}_\alpha = |\vec{C}_1 \cdot \vec{X}_\alpha - \vec{X}|, \vec{D}_\beta = |\vec{C}_2 \cdot \vec{X}_\beta - \vec{X}|, \vec{D}_\delta = |\vec{C}_3 \cdot \vec{X}_\delta - \vec{X}| \quad (5)$$

$$\vec{X}_1 = |\vec{X}_\alpha - \vec{A}_1 \cdot \vec{D}_\alpha|, \vec{X}_2 = |\vec{X}_\beta - \vec{A}_2 \cdot \vec{D}_\beta|, \vec{X}_3 = |\vec{X}_\delta - \vec{A}_3 \cdot \vec{D}_\delta| \quad (6)$$

$$\vec{X}(t+1) = \frac{\vec{X}_1 + \vec{X}_2 + \vec{X}_3}{3} \quad (7)$$

A final note about the GWO is the updating of the parameter \vec{a} that controls the tradeoff between exploitation and exploration. The parameter \vec{a} is linearly updated in each iteration to range from 2 to 0 according to the equation 8:

$$\vec{a} = 2 - t \cdot \frac{2}{Max_{iter}} \quad (8)$$

where t is the iteration number and Max_{iter} is the total number of iteration allowed for the optimization.

3 The Proposed Algorithm

In this section, we present the proposed GWO optimizer based on K-nearest neighbor for feature selection; see figure 1. We used the principles of gray wolf optimization for the optimal feature selection problem. Each feature subset can be seen as a position in such a space. If there are N total features, then there will be 2^N different feature subset, different from each other in the length and features included in each subset. The optimal position is the subset with *least length* and *highest classification accuracy*. We used gray wolf optimization for selecting the optimal feature set. Eventually, they should converge on good, possibly optimal, positions. The GWO makes iterations of exploration of new regions in the feature space and exploiting solution until reaching near-optimal solution.

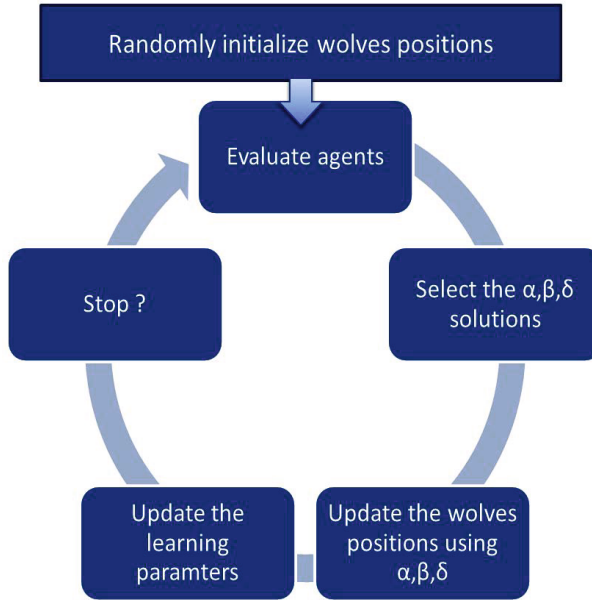


Fig. 1. The overall feature selection algorithm

The solution space in here represents all possible selections of features and hence bat positions represents binary selection of feature sets. Each feature is considered as an individual dimension ranging from -2 to 2. To decide if a feature will be selected or not its position value will be threshold with a constant threshold. The used fitness function is the classification accuracy for k-nearest neighbors (KNN) classifier on the validation set. Each individual data set is divided into three equal subsets namely training, validation and test portions. The training set and validation set are used inside the fitness function to evaluate the selection classification accuracy while the test set is used in the end of optimization to evaluate the final selection classification performance.

We made use of two fitness functions in gray-wolf optimization (GWO) for feature selection, which are *KNN*, and *KNN_{size}* resembling the well-known forward selection. Forward selection starts with an empty feature set (no features) and searches for a feature subset(s) with one feature by selecting the feature that achieves the highest classification performance. Then the algorithm selects another feature from the candidate features to add to *S*. Feature *i* is selected if adding *i* to *S* achieves the largest improvement in classification accuracy. While, backward selection starts with all the available features, then candidate features are sequentially removed from the feature subset until the further removal of any feature does not increase the classification performance. Small initialization resembles forward selection, large initialization motivated by backward selection and mixed initialization aiming to take the advantages of forward and backward selection to avoid their disadvantages.

4 Experimental Results and Discussion

4.1 Data Sets and Parameters Setting

Table 1 summarizes the 8 used data set for further experiments. The data set are drawn from the UCI data repository [27]. The data is divided into 3 equal parts one for *training*, the second part is for *validation* and the third part is for *testing*. We implement the GWO feature selection algorithms in MatLab R2009a. The computer used to get results is Intel (R), 2.1 GHz CPU; 2 MB RAM and the system is Windows 7 Professional. The parameter setting for the GWO algorithm is outlined in table 2. Same number of agents and same number of iterations are used for GA and PSO.

Table 1. Description of the data sets used in experiments

Dataset	No. of features	No. of samples
Lymphography	18	148
Zoo	16	101
Vote	16	300
Breastcancer	9	699
M-of-N	13	1000
Exactly	13	1000
Exactly2	13	1000
Tic-tac-toe	9	958

Table 2. Parameter setting for gray-wolf optimization

parameter	value(s)
No of wolves	5
No of iterations	100
problem dimension	same as number of features in any given database
Search domain	[0 1]

4.2 Results and Discussion

Four scenarios has been considered when we evaluate the proposed approach. They are: (1) **Scenario 1:** GWO, GA, and PSO features selection techniques using *normal* initialization, (2) **Scenario 2:** GWO, GA, and PSO features selection techniques using *large* initialization, (3) **Scenario 3:** GWO, GA, and PSO features selection techniques using *mixed* initialization, and (4) **Scenario 4:** GWO, GA, and PSO features selection techniques using *small* initialization.

Tables 3, 4, 5, and 6 are showing the performance of GWO, GA, PSO algorithms on the different eight data sets. Every algorithm is applied for 5 times on every data set to be sure about the algorithm robustness and we display the average result of all solutions. Gray wolf optimization (GWO) algorithm achieves high accuracy with the different data sets and initialization methods as showing in figures 2, 3, 4, and 5.

This demonstrates that GWO shows a good balance between exploration and exploitation that results in high local optima avoidance. This superior capability is due to the adaptive value of A . As mentioned above, half of the iterations are devoted to exploration ($|A| \geq 1$) and the rest to exploitation ($|A| < 1$). This mechanism assists GWO to provide very good exploration, local minima avoidance, and exploitation simultaneously.

Table 3. Normal initialization results for different datasets

Dataset No.	GWO	GA	PSO
1	0.980952	0.976190	0.985714
2	0.733333	0.970000	1.000000
3	0.773333	0.766667	0.740000
4	0.863636	0.886364	0.886364
5	1.000000	0.943333	0.946667
6	0.797909	0.797909	0.808362
7	0.977778	0.966667	0.944444
8	1.000000	0.933333	1.000000

Table 4. Large initialization results for different datasets

Dataset No.	GWO	GA	PSO
1	0.976190	0.990476	0.985714
2	1.000000	0.760000	0.723333
3	0.773333	0.756667	0.780000
4	0.886364	0.863636	0.795455
5	1.000000	0.920000	0.883333
6	0.839721	0.843206	0.832753
7	0.966667	0.966667	0.966667
8	1.000000	0.966667	0.933333

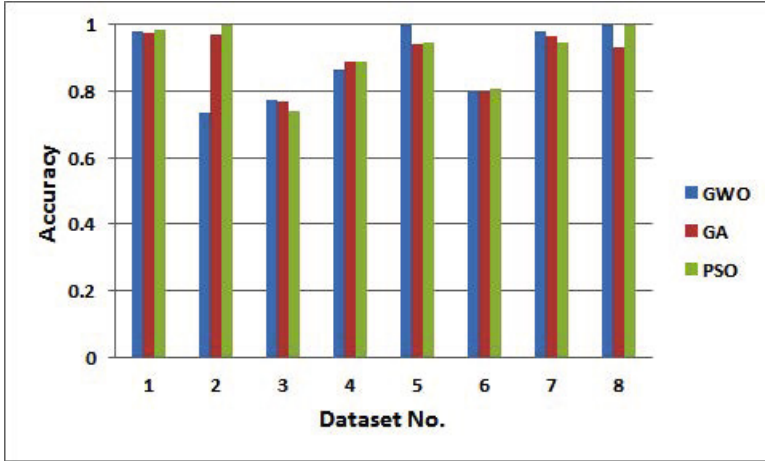


Fig. 2. Comparison of normal initialization results for GWO, GA, PSO with different datasets

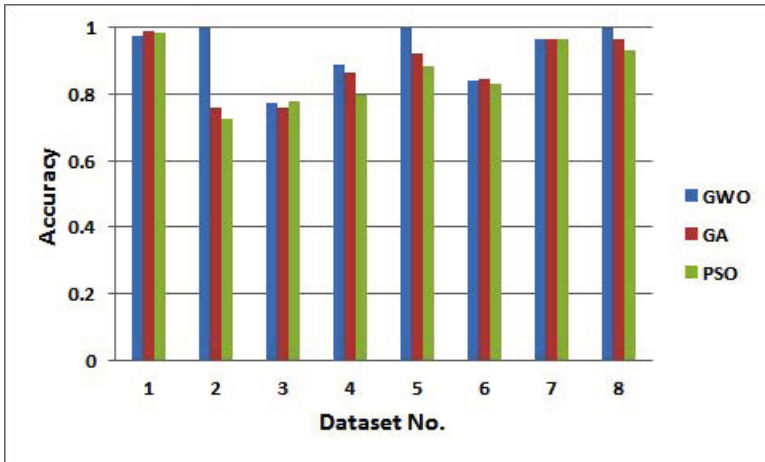
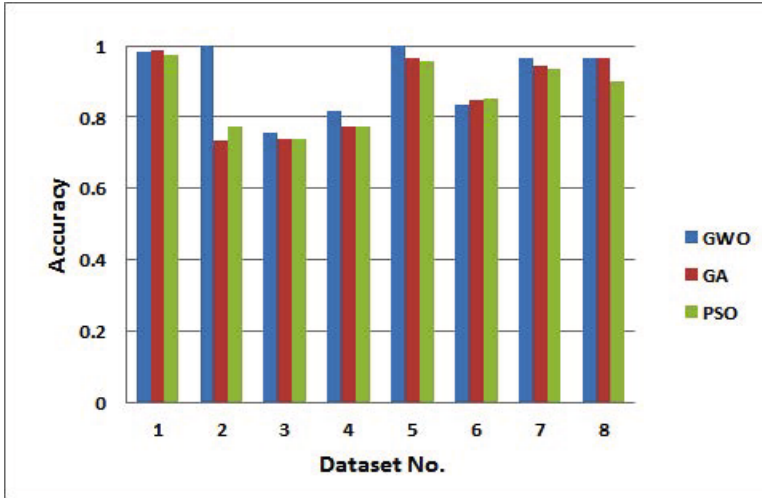


Fig. 3. Comparison of large initialization results for GWO, GA, PSO with different datasets

Figures 6 (a), (b), (c), and (d) present the standard deviation for the obtained fitness functions after running the each optimizer for 5 runs. The obtained standard deviation for GWO is much less than the obtained standard deviation for the PSO and GA which can be considered a prove for algorithm robustness regardless of the initialization method. SO, GWO always converge to the optimal solution or near optimal one regardless of its initialization method. GWO has abrupt changes in the movement of search agents over the initial steps of optimization. This assists a meta-heuristic to explore

Table 5. Mixed initialization results for different datasets

Dataset No.	GWO	GA	PSO
1	0.980952	0.985714	0.976190
2	1.000000	0.733333	0.773333
3	0.756667	0.736667	0.740000
4	0.818182	0.772727	0.772727
5	1.000000	0.966667	0.956667
6	0.832753	0.846690	0.850174
7	0.966667	0.944444	0.933333
8	0.966667	0.966667	0.900000

**Fig. 4.** Comparison curve of mixed initialization results for GWO, GA, PSO with different datasets

the search space extensively. Then, these changes should be reduced to emphasize exploitation at the end of optimization. In order to observe the convergence behavior of the GWO algorithm. The exploration power of GWO allow it to tolerate for initial solution as it quickly explore many regions and select the promising ones for further search.

Table 6. Small initialization results for different datasets

Dataset No.	GWO	GA	PSO
1	0.980952	0.966667	0.985714
2	1	0.86	0.753333
3	0.743333	0.773333	0.726667
4	0.863636	0.840909	0.863636
5	1	0.986667	0.94
6	0.850174	0.860627	0.84669
7	0.977778	0.955556	0.944444
8	0.966667	0.933333	0.966667

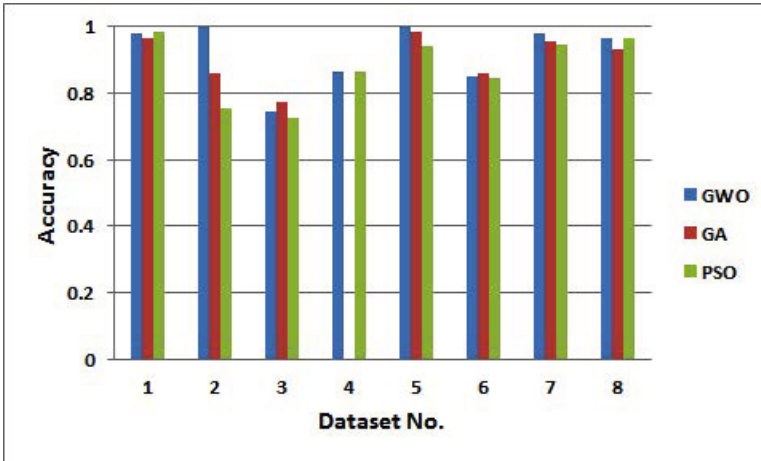


Fig. 5. Comparison curve of small initialization results for GWO, GA, PSO with different datasets

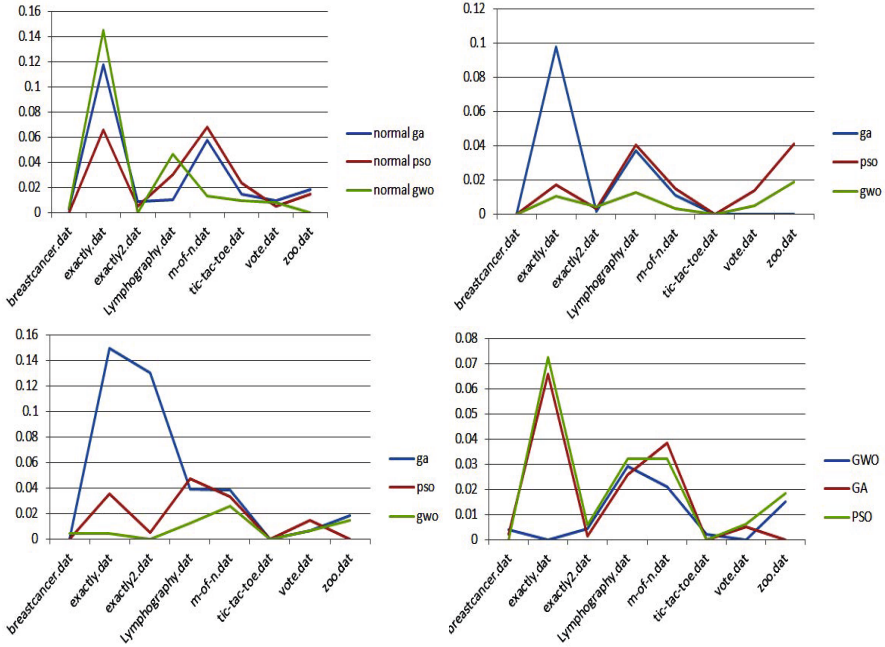


Fig. 6. (a) Standard deviation obtained for running the different methods on the different data sets [initialization normal], (b) Standard deviation obtained for running the different methods on the different data sets [initialization large] (c) Standard deviation obtained for running the different methods on the different data sets [initialization mixed] and (d) Standard deviation obtained for running the different methods on the different data sets [initialization small]

5 Conclusions and Future Work

In this paper, a system for feature selection based on intelligent search of gray wolf optimization has been proposed. Compared with PSO and GA over a set of UCI machine learning data repository, GWO proves much better performance as well as its proves much robustness and fast convergence speed. Moreover, the gray wolf optimization approach proves much robustness against initialization in comparison with PSO and GA optimizers. Other improvement to our work may involve applying some other feature selection algorithms and different fitness functions in the future which are expected to further enhance the results.

Acknowledgment. This work was partially supported by the IPROC/Marie Curie initial training network, funded through the People Programme (Marie Curie Actions) of the European Union’s Seventh Framework Programme FP7/2007-2013/ under REA grant agreement No. 316555. This fund only apply for two co-authors (Hossam M. Zawbaa and Crina Grosan).

References

1. Yang, C.-H., Tu, C.-J., Chang, J.-Y., Liu, H.-H., Ko, P.-C.: Dimensionality Reduction using GA-PSO. In: Proceedings of the Joint Conference on Information Sciences (JCIS), October 8-11. Atlantis Press, Kaohsiung (2006)
2. Cannas, L.M.: A framework for feature selection in high-dimensional domains. Ph.D. Thesis, University of Cagliari (2012)
3. Dash, M., Liu, H.: Feature selection for Classification. *Intelligent Data Analysis* 1(3), 131–156 (1997)
4. Jensen, R.: Combining rough and fuzzy sets for feature selection. Ph.D. Thesis, University of Edinburgh (2005)
5. Liu, H., Motoda, H.: *Feature Selection for Knowledge Discovery and Data Mining*. Kluwer, Boston (1998)
6. Mirjalili, S., Mirjalili, S.M., Lewis, A.: Grey Wolf Optimizer. *Adv. Eng. Softw.* 69, 46–61 (2014)
7. Zhong, N., Dong, J.Z.: Using rough sets with heuristics for feature selection. *J. Intell. Inform. Systems* 16, 199–214 (2001)
8. Raymer, M.L., Punch, W.E., Goodman, E.D., et al.: Dimensionality reduction using genetic algorithms. *IEEE Trans. Evol. Comput.* 4(2), 164–171 (2000)
9. Lai, C., Reinders, M.J.T., Wessels, L.: Random subspace method for multivariate feature selection. *Pattern Recognition Lett.* 27, 1067–1076 (2006)
10. Kohavi, R.: Feature subset selection using the wrapper method, Overfitting and dynamic search space topology. In: Proc. AAAI Fall Symposium on Relevance, pp. 109–113 (1994)
11. Gasca, E., Sanchez, J.S., Alonso, R.: Eliminating redundancy and irrelevance using a new MLP-based feature selection method. *Pattern Recognition* 39(2), 313–315 (2006)
12. Neumann, J., Schnorr, C., Steidl, G.: Combined SVM-based feature selection and classification. *Machine Learning* 61, 129–150 (2005)
13. Kira, K., Rendell, L.A.: The feature selection problem: Traditional methods and a new algorithm. In: Proc. AAAI 1992, San Jose, CA, pp. 129–134 (1992)
14. Modrzejewski, M.: Feature selection using rough sets theory. In: Proceedings of the European Conference on Machine Learning, Vienna, Austria, pp. 213–226 (1993)
15. Dash, M., Liu, H.: Consistency-based search in feature selection. *Artif. Intell.* 151, 155–176 (2003)
16. Guyon, I., Weston, J., Barnhill, S., Vapnik, V.: Gene selection for cancer classification using support vector machines. *Machine Learning* 46, 389–422 (2002)
17. Xie, Z.-X., Hu, Q.-H., Yu, D.-R.: Improved feature selection algorithm based on SVM and correlation. In: Wang, J., Yi, Z., Žurada, J.M., Lu, B.-L., Yin, H. (eds.) ISNN 2006. LNCS, vol. 3971, pp. 1373–1380. Springer, Heidelberg (2006)
18. Yao, Y.Y.: Information-theoretic measures for knowledge discovery and data mining. In: Karmeshu (ed.) *Entropy Measures, Maximum Entropy and Emerging Applications*, pp. 115–136. Springer, Berlin (2003)
19. Deogun, J.S., Raghavan, V.V., Sever, H.: Rough set based classification methods and extended decision tables. In: Proc. of the Int. Workshop on Rough Sets and Soft Computing, pp. 302–309 (1994)
20. Zhang, M., Yao, J.T.: A rough sets based approach to feature selection. In: IEEE Annual Meeting of the Fuzzy Information, Processing NAFIPS 2004, June 27-30, vol. 1, pp. 434–439 (2004)
21. Hu, X.: Knowledge discovery in databases: an attribute-oriented rough set approach. Ph.D thesis, University of Regina, Canada (1995)

22. Blackwell, T., Branke, J.: Multiswarms, “exclusion, and anti-convergence in dynamic environments”. *IEEE Transactions on Evolutionary Computation* 10, 459–472 (2006)
23. Parrott, D., Li, X.D.: Locating and tracking multiple dynamic optima by a particle swarm model using speciation. *IEEE Transactions on Evolutionary Computation* 10, 440–458 (2006)
24. Yang, X.-S.: A New Metaheuristic Bat-Inspired Algorithm. In: González, J.R., Pelta, D.A., Cruz, C., Terrazas, G., Krasnogor, N. (eds.) *NICSO 2010. SCI*, vol. 284, pp. 65–74. Springer, Heidelberg (2010)
25. Mirjalili, S., Mirjalili, S.M., Lewis, A.: Grey Wolf Optimizer. *Advances in Engineering Software* 69, 46–61 (2014)
26. Xue, B., Zhang, M., Browne, W.N.: Particle swarm optimisation for feature selection in classification: Novel initialisation and updating mechanisms. *Applied Soft Computing* 18, 261–276 (2014)
27. Bache, K., Lichman, M.: *UCI Machine Learning Repository*. University of California, School of Information and Computer Science, Irvine (2013), <http://archive.ics.uci.edu/ml>

Optimization of Wind Direction Distribution Parameters Using Particle Swarm Optimization

Jana Heckenbergerova¹, Petr Musilek^{2,3}, and Pavel Krömer^{2,3}

¹ Department of Mathematics and Physics,
Faculty of Electrical Engineering and Information Science,
University of Pardubice, Czech Republic

`jana.heckenbergerova@upce.cz`

² Department of Computer Science
VŠB Technical University of Ostrava
Ostrava, Czech Republic

`{petr.musilek,pavel.kromer}@vsb.cz`

³ Department of Electrical and Computer Engineering,
University of Alberta, Edmonton AB T6G 2V4, Canada

`{pmusilek,pavel.kromer}@ualberta.ca`

Abstract. Data describing various natural and industrial phenomena can be modeled by directional statistical distributions. In the field of energy, wind direction and wind speed are the most important variables for wind energy generation, integration, and management. This work proposes and evaluates a new method for accurate estimation of wind direction distribution parameters utilizing the well-known Particle Swarm Optimization algorithm. It is used to optimize the parameters of a site-specific wind direction distribution model realized as a finite mixture of circular normal von Mises statistical distributions. The evaluation of the proposed algorithm is carried out using a data set describing annual wind direction on two distinct locations. Experimental results show that the proposed method is able to find good model parameters corresponding to input data.

Keywords: von Mises distribution, circular normal distribution, finite mixture model, circular data statistics, Particle Swarm Optimization, experiments

1 Introduction

Directional data plays an important role in many scientific and applied areas. It has been used to model a wide range of phenomena from diverse areas including medical science (sudden infant death syndrome [1]), image analysis and video surveillance (detection of abnormal behaviour [2]), gas load forecasting [3], and various applications of wind properties modelling (pollutant source identification [4] annual wind direction modelling [5–8], wind energy potential assessment [9]).

In the field of energy, wind direction and wind speed are the most important variables for wind energy generation, integration, and management. Accurate modelling of wind speed and direction is necessary for predicting the volume of generated wind energy. Such predictions are crucial for effective and safe operations of stochastic renewable energy sources like wind turbines and wind farms. In addition, wind properties affect the current-carrying capacity of overhead power transmission lines and determine wind energy potential of a specific site. The estimation of wind energy generation is an essential part of green energy and gains on importance in both developed and developing regions [10, 11]. The analysis of wind power potential of a specific site plays an important role when selecting the location of new wind energy facilities, and has impact on both, price of generated energy and return of investment [12].

Wind direction modelling is also important for power delivery. Dynamic thermal rating (DTR) of overhead power transmission lines estimates their capacity with respect to the actual operating conditions to increase system throughput [13, 14]. It is known that wind parameters have significant impact on thermal rating so the knowledge of wind direction probabilities on a specific location contributes to accurate estimation of real power transmission limits and capacity [15].

To assess the wind-related properties of a specific site, various models analyzing wind characteristics can be used [5]. Statistical wind modelling aims at identifying the probability density function (*pdf*) [16, 17] describing site-specific wind properties (speed, direction, speed-direction tuple [18]). Statistical analysis of directional wind speed is usually performed using discrete models based on histograms of annual wind data divided into a small number of groups (sectors) [5].

It is known that directional data cannot be analyzed by standard linear statistical methods. Instead, periodical and directional statistical distributions such as uniform, cardioid, wrapped normal, and von Mises distribution, can be used for accurate circular data modelling [1, 2, 5]. A popular approach to circular data *pdf* modelling is the use of a finite mixture of von Mises distributions [1, 2, 5, 7, 8].

This work introduces a new meta-heuristic approach to the estimation of parameters of such finite mixture of von Mises distributions based on the Particle Swarm Optimization (PSO) algorithm. The proposed method is general and can be used to model circular data from any application domain. In this study, however, it is used for wind direction modelling and the experiments conducted in scope of this study performed the search for accurate parameters of annual wind direction *pdfs* for two locations in Canada. The quality of the obtained *pdfs* is measured by Pearson's chi-squared goodness-of-fit function which was also used as the fitness criterion.

The rest of this paper is organized in the following way. Section 2 summarizes relevant recent work on wind direction estimation and forecasting. Section 3 summarizes the principles of circular data statistics. Section 4 outlines the basics of Particle Swarm Optimization and details the proposed approach

to wind direction distribution parameters optimization. Computational experiments evaluating the proposed method are shown in section 5 and conclusions are drawn in section 6.

2 Related Work

There are many approaches to wind properties modelling and analysis. Most of them aim at long and short-term modelling of wind speed or joint modelling of both, wind speed and direction [6, 18–21]. However, studies focusing on wind direction analysis, modelling, and estimation can be found as well [4, 5, 7, 8, 22]. The algorithms used for wind direction modelling include detection of local gradients on synthetic aperture radar (SAR) images [22], fitting of a finite mixture of von Mises distributions by the least square method [5] or by the expectation-minimization algorithm [8], autoregressive moving average (ARMA) based approaches [18], and various meta-heuristic and bio-inspired methods [4, 7, 19, 20].

Wind direction determination was a part of complex method for pollutant source identification proposed by Allen et al. [4]. The main aim of the research was the identification of surface wind directions and other detailed characteristics of pollutant sources. It used Genetic Algorithms (GA) to match pollutant dispersion model parameters with real data recorded by field sensors.

The use of genetic fuzzy systems for local wind conditions modelling was investigated by de la Rosa et al. [19]. A fuzzy system evolved by the GA was used to adjust a statistical regional wind properties model to better match local conditions on a particular site. Heckenbergerova et al. [7] used a simple GA to find the parameters of finite mixture of von Mises distributions to model annual wind direction probability on a specific location.

Four models of wind speed and direction based on adaptive neuro-fuzzy inference systems (ANFIS) and different probability distributions (Weibull, Frechet, Gumbel, and joint distribution) were described by Shamshirband et al. [20]. The study concluded that ANFIS can learn the wind speed and direction *pdf* well.

In this work, a popular bio-inspired real-parameter optimization meta-heuristic algorithm called Particle Swarm Optimization is used to optimize parameters of finite mixture of von Mises distributions describing the *pdf* of annual wind directions on two sites in Canada. The Particle Swarm Optimization algorithm was selected because it is, in contrast to e.g. GAs, a method designed for real-parameter optimization and has an excellent record of real-world applications [23, 24].

3 Circular Data Statistics

In this paper, finite Mixture of simple von Mises distributions (MvM) has been selected for representing multimodal directional data. Simple von Mises distribution (SvM), defined for random circular variable θ , has two parameters:

μ represents the prevailing wind direction, and κ (concentration parameter) indicates the variance around the mean. This distribution is also known as circular normal distribution. Its probability density function (SvM-*pdf*) is given by

$$f(\theta; \mu, \kappa) = \frac{1}{2\pi I_0(\kappa)} e^{\kappa \cos(\theta - \mu)}, \quad (1)$$

where $\kappa \geq 0$, $0 \leq \mu \leq 2\pi$, $0 \leq \theta \leq 2\pi$ and $I_0(\kappa)$ is the modified Bessel function of the first kind and order zero

$$I_0(\kappa) = \frac{1}{\sqrt{2\pi}} \int_0^{2\pi} e^{\kappa \cos(\theta)} d\theta = \sum_{k=0}^{\infty} \frac{1}{(k!)^2} \left(\frac{\kappa}{2}\right)^{2k}. \quad (2)$$

The SvM-*pdf* is symmetrical and unimodal and therefore can approximate directional data with one prevailing direction only. For $\kappa = 0$, SvM distribution becomes uniform around the circle with all directions equally probable. When the modelled data contains more than one prevailing direction, it is necessary to use a mixture of such distributions.

Finite mixture model of simple von Mises distributions is defined by probability density function (MvM-*pdf*)

$$\phi(\theta; \boldsymbol{\nu}) = \sum_{j=1}^k \omega_j \cdot f_j(\theta; \mu_j, \kappa_j), \quad (3)$$

where k is the number of functions creating the mixture, j is the index of particular SvM-*pdf* with parameters μ_j and κ_j , θ is an angular variable ($0 \leq \theta \leq 2\pi$), and $\boldsymbol{\nu}$ is a vector parameter

$$\boldsymbol{\nu} = (\boldsymbol{\mu}, \boldsymbol{\kappa}, \boldsymbol{\omega}) = (\mu_1, \dots, \mu_k, \kappa_1, \dots, \kappa_k, \omega_1, \dots, \omega_k). \quad (4)$$

Weight of each mixture member has to be nonnegative and satisfy the following conditions

$$0 \leq \omega_j \leq 1 \quad \forall j \in \{1, \dots, k\}, \quad \sum_{j=1}^k \omega_j = 1. \quad (5)$$

The probability that wind direction is within particular angular sector $\langle \theta_l; \theta_u \rangle$ can be obtained by the integration of MvM-*pdf* between the boundary values of θ ,

$$P(\theta \in \langle \theta_l; \theta_u \rangle) = \int_{\theta_l}^{\theta_u} \phi(\theta; \boldsymbol{\nu}) d\theta, \quad (6)$$

and has to be evaluate numerically [25].

Maximum likelihood estimates of MvM distribution vector parameter $\boldsymbol{\nu}$ lead to a system of nonlinear equations that has to be solved numerically as well [26]. The parameter estimates for the MvM distribution can be evaluated using an expectation maximization algorithm fully described in [27]. In this work, a bio-inspired meta-heuristic approach is adopted in place of commonly used numerical methods.

3.1 Analytical Estimation of Distribution Parameters

The distribution parameters μ_j , κ_j , and ω_j , $j \in \{1 \dots k\}$ can be numerically approximated from data. Let us assume that histogram of annual wind data is composed of T classes with frequencies O_i and centers θ_i , $i = 1, \dots, T$, respectively. The frequency classes have to be divided into k same-length sectors corresponding to the number of major wind directions.

Prevailing wind direction μ_j for each sector is estimated by

$$\mu_j = \begin{cases} \arctan\left(\frac{s_j}{c_j}\right), & s_j \geq 0, \quad c_j > 0 \\ \frac{\pi}{2}, & s_j > 0, \quad c_j = 0 \\ \pi + \arctan\left(\frac{s_j}{c_j}\right), & c_j < 0 \\ \pi, & s_j > 0, \quad c_j = -1 \\ 2\pi + \arctan\left(\frac{s_j}{c_j}\right), & s_j < 0, \quad c_j > 0 \\ 3\frac{\pi}{2}, & s_j < 0, \quad c_j = 0 \end{cases}, \quad (7)$$

where s_j and c_j represent the average sine and cosine of selected sector data.

The concentration parameter κ_j is for each sector estimated iteratively by solving the equation

$$\frac{I_1(\kappa_j)}{I_0(\kappa_j)} = \sqrt{s_j^2 + c_j^2}, \quad (8)$$

where $I_1(\kappa_j)$ is the modified Bessel function of the first kind and order one defined by

$$I_1(\kappa_j) = \frac{1}{\sqrt{2\pi}} \int_0^{2\pi} \cos(\theta) \exp^{\kappa_j \cos(\theta)} d\theta = \sum_{k=0}^{\infty} \frac{1}{\Gamma(k+2)\Gamma(k+1)} \left(\frac{\kappa_j}{2}\right)^{2k+1}, \quad (9)$$

and $\Gamma(n) = (n-1)!$ is the gamma function. Alternatively, κ_j can be approximated by

$$|\kappa_j| = \{23.29041409 - 16.8617370 \sqrt[4]{s_j^2 + c_j^2}\}. \quad (10)$$

The latter approach is used in this study.

Initial weight ω_j is approximated as

$$\omega_j = \frac{\sum_{i=J_l}^{J_u} O_i}{\sum_{i=1}^T O_i}, \quad (11)$$

where J_l and J_u are index boundary values of j -th sector.

3.2 Test of Goodness-of-Fit

Pearson's chi-squared goodness-of-fit test, χ^2 , is a statistical test applied to collections of categorical data, such as histograms. Its properties were first described by Karl Pearson in 1900 [28]. It can be used to establish whether or not an observed frequency distribution differs from a theoretical distribution. In this study, χ^2 is directly used as the objective function, f_{obj} , to be minimized by the optimization procedure described in the next section.

The chi-squared test statistic resembles a normalized sum of squared deviations between the observed and theoretical frequencies

$$\chi^2 = \sum_{i=1}^T \frac{(O_i - np_i)^2}{np_i}, \quad (12)$$

where T is the number of frequency classes, n is the sum of all observed frequencies, O_i , and p_i represents theoretical probabilities of each frequency class. The theoretical probabilities are evaluated either directly from the expected cumulative distribution function

$$p_i = F(u_i) - F(l_i) \quad (13)$$

or as a definite integral of the expected *pdf*

$$p_i = \int_{l_i}^{u_i} f(x)dx, \quad (14)$$

where u_i and l_i are boundary values of corresponding frequency class [29].

4 Particle Swarm Optimization

Particle swarm optimization is a global, population-based search and optimization algorithm based on simulation of swarming behaviour of bird flocks, fish schools and even human social groups [30–32]. PSO uses a population of motile candidate particles characterized by their position, x_i , and velocity, v_i , inside an n -dimensional search space they collectively explore. Each particle remembers the best position (in terms of fitness function) it visited, y_i , and is aware of the best position discovered so far by the entire swarm, \bar{y} . In each iteration, the velocity of particle i is updated [31] according to

$$v_i^{t+1} = v_i^t + c_1 r_1^t (y_i - x_i^t) + c_2 r_2^t (\bar{y}^t - x_i^t), \quad (15)$$

where c_1 and c_2 are positive acceleration constants that influence the tradeoff between exploration and exploitation. Vectors r_1 and r_2 contain random values sampled from a uniform distribution on $[0, 1]$. The position of particle i is updated given its velocity [31] as follows

$$x_i^{t+1} = x_i^t + v_i^{t+1}. \quad (16)$$

Create population of M particles with random position and velocity;
 Evaluate an objective function f_{obj} ranking the particles in the population;
while *Termination criteria not satisfied* **do**
 for $i \in \{1, \dots, M\}$ **do**
 Set personal and global best position:
 if $f_{obj}(x_i) < f_{obj}(y_i)$ **then**
 | $y_i = x_i$
 end
 if $f_{obj}(x_i) < f_{obj}(\bar{y})$ **then**
 | $\bar{y} = x_i$
 end
 Update velocity of particle i by (15) and its position by (16);
 end
end

Algorithm 1. Particle swarm optimization PSO *gbest*

A basic global PSO (*gbest*) according to [31, 32] is summarized in Algorithm 1.

PSO is useful for dealing with problems whose solution can be represented as a point or surface in an n -dimensional search space. Candidate solutions (particles) are placed in this space and provided with a random initial velocity. The particles then move through the search space and are periodically evaluated using a fitness function. Over time, particles are accelerated towards those locations in the problem space that have relatively better fitness values.

In addition to the basic model, there is a number of alternative versions of PSO algorithm including self-tuning PSO, niching PSO, and multiple-swarm PSO. These variants have been developed to improve the convergent properties of the algorithm, or to solve other specific problems [30, 31].

4.1 Particle Swarm Optimization for Optimization of Wind Direction Distribution Parameters

The PSO for wind direction distribution parameter optimization is defined by candidate solution representation, fitness function, and population initialization algorithm.

In the proposed algorithm, the mixture of k von Mises distributions is represented by a candidate vector $\mathbf{v} = (v_1, \dots, v_n)$, $v_i \in [0, 1]$ with three parts encoding the vector parameter $\boldsymbol{\nu} = (\boldsymbol{\mu}, \boldsymbol{\kappa}, \boldsymbol{\omega})$ respectively

$$\mathbf{v} = \overbrace{(v_1, \dots, v_k)}^{\boldsymbol{\mu}}, \underbrace{(v_{k+1}, \dots, v_{2k})}_{\boldsymbol{\kappa}}, \overbrace{(v_{2k+1}, \dots, v_n)}^{\boldsymbol{\omega}}. \quad (17)$$

The decoding of $\boldsymbol{\mu}$ involves scaling of v_i , $i \in \{1, \dots, k\}$ to $[0, 2\pi]$ and the decoding of $\boldsymbol{\kappa}$ requires scaling of v_i , $i \in \{k+1, \dots, 2k\}$ to $[0, 700]$. The upper bound of κ_j was in this work chosen with respect to precision of the numerical algorithm used to implement I_0 .

The following decoding rule is applied to satisfy the requirements for mixture member weights ω_j defined in (5).

$$\omega_j = \begin{cases} v_i, & j = 1, \quad i = 2k + 1 \\ v_i \left(1 - \sum_{l=1}^{j-1} w_l \right), & j > 1, \quad i = 2k + 1 + j \end{cases}. \quad (18)$$

It guarantees that the sum of all weights ω_j is equal to 1 and does not impose any additional constraints on candidate vector handling. All vectors created in scope of the optimization process are valid candidate solutions representing a mixture of k von Mises distributions.

In the proposed algorithm, the initial population of particles is formed from randomly perturbed analytical estimates of $\boldsymbol{\mu}$, $\boldsymbol{\kappa}$, and $\boldsymbol{\omega}$ created according to the principles outlined in section 3.1. The fitness of each candidate solution is evaluated using Pearson's χ^2 function (12). The use of a standard statistical goodness-of-fit measure allows direct statistical interpretation of the results.

5 Experiments

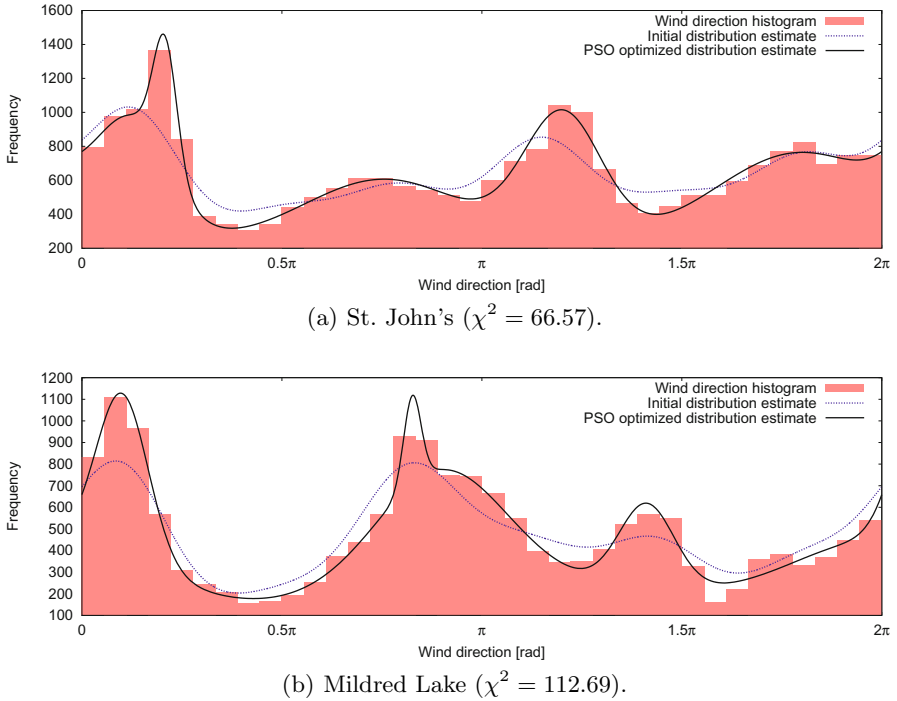
A series of computational experiments with wind direction data from two distinct locations was conducted in order to evaluate the proposed approach. The first data set was collected at St. John's airport located in Newfoundland, Canada. The wind direction records cover the period from 01. 01. 2007 to 31. 12. 2009 with time resolution of one hour and angular resolution 10 degrees. However, 2,920 hourly measurements are missing from this data set so it contains only 23,384 records. These data were used in an earlier study on wind direction modelling by Heckenbergerova et al. [7]. The second data set describes hourly wind direction measured at Mildred Lake, Alberta, Canada and covers the period between 1. 1. 2012 and 31. 12. 2013. This collection contains more than 17,200 records with angular resolution of 0.1 degree. 330 records are missing from this collection due to sensor reading errors.

PSO parameters used for the experiments were: population size 100, 50000 iterations, inertia weight $w = 0.89$, local PSO weight $c_1 = 0.50$, and global PSO weight $c_2 = 0.70$. The parameters were selected on the basis of best practices, previous experience with the PSO algorithm, and initial trial-and-error runs. Each experiment was repeated 30 times due to the stochastic nature of the PSO algorithm. Results of the experiments are shown in table 1 and visual illustrations of the best optimized wind direction distribution estimates found for St. John's and Mildred Lake are shown in fig. 1(a) and fig. 1(b) respectively. Indeed, the optimized distribution estimates fit the data more accurately than the initial analytical estimates.

The results indicate that the proposed PSO can be used for optimization of wind direction distribution estimates and that the distribution with optimized parameters approximates real wind directions better than the original analytical estimate. Moreover, the goodness-of-fit χ^2 for the data from St. John's is in all

Table 1. Wind direction distribution estimates optimized by the PSO

Location	Goodness-of-fit (χ^2)		
	best	average	worst
St. John's	66.57	133.76	161.58
Mildred Lake	112.69	141.64	194.90


Fig. 1. Best wind direction distribution estimates found by the PSO

cases (best, average, worst) better than that of the estimates obtained by GA in [7]. It means that the worst result obtained by the PSO proposed in this work is better than the best estimate (with $\chi^2 = 181.64$) obtained by the GA in [7]. Nevertheless, $\chi^2 = 66.57$ for St. John's and $\chi^2 = 112.69$ for Mildred Lake means that, from the statistical point of view, the PSO optimized MvM-*pdf* estimate still does not pass the goodness-of-fit test with the level of significance $\alpha = 0.05$.

6 Conclusions

This study introduced a new bio-inspired meta-heuristic method for circular data modelling based on the popular PSO algorithm. PSO is suitable for the

search for MvM-*pdf* due to its ability to solve floating-point valued problems and good results solving real-world optimization. A new representation of MvM-*pdf* vector parameter, ν , suitable for the use with PSO was defined and the proposed algorithm was evaluated on directional wind data collected on two distinct locations in Canada. The results show that PSO is able to improve initial analytical ν estimates and can be used for wind direction modelling. Moreover, the model optimized by PSO exhibits on the same data set better accuracy and lower error than a previous model based on GA. Although more accurate than previous ones, the optimized models still do not pass the χ^2 test. Thus, from a strict statistical perspective, these models do not match the data.

There is a plenty of opportunities for further work in this area. Different fitness functions (e.g. R^2 , RMSE) as well as other real-parameter optimization methods can be used for meta-heuristic wind direction distribution estimation. Another research direction is the design of an algorithm that would automatically select the number of distributions needed in the mixture to represent given data well. From the application point of view, the proposed approach can be used to model different types of circular data, in addition to wind direction data used as an example in this contribution.

Acknowledgement. This work has been supported by the Natural Sciences and Engineering Council of Canada (NSERC), Helmholtz-Alberta Initiative (HAI), the Ministry of Education, Youth and Sports of Czech Republic project CZ.1.07/2.3.00/30.0058, and by institutional support of the University of Pardubice. It was also partially supported by the European Regional Development Fund in the IT4Innovations Centre of Excellence project (CZ.1.05/1.1.00/02.0070) and by the Bio-Inspired Methods: research, development and knowledge transfer project, reg. no. CZ.1.07/ 2.3.00/20.0073 funded by Operational Programme Education for Competitiveness, co-financed by ESF and state budget of the Czech Republic. VŠB - Technical University of Ostrava supported this work via SGS under the project no. SP2014/110.

References

1. Mooney, J.A., Helms, P.J., Jolliffe, I.T.: Fitting mixtures of von Mises distributions: a case study involving sudden infant death syndrome. *Computational Statistics & Data Analysis* 41(3-4), 505–513 (2003); *Recent Developments in Mixture Model*
2. Calderara, S., Cucchiara, R., Prati, A.: Detection of Abnormal Behaviors Using a Mixture of Von Mises Distributions. In: *Proceedings of the 2007 IEEE Conference on Advanced Video and Signal Based Surveillance, AVSS 2007*, pp. 141–146. IEEE Computer Society, Washington, DC (2007)
3. Vejmelka, M., Musilek, P., Paluš, M., Pelikán, E.: K-means clustering for problems with periodic attributes. *International Journal of Pattern Recognition and Artificial Intelligence* 23(4), 721–743 (2009)
4. Allen, C.T., Young, G.S., Haupt, S.E.: Improving pollutant source characterization by better estimating wind direction with a genetic algorithm. *Atmospheric Environment* 41(11), 2283–2289 (2007)

5. Carta, J.A., Bueno, C., Ramírez, P.: Statistical modelling of directional wind speeds using mixtures of von mises distributions: Case study. *Energy Conversion and Management* 49(5), 897–907 (2008)
6. Carta, J.A., Ramírez, P., Bueno, C.: A joint probability density function of wind speed and direction for wind energy analysis. *Energy Conversion and Management* 49(6), 1309–1320 (2008)
7. Heckenbergerova, J., Musílek, P., Mejznar, J., Vancura, M.: Estimation of wind direction distribution with genetic algorithms. In: CCECE, pp. 1–4. IEEE (2013)
8. Masseran, N., Razali, A., Ibrahim, K., Latif, M.: Fitting a mixture of von Mises distributions in order to model data on wind direction in Peninsular Malaysia. *Energy Conversion and Management* 72, 94–102 (2013); The III. International Conference on Nuclear and Renewable Energy Resources {NURER2012}.
9. Jung, S., Kwon, S.D.: Weighted error functions in artificial neural networks for improved wind energy potential estimation. *Applied Energy* 111, 778–790 (2013)
10. Bazilian, M., Nussbaumer, P., Rogner, H.H., Brew-Hammond, A., Foster, V., Kammen, D.M., Pachauri, S., Williams, E., Howells, M., Niyongabo, P., Lawrence, M.: O Gallachoir, B., Radka, M.: *Energy Access Scenarios to 2030 for the Power Sector in Sub-Saharan Africa*. Utilities Policy 20, 1–16 (2012)
11. Pereira, M.G., Camacho, C.F., Freitas, M.A.V., da Silva, N.F.: The renewable energy market in Brazil: Current status and potential. *Renewable and Sustainable Energy Reviews* 16(6), 3786–3802 (2012)
12. Bekele, G., Tadesse, G.: Feasibility study of small Hydro/PV/Wind hybrid system for off-grid rural electrification in Ethiopia. *Applied Energy* 97, 5–15 (2012); *Energy Solutions for a Sustainable World - Proceedings of the Third International Conference on Applied Energy, Perugia, Italy, May 16-18 (2011)*
13. Davis, M.W.: A new thermal rating approach: The real time thermal rating system for strategic overhead conductor transmission lines – Part I: General description and justification of the real time thermal rating system. *IEEE Transactions on Power Apparatus and Systems* 96(3), 803–809 (1977)
14. Douglass, D.: Weather-dependent versus static thermal line ratings [power overhead lines]. *IEEE Transactions on Power Delivery* 3(2), 742–753 (1988)
15. Heckenbergerová, J., Musílek, P., Filimonenkov, K.: Quantification of gains and risks of static thermal rating based on typical meteorological year. *International Journal of Electrical Power & Energy Systems* 44(1), 227–235 (2013)
16. Ettoumi, F., Sauvageot, H., Adane, A.E.H.: Statistical bivariate modelling of wind using first-order Markov chain and Weibull distribution. *Renewable Energy* 28(11), 1787–1802 (2003)
17. García-Rojo, R.: Algorithm for the Estimation of the Long-term Wind Climate at a Meteorological Mast Using a Joint Probabilistic Approach. *Wind Engineering* 28(2), 213–223 (2004)
18. Erdem, E., Shi, J.: {ARMA} based approaches for forecasting the tuple of wind speed and direction. *Applied Energy* 88(4), 1405–1414 (2011)
19. de la Rosa, J.J.G., Pérez, A.A., Salas, J.C.P., Leo, J.G.R., Muñoz, A.M.: A novel inference method for local wind conditions using genetic fuzzy systems. *Renewable Energy* 36(6), 1747–1753 (2011)
20. Shamshirband, S., Iqbal, J., Petković, D., Mirhashemi, M.A.: Survey of four models of probability density functions of wind speed and directions by adaptive neuro-fuzzy methodology. *Advances in Engineering Software* 76, 148–153 (2014)
21. Tascikaraoglu, A., Uzunoglu, M.: A review of combined approaches for prediction of short-term wind speed and power. *Renewable and Sustainable Energy Reviews* 34, 243–254 (2014)

22. Koch, W.: Directional analysis of SAR images aiming at wind direction. *IEEE Transactions on Geoscience and Remote Sensing* 42(4), 702–710 (2004)
23. AlRashidi, M.R., El-Hawary, M.: A Survey of Particle Swarm Optimization Applications in Electric Power Systems. *IEEE Transactions on Evolutionary Computation* 13(4), 913–918 (2009)
24. Onwunali, J., Durlofsky, L.: Application of a particle swarm optimization algorithm for determining optimum well location and type. *Computational Geosciences* 14(1), 183–198 (2010)
25. Mardia, K., Jupp, P.: *Directional Statistics*. John Wiley & Sons (2000)
26. Fisher, N.I.: *Statistical analysis of circular data*. Cambridge University Press (1995)
27. Banerjee, A., Dhillon, I.S., Ghosh, J., Sra, S.: Clustering on the unit hypersphere using von Mises-Fisher distributions. *Journal of Machine Learning Research* 6 (2005)
28. Pearson, K.: On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. *Philosophical Magazine Series* 5 50, 157–175 (1900)
29. Anderson, T.W., Darling, D.A.: A Test of Goodness of Fit. *Journal of the American Statistical Association* (49), 765–769 (1954)
30. Clerc, M.: *Particle Swarm Optimization*. ISTE. Wiley (2010)
31. Engelbrecht, A.: *Computational Intelligence: An Introduction*, 2nd edn. Wiley, New York (2007)
32. Kennedy, J., Eberhart, R.: Particle Swarm Optimization. In: *Proceedings of the IEEE International Conf. on Neural Networks*, vol. 4, pp. 1942–1948 (1995)

A Perspective of the Cellular Network of the Future: Cloud-RAN

Santhi Kumaran

School of ICT,
College of Science and Technology, University of Rwanda,
Kigali, Rwanda
santhikr@yahoo.com

Abstract. The increasingly widespread use of smartphones capable of delivering high bandwidth content is leading to an exponential growth in the volume of traffic carried by mobile networks. The solution to this can be found in a new distributed architecture called Cloud Radio Access Network (C-RAN) which offers a new paradigm in base station architecture. It is less costly for mobile operators to deploy because it could run on smaller cell sites. In a C-RAN setup, the baseband processing elements of a traditional base station are centrally located and connected to smaller, distributed radios, often via fiber. The concept of C-RAN lowers cost and power dramatically and has recently generated significant interests among the research community. Though some mobile operators across the globe are evaluating, or even already migrating to, C-RAN architectures, this paper will explain the need for C-RAN architecture, the drivers behind it and about the recent advances in mobile fronthaul solutions. This paper also contains comprehensive details of the various C-RAN deployment methodologies based on the mobile traffic needs and describes about the perception of C-RAN technology by different mobile operators.

Keywords: C-RAN, BU pool, RRU, centralized baseband processing, mobile fronthaul.

1 Introduction

As smartphone users are more habitual to mobile internet with smart applications handling streaming video and audio, there is an explosion of the mobile traffic. The conventional mobile network cannot meet the demand of large amounts of mobile data traffic unless until very dense and low-power small-cell networks with a very high spatial reuse are being deployed. But these small cells face varied problems like high resource requirements, which in turn increases Operation and Maintenance (OAM) costs, makes network planning difficult with more cells & less intervals, very high inter-cell interference leading to inefficient spectrum usage, inefficient resource sharing due to tidal effect etc.

The above problems could be addressed by an innovative and environment friendly radio access network architecture called Cloud-Radio Access Network (C-RAN).

Traditional base stations are called Radio Access Networks (RAN), but C-RAN will move conventional digital base station electronics away from the cell tower, 20 to 50 km away, back in the data center. In this paper, section 2 describes the C-RAN concept, sections 3, 4 explains the need and the characteristics of C-RAN, section 5 gives the various perceptions of C-RAN and section 6 described the deployment use cases. Section 7, 8 provides the advantages and suggestion for improvement of C-RAN respectively. Section 9 concludes the paper.

2 Cloud Radio Access Network (C-RAN)

Conventional base stations are called Radio Access Networks (RAN) and each base station should be placed in the geographical center of its coverage area and it serves all the mobile devices within its reach. Each base station has its digital component to manage its radio resources, handoff, data encryption and decryption and an RF component which transforms the digital information into analog RF. The RF elements are connected to a passive antenna that transmits the signals to the air [1]. Distributed Base Station Architecture is the fundamental building block of C-RAN Architecture [2]. In open platform C-RAN setup, the base station is split in to a Baseband Unit (BU) and a tiny Remote Radio Unit (RRU) and Antenna as shown in figure 1 below.

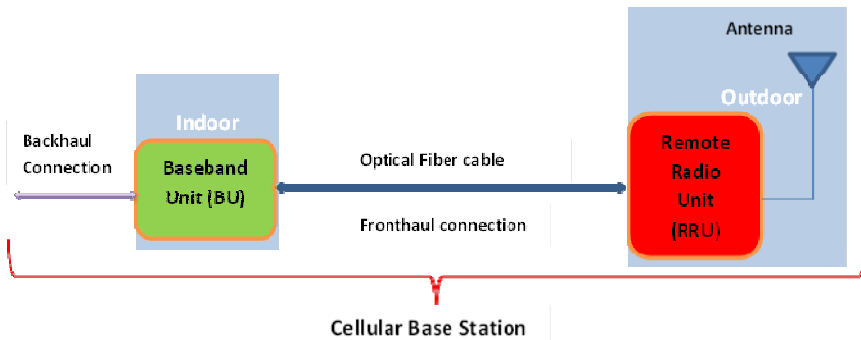


Fig. 1. Cloud RAN Concept

BU implements the MAC PHY and Antenna Array System (AAS) functionality. The BUs of various terrestrial cell sites are centrally located to form clusters of Virtual Base Station Pool also called as Baseband processing Cloud where the baseband processing takes place. The centralized processing servers can perform the baseband processing for a high number of cell sites for example 100 to 1000 cell sites. As the BUs are centralized, RRUs with antennas only remain at the cell sites for radio transmission. RRUs are connected to the centralized Base Station Pool via fiber as shown in figure. RRU is an active unit that obtains the digital (optical signal) signals through fiber, converts digital signals to analog, amplifies the power, and sends the actual transmission. The idea is that such a distributed network could be less costly to deploy for mobile operators because they can now place numerous BUs in a single

geographical point while distributing and densely deploying the RRUs. RRUs with baseband processing virtualized at the cloud is called as the C-RAN concept, which is intended to increase the capacity and provides ubiquitous coverage for cellular network of the future.

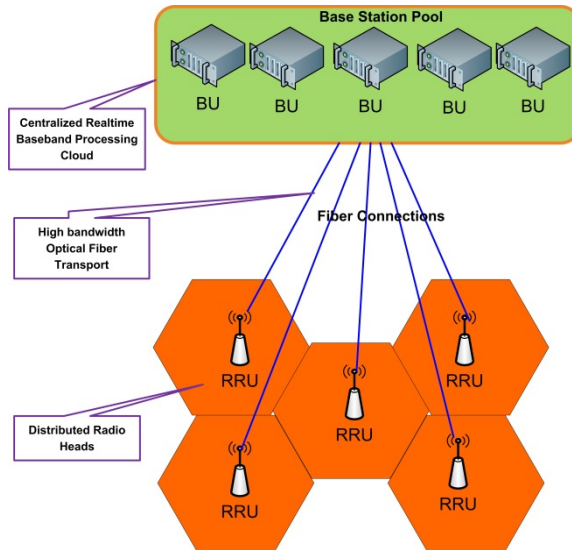


Fig. 2. Cloud RAN Architecture

3 Need for C-RAN Architecture

There is a need for a single mobile network architecture such as C-RAN that consolidates the workloads in to a more scalable and simplified solution in a cost effective manner because today's conventional mobile networks tends to limit flexibility and increases cost. Moreover conventional mobile operators struggle to attain a stable growth in this mobile internet era due to the challenges listed below:

(i) Increasing mobile data Traffic: Smartphone and Tablet PC users require ubiquitous coverage for mobile internet usage. According to [3] the mobile broadband traffic is expected to increase from 57 exabytes in 2013 to almost 335 exabytes in 2020 and the conventional macro cell based cellular network design restricts the development of high-speed mobile internet demand.

(ii) Surging Power consumption: As stated by Huang in [4], in the conventional mobile architecture, outdoor base stations consume 67% of the network's energy, which is expected to increase, thus leading to low power efficiency.

(iii) High Capital Expenditure (CAPEX) and Operational Expenditure (OPEX) of Traditional RAN: In a conventional cellular network the increasing demands on

mobile traffic, can be handled by increasing the number of cells by cell splitting. But, microcell deployment in turn increases the cost of infrastructure equipment to 30% of Total Cost of Ownership (TCO) [5] and OAM costs.

(iv) Space limitation: In traditional RAN each base station should be placed in the geographical center of its coverage area. Operators find difficult to get the space to install the base station enclosures in dense urban areas [6]. Even if site locations are available, the mobile operators have difficulties in renting the real estate, finding proper powering options, securing the location and protecting the equipment from weather conditions etc. which incur high CAPEX and OPEX.

(v) Low Resource utilization rate: Generally mobile data traffic is bursty and fluctuating in nature because mobile users are moving from one place to another. Therefore, each of the base stations are designed to handle the maximum traffic expected by the network. So there is a high likelihood that many of the base stations carry no or little traffic, and, these processing resources cannot be shared with other BTS leading to high energy waste [7].

(vi) Inter cell interference problem: The main challenge of a traditional RAN design is to deal with radio frequency (RF) interference, which limits network capacity. In order to improve system capacity many small cell base stations have to be installed and hence base stations using the same frequency or channel will be close to each other causing interference which leads to low bandwidth.

(vii) Support of heterogeneous networks: Traditional RAN lacks multisystem convergence, which requires heterogeneous backhaul solutions to connect the small-cells to the core network which leads to high infrastructure cost.

In order to overcome the above challenges, comes in the innovative C-RAN architecture for small cells with RRUs connected through optical fiber to a central base station housing the baseband processing [2].

4 Characteristics of C-RAN Technology

Here we provide the functionalities and characteristics of C-RAN technology and the tools for its deployment.

4.1 Multi-antenna Array& Remote Radio Unit (RRU)

Antennas with built in RRUs are already available in the market. RRUs transmit the digitized waveforms to the base station cloud through an Ir interface or Common Public Radio Interface (CPRI) [8] or Open Base Station Architecture Initiative (OBSAI) [9]. AIRS (Antenna Integrated Radio Unit) from Ericsson, Alcatel-Lucent's LightRadio [6], and Nokia Siemens' Liquid Radio [10] are examples of RRUs.

4.2 Common Base Station (BS) Pool/ Cloud

Servers with more powerful communications processors that handle the more sophisticated processing functions are replacing the base stations and form a Virtual BS pool. For example, multicore general purpose processors (GPP) like Intel's bank of x86 servers running software-defined radio (SDR) application, and IBM's POWER processors are well matched to the computation - intensive portions such as the physical (PHY) layer. Processors like IBM PowerEN and Sun Niagara, which are heavily multithreaded, as well as low power ARM servers are well suited for network processing (MAC and transport layers). These servers require software-defined network implementations to give real-time responses in a few milliseconds [4].

4.3 Mobile Fronthaul Connection

Base Station Cloud servers, however, would require connectivity to thousands of tiny RRUs. This introduces a new transmission network into the C-RAN infrastructure - Mobile Fronthaul. It is a high bandwidth optical fiber where signals are transmitted at a typical bandwidth of several gigabits per second to sites at a distance as long as 40 km. C-RAN concept is expected to deliver high-data rates of 4G networks of about 300 Mbps and more. Therefore at least 1 Gbps fronthaul connection is required per base station. In many cases two 1Gbps connection is installed for redundancy purposes. If higher bandwidth is required these two 1Gbps connection can be unified using LAG [1]. If data compression adopted, would reduce bandwidth and transport costs [6].

4.4 Backhaul Connection

This is a connection from the Baseband pool toward the core mobile provider infrastructure. C-RAN architecture requires an enormous backhaul capacity that the mobile provider requires multiple 10Gig backhaul in order to ensure a non-blocking architecture. The newer generation mobile technologies (e.g. HSPA+, 4G, WiMAX etc.) uses a high capacity Carrier Ethernet- based backhaul with Synchronous Ethernet and 1588v2 Timing over Packet replacing previous synchronization schemes [1].

4.5 Co-operative Radio

In C-RAN Carrier Aggregation (CA) will be supported for cells served by the same base station enabling flexible deployment of add on cells. It is expected C-RAN baseband pool share more than 1000 carriers. A single RRU supports 12 or more carriers and cross cell configuration. [5]. When combined with Multi Input Multi Output (MIMO), the CA techniques may lead up to 100 megabits per second between the radio and the baseband units (i.e., the front-haul) [12]. LTE-A offers aggregation of up to five 20-megahertz LTE carriers.

4.6 Real Time Processing Cloud

In order to support real time processing, massive and concurrent baseband calculations done by general purpose CPUs were offloaded to highly specialized Modem Processing Units (MPU) [14]. MPU solution which is re-configurable at runtime has a great advantage over traditional hard-wired designs. Intel's software enhanced generic Xeon processors, IBM's PowerPC or ARM cores [13] could be used. MPUs support a wide range of system partitioning, and advanced algorithms for interference management and coverage, including Collaborative Multi-point communications (CoMP), Enhanced Inter-Cell Interference Coordination (eICIC) and massive MIMO [14]. CoMP can improve throughput by as much as 80% at the cell edge, for both uplink and downlink and requires very low latency (less than 1-2 microseconds) to provide a big boost in capacity. eICIC coordinates blank and almost-blank sub-frames on a real-time basis, so locating the baseband processing in one place greatly reduces the complexity of coordinating changes [15]. Moreover, Virtual BS pool framework provides the interface between SW components to control the overall jitter and latency, and improve the precision of the synchronization between the RRH and Virtual Base Station (VBS). The SW component works in a data-driven mode and the latency caused by data transmission between the data interfaces is minimized [9].

4.7 Synchronization Schemes

Multiple synchronization schemes are required to support frequency and phase synchronization between the VBS and the RRU. Most common synchronization schemes are:

- IEEE 1588v2 – 1588v2 is a packet-based protocol and is carried in-band with user traffic. It can support highly accurate frequency and phase synchronization but is affected by network congestion if not properly prioritized across the path [1].
- SyncE (AKA G.8261) is a physical layer technology that supports frequency synchronization and is not affected by network congestion. However it requires that every node in the path have hardware support for SyncE [1].
- GPS used by many mobile providers, meets both phase and frequency time synchronization requirements but involves high CAPEX and low OPEX [1].
- Transmode's active and passive WDM mobile fronthaul architectures provide good synchronization and transparent transport with its low power and low latency capabilities.

4.8 Strong QoS (Quality of Service)

C-RAN architecture supports multiple services for customers that have differing quality of service priorities. Network failure is not accepted. Therefore, C-RAN uses HQOS (Hierarchical QoS) that allows for higher quality of service granularity, enabling the backhaul provider to ensure the availability and quality of service parameters (latency, delay, amount of bandwidth) per user, per service [1].

4.9 Strong OAM (Operation and Maintenance)

C-RAN architecture uses various tool sets specified by IETF RFC 2544, Y.1731, IEEE 802.1ag CTM (connectivity Fault Management), and Service OAM etc. for effective proactive troubleshooting. Moreover, the idle field of CPRI protocol between BU and RRU connections is used to develop an array of functions for centralized OAM. This makes network controllable and manageable [5].

5 C-RAN Perception by Leading Telecom Operators

Different operators perceive the cloud-based mobile networks in a different manner but the core idea is to move the highly-specialized wireless architectures into more general purpose computing platforms ("the cloud") which takes care of performing the baseband processing.

5.1 China Mobile's C-RAN

China Mobile calls its cloud-based mobile network architecture as Cloud-RAN (C-RAN). It splits the Base Station into BaseBand Unit (BBU) and into Remote Radio Head (RRH). BBUs are accumulated into a BBU Pool which does the baseband processing and connected to RRHs via fiber which converts the Radio signal into an RF signal.

5.2 DOCOMO's Advanced Cloud RAN

NTT Docomo refers its cloud based mobile network as Advanced Centralized RAN. Advanced C-RAN will enable small 'add-on' cells for localized coverage to cooperate with macro cells that provide wider area coverage and increase throughput and overall system capacity as shown in Figure 3. This will be achieved with carrier aggregation technology, and will rely on some HetNet principles. The cell site equipment, consisting of radio and antenna, is compact and low power as that of C-RAN. High-capacity base stations utilizing advanced C-RAN architecture will serve as master base stations both for multiple macro cells covering broad areas and for add-on cells in smaller, high-traffic areas [16].

5.3 South Korean's CCC

In South Korea, both KT and SK calls in cloud-based mobile network similar to C-RAN architecture as the Cloud Communication Center (CCC) and Smart Cloud Access Network (SCAN) respectively [7]. With Intel's servers and data centers integrated with Samsung modems create a centralized exchange for baseband communications processing the CCC [10]. CCC is linked by fiber to the cell sites which has only the RRH and antennas. Alcatel Lucent's LightRadio Cube [6] and Nokia Siemens' Liquid Radio [13] are used at the cell sites. Similar to C-RAN, CCC harnesses virtualization technology so that the central processing resources can be allocated flexibly according to the peaks and troughs of demand in different sites. The CCC architecture can manage 144 base stations per server and accommodate 1,000 servers in each data center, all of them acting as a central processing entity [10].

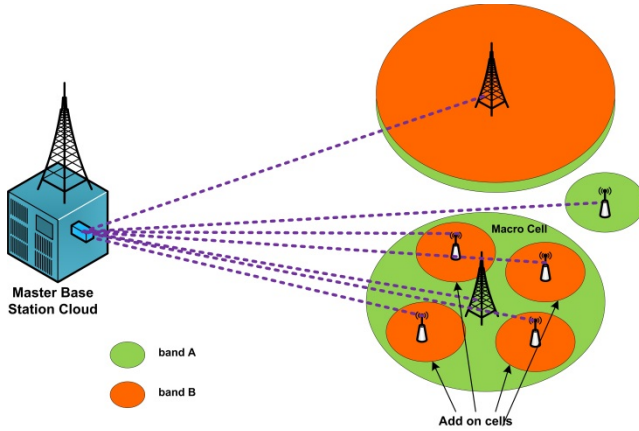


Fig. 3. Advanced Cloud RAN Architecture

6 C-RAN Deployment Methodologies

6.1 Based on Network Traffic

This section describes three different deployment uses cases based on the network traffic requirements and the suitability of adopting C-RAN architecture for them.

(i) Deployment of combined Macro and Micro cellular Network (Suburban to urban flow): In a business district mobile data traffic will be very high in the mornings and evenings. Outside of those peak times, that capacity just goes to waste. But with Cloud-RAN, operators can allocate capacity where and when needed, following the flow of network congestion from the suburbs to the central city and back again [13]. Advanced Centralized RAN concept could be best suitable for this scenario.

(ii) Deployment of dense Micro Cells (Heavy Urban Load): C-RAN will make dense urban networks possible. The C-RAN could solve lot of problems in urban geographic where the demand could be switching dynamically across different areas. It can be deployed in different ways. By splitting an omni cell in to several small sector cells or by splitting a big macro cells into several small cells by adding more light RRUs and connecting them to BU cloud/legacy cell site or by overlay techniques.

(iii) Deployment of Pico Cells Indoor (Enterprise Building): A large financial center in a city can have tens of thousands of mobile devices eager to exchange gigabytes of data with corporate servers and the cloud. Similarly, a large suburban mall can contain tens of thousands of customers who want to access social-media websites. But today's macro-cellular networks are not capable of providing reliable coverage and capacity inside buildings. With C-RAN Architecture a number of RRUs with antenna array units can be deployed on nearby ceilings and can be connected to BU or BU pool by fiber or wireless links. Cooperative radio algorithms can be implemented to reduce interference between cells and improve indoor capacity [18].

(iv) Deployment of Super-Hotspots: C-RAN is suitable for special situations such as a stadium with very high density. The same scenario is seen in major transportation hubs like metro/bus/subway stations. C-RAN provides uniform coverage in this kind of environment with interference management algorithms such as eICIC, CoMP etc.[18]. If fiber bandwidth is free, then the cost of transport is not big problem, and C-RAN is attractive [6].

6.2 Based on Fronthaul Fiber Availability

The availability of low cost fiber is a gating factor in deploying cloud RAN architecture. C-RAN deployment methodology is typically influenced by the availability of fibers.

(i) Area with Abundant Fiber Resource: Fiber connection is directly used for fast C-RAN deployment in areas with abundant fibers resources. The benefits of the solution are: fast deployment and low cost because no additional optical transport network equipment is needed [19]. For multi-RAT (Radio Access Technology) star topology is used as shown in Figure 4 and for single-RAT ring topology can be used.

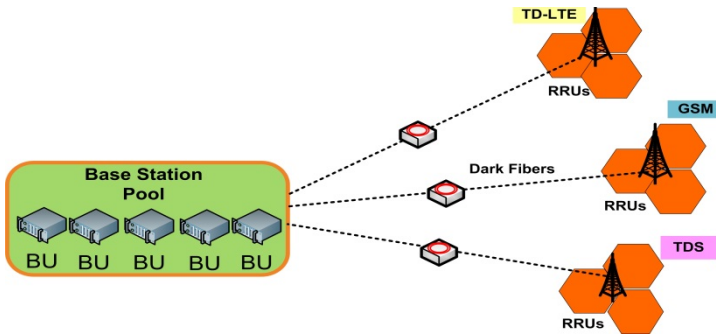


Fig. 4. Star Topology for Fronthaul with multi-RAT

(ii) Area with Available but limited Fiber resource: For most operators around the world, fiber cost is a significant chunk of their operating budget [15]. When the fiber resource available is very limited or adding new fiber cost is high, should go for alternate methodologies.

For this scenario, ring topology can be used for single-RAT and Wavelength Division Multiplexing (WDM)/OTN (Optical Transport Network) ring as shown in Figure 5 for multi-RAT sites [18]. WDM/OTN option, improves the bandwidth of CPRI/Ir/OBRI interface between BU pool and RRU considerably. OTN is a WDM-based integrated bearer device. CPRI can be a service type carried by OTN. Because OTN is introduced into the transmission layer, it provides perfect protection, OAM, and fault diagnosis and can support ring, tree, and mesh topologies [20] and considerably reduces costs.

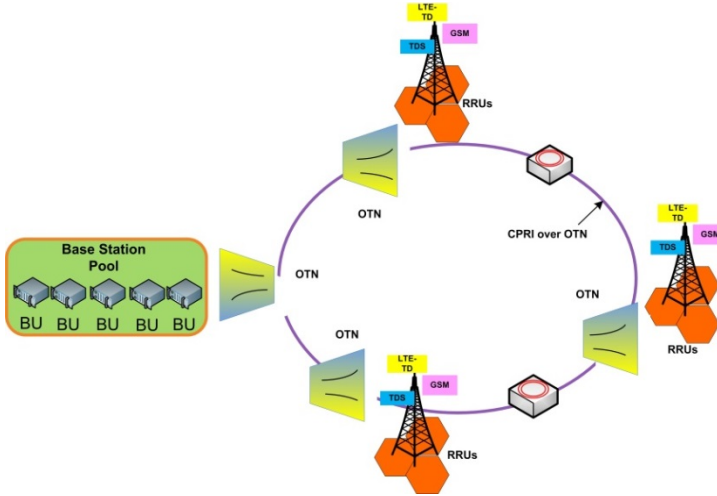


Fig. 5. WDN/OTN Ring Fronthaul Solution

(iii) Area with limited Rental fiber resource: In this case the option is based on Coarse Wavelength Division Multiplexing (CWDM) technology. It combines the fixed broadband and mobile access network transmission at the same time for indoor coverage with passive optical technology, thus named as Unified PON (UniPON). It can provide both PON services and CPRI/Ir/OBRI transmission on the same fiber and supports up to 14 wavelengths. It is suitable for C-RAN centralized baseband pool deployment of indoor coverage [19].

Transmode offers various WDM topologies to cater for the limited fiber resource areas. There are three WDM options proposed by Transmode [21]. i) Active WDM-Transparent, ii) Active WDM-Forward Error Correction (FEC) and enhanced management, and iii) Passive WDM. Figure 6 shows the Active and passive WDM scenarios.

- (a) The Transparent WDM option, uses 1G/2.5G and 10G Transponders (TP), brings important values such as network latency and sync performance, both factors that are critically important in fronthaul networks by adding a small portion of power consumption
- (b) The WDM with Forward Error Correction (FEC) and enhanced management option is based on Transmode’s Transponders and Muxponders (MXP), and adds capabilities such as in-band management, sophisticated performance monitoring and extended reach and address more demanding networking applications [21].
- (c) The Passive Mobile Fronthaul option is based on Transmode’s widely deployed TG-Series platform, which has an extremely broad range of networking options including CWDM or DWDM, single fiber or fiber pair, point-to-point or ring architectures. Protected CPRI transport can also be offered. It supports optical paths up to 80 km.

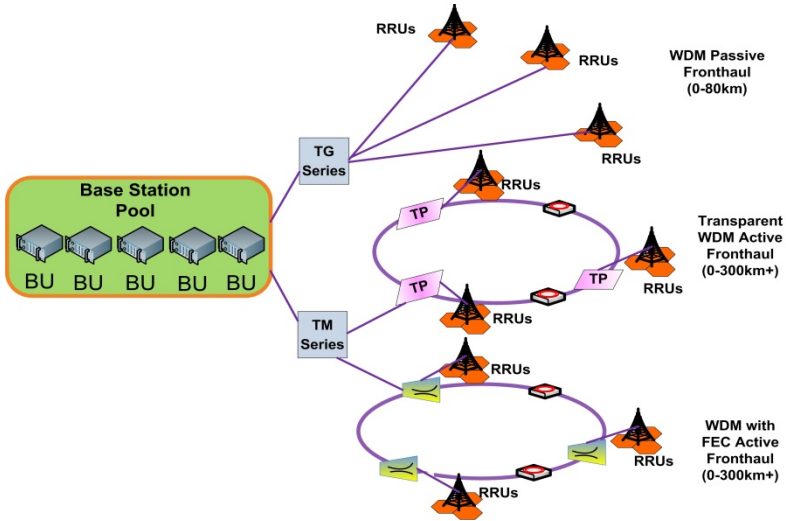


Fig. 6. Transmode's Passive and Active Fronthaul Networks

Moreover, WDM topology has various schemes to use few fibers and to build an independent C-RAN. It is by means of enhanced fiber connection or colored fiber connection. All these solutions use a very few fibers to meet C-RAN bearer requirements in different scenarios.

(a) Colored Fiber Connection

The colored connection is similar to Transmode's Active Mobile fronthaul solution with FEC and enhanced Management. In colored fiber connection [20], WDM is used to configure optical multiplexers/demultiplexers (OMDs) in the BU pool and optical add/drop (OAD) devices on the RRU nodes. OAD is a passive optical device that can be placed in an outdoor power cabinet or in an optical cross-connecting box. OMD and OA can be installed in the equipment room where a BBU pool is located. Figure 7 shows a configuration of four RRU nodes, each with an OAD. Optimal modules on the BBUs and RRUs are colored WDM modules with certain wavelengths. Coarse WDM is used to support a maximum of 18 wavelengths, and Dense WDM is used to support a maximum of 80 wavelengths.

Through colored fiber connection, both TD-SCDMA and TD-LTE C-RAN bearer schemes can be implemented with only one pair of fibers.

(b) Enhanced Fiber

Enhanced fiber connection [20] allows for cascading of up to 18 RRUs through the CPRI interface. A mature 6G optical module can be used for each RRU. Equipment at the radio side provides protection and OAM for enhanced fiber connection using CPRI signals. With enhanced fiber connection, six GSM sites can be networked with only one pair of fibers. This is an optimal transmission scheme for GSM C-RAN.

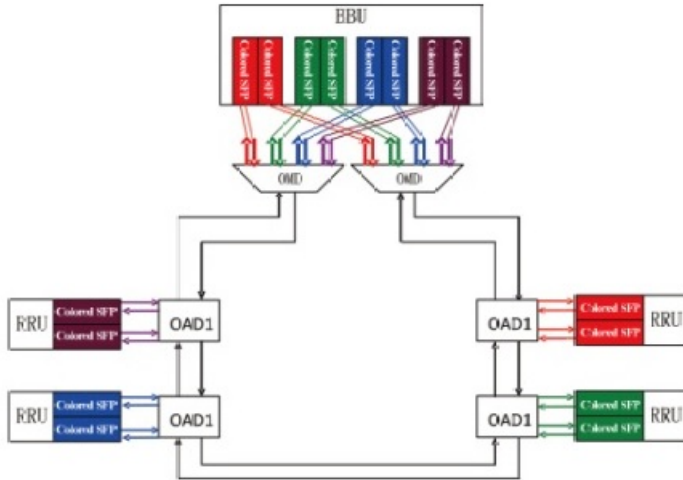


Fig. 7. Colored WDM Connection (source [20])

7 Advantages of C-RAN

(i) **Energy savings:** Operators could drastically cut the processing power necessary (as much as 40 percent) [22,13] to run the network as a whole by putting their base stations in the cloud. Additionally, as the data from all base stations is centralized in the same location, it provides the ability to globally control the access network, when necessary, e.g. to selectively turn off those base stations that are not needed, which results in energy savings [2].

(ii) **Efficient Spectrum Usage:** C-RAN provides greater potential data capacity than a traditional architecture because the radio units could be more densely deployed and spectrum reuse is enhanced.

(iii) **Efficient resource Sharing:** The capacity of a mobile network varies from time to time and from place to place and causes a tidal effect [5]. With C-RAN, operators can allocate capacity where and when needed, following the flow of network congestion [22] and no longer have to build networks to meet peak demands at every tower.

(iv) **Real Estate Cost:** Traditional cell sites carry with them a continuous stream of CAPEX, OPEX to address the high rental rates for real estate, electrical expenses, cost of backhaul for the cell site and security measures to protect the location from intruders and equipment from weather conditions etc. [23]. C-RAN solves these problems.

(v) **Multi System Convergence:** C-RAN is expected to support multiple technologies, including GSM, TD-SCDMA, LTE TDD & LTE TDD-Advanced [12]. Baseband processing units of GSM, TD-SCDMA, and LTE networks based on the same platform architecture can be placed on the same BTS rack. They can share the same backboard, power supply, and main control, clock, and transmission facilities [5].

(vi) Clean Systems: C-RAN consumes less power (71% compared to traditional RAN) and has low OPEX and supports fast system roll out. Due to simpler RRUs, system roll out can be up to 1/3 the usual time [19].

(vii) Easy Migration: As wireless access technologies evolve in a fast pace from 3G to 4G to 5G, operators cannot migrate from one standard to another overnight. Switching over to new networks in a parallel manner is not cost efficient. C-RAN can ease this migration with a virtualization layer and by coupling efficient accelerators that allocates resources to standards as necessary [24].

(viii) Reduced CAPEX & OPEX: The centralized location of the baseband processing, allows maintenance simpler and more efficient by a technician [15] thus saving OPEX. The base station is the most expensive element of the wireless network. C-RAN has shifted the hardware infrastructure at the bottom of every tower to a Centralized Network Computing unit leaving only the radios and antennas at the cell site and thus has reduced CAPEX. The CAPEX estimated 15% and OPEX estimated 50% savings could be translated to savings for the end-users [25]. If an operator builds six sites in rural areas using C-RAN architecture, they can save equipment costs by around 9%, construction costs by around 30%, and OAM costs by around 76% [5].

8 Suggestions to Improve

(i) Standardization initiatives: By adopting the more open standards and common platform the wireless industry could innovate at the faster pace, which would ultimately save development and deployment costs.

(ii) Power Efficient Systems: The equipment suppliers must provide cost-effective power efficient solution so that the operators can gracefully migrate to higher capacity systems.

(iii) Data compression: The data between baseband and a remote radio is typically oversampled I/Q streams for different MIMO antennas. The higher the numbers of antennas in future base station, the higher the resulting data rate and bandwidth. One possible way to achieve these savings is by developing various frequency and time domain data compression schemes that will compress to I/Q streams.

(iv) Base Station virtualization: Virtual BS concept should be developed. It will allow efficient resource utilization of the centralized baseband pool so that operators can dynamically allocate resources to various virtual base stations and air interfaces balancing the tidal effect.

(v) Interference cancellation and load management: More sophisticated interference cancellation and load management techniques should be developed for good scalability and increasing capacity of the C-RAN system.

9 Conclusion

C-RAN and mobile fronthaul are emerging technologies that helps to address the ongoing rapid growth of mobile internet traffic demand. It is also estimated that the

processing power necessary to run the mobile network as a whole is cut as much as 40 percent by C-RAN. Though, C-RAN is predicated on having bunch of dark fibers and have strict latency and synchronization requirements, it offers many advantages as listed in section 5. This paper has briefly dealt with the various deployment use cases that could be adopted as fronthaul solutions in case of limited fibers. Enhanced and colored fiber connection has advantages over traditional fiber connection and widespread C-RAN deployment is feasible. Some advantages of C-RAN were also highlighted in this paper.

References

- [1] Halachmi, N.: How C-RAN architecture can reduce costs for mobile backhaul deployments. Telco Systems (November 18, 2011)
- [2] Rayal, F.: Cloud RAN vs. Picocells: The Need for Integrative Approach in Next Generation Network Design (April 18, 2013), <http://frankrayal.com/2013/04/18/cloud-ran-vs-picocells-the-need-for-integrative-approach-in-next-generation-network-design/>
- [3] The HetNet Bible (Small Cells and Carrier WiFi) - Opportunities, Challenges, Strategies and Forecasts: 2013 – 2020 – With an Evaluation of DAS & Cloud RAN. Signals and Systems Telecom, Market Publishers Ltd. London, UK (PRWEB) (June 12, 2013), http://marketpublishers.com/report/wireless_technology/wifi-wimax/hetnet-bible-small-cells-n-carrier-wifi-opportunities-challenges-strategies-n-forecasts-2013-2020-with-an-evaluation-of-das-cloud.html
- [4] Merritt, R.: China drives servers to base station role (January 3, 2012), http://www.eetimes.com/document.asp?doc_id=1261296
- [5] Kaiwei, H.: Building New-Generation GSM Networks with C-RAN, Article No.5 (September 19, 2011), <http://www.en.zte.com.cn/en/about/publications/>
- [6] EE Daily News, Will 4G wireless networks move base stations to the cloud? (June 20, 2011), <http://www.eedailynews.com/2011/06/will-4g-wireless-networks-move.html>
- [7] Cloud technologies to improve performance and efficiency of mobile networks (July 11, 2013), <http://phys.org/news/2013-07-cloud-technologies-efficiency-mobile-networks.html>
- [8] CPRI, Cpri specification v4.0 (June 30, 2008), <http://www.cpri.info/downloads/CPRIv402008-06-30.pdf>
- [9] Zhu, Z., et al.: Virtual Base Station Pool: Towards A Wireless Network Cloud for Radio Access Networks. IBM Research
- [10] Gabriel, C.: Korea Telecom plans world's first commercial Cloud-RAN (December 8, 2011), <http://www.rethink-wireless.com/2011/12/08/korea-telecom-plans-worlds-commercial-cloud-ran.htm>

- [11] Baldry, J.: Mobile fronthaul -mobile's new kid on the transport network block. Transmode, RCR Wireless (July 22, 2013), <http://www.rcrwireless.com/article/20130722/opinion/reader-forum-mobile-fronthaul-mobiles-new-kid-transport-network-block/>
- [12] Akhter, M.: Reader Forum: Front-haul compression for emerging C-RAN and small cell networks. Integrated Device Technology Online (April 22, 2013)
- [13] Fitchard, K.: Intel's next big wireless play: It's not smartphones. GIAOM Online Magazine (January 23, 2012), <https://gigaom.com/.../intels-next-big-wireless-play-its-not-smartphones/>
- [14] ASOCS and CMCC Team up to Deliver Mass Market High, Commercially Viable C-RAN Solutions (February 25, 2013), <http://uk.reuters.com/article/2013/02/25/bc-asocs-idUSnPn2255113+100+PRN20130225>
- [15] Madden, J.: Cloud RAN or small cells?. Fierce Broadband Wireless Online Magazine (April 30, 2013)
- [16] <http://www.nttdocomo.com/info/index2.html>
- [17] Ghadialy, Z.: A presentation looking at Small Cell Standardization in 3GPP Rel-12 (February 20, 2013), <http://blog.3g4g.co.uk/2013/02/small-cell-standardization-in-3gpp.html>
- [18] NGMN Alliance, Suggestions on Potential Solutions to C-RAN, ver 4.0 (January 03, 2013)
- [19] Chen, C.: C-RAN: the Road Towards Green Radio Access Network, ver 2.5, China Mobile (October 2012)
- [20] Huitao, W., Yong, Z.: C-RAN Bearer Network Solution (November 18, 2011), http://www.zte.com.cn/endata/magazine/zte technologies/2011/No6/articles/201111/t20111118_263975.html
- [21] Transmode, Mobile Fronthaul, <http://www.transmode.com>
- [22] Fitchard, K.: Intel's next big wireless play: It's not smartphones (January 23, 2012), <http://gigaom.com/2012/01/23/intels-next-big-wireless-play-its-not-smartphones/>
- [23] Donega, M.: LR Mobile News Analysis China Mobile Steps Up Cloud RAN Efforts (September 03, 2013), <http://www.lightreading.com/mobile/china-mobile-steps-up-cloud-ran-efforts/>
- [24] Flanagan, T.: Creating Cloud base stations with TI's Key Stone multicore architecture. White Paper, Texas Instruments (October 2011)
- [25] Akhter, M.: Front-haul compression for emerging C-RAN and small cell networks. RCR Wireless Online (April 22, 2013), <http://www.rcrwireless.com/article/20130422/opinion/reader-forum-front-haul-compression-emerging-c-ran-small-cell-networks/>
- [26] Rayal, F.: Cloud RAN vs. Small Cells: Trading Processing for Transport Cost (March 17, 2012), <http://frankrayal.com/2012/03/17/cloud-ran-vs-small-cells-trading-processing-for-transport-cost/>
- [27] Merritt, R.: Cloud RAN Attracts Asian, European Carriers (February 7, 2013), http://www.eetimes.com/document.asp?doc_id=1318782

Reliable Data Transmission for Multi-source Multi-sink Wireless Sensor Networks

Kassahun Tamir¹ and Menore Tekeba²

¹ Adama Science and Technology University,
Adama, Ethiopia

kassahuntamir@yahoo.com

² Addis Ababa Institute of Technology,
Addis Ababa, Ethiopia

menoretekeba@yahoo.com

Abstract. Multi-source multi-sink wireless sensor networks (WSNs) have got variety of application in areas that need to detect multiple environmental monitored parameters using a single sensor field/network. This type of WSN, beside its limited resources such as energy limitation like any WSNs are, has unique characteristics for instance network congestion, since multiple sensor nodes can send data to multiple/single sink at the same time. Research works on WSNs are not matured and well-developed. There are a number of research issues that are not yet addressed. In particular reliable data transmission is among those that are the decisive and unsolved ones as much as required. In this research work, reliable data transmission is ensured using hop-by-hop(in every hop) loss detection and recovery and it is tested using NS2 simulator. This paper provides 100% reliable data transmission. Reliability is assured by employing hop-by-hop loss detection and recovery using a hybrid of NACK and ACK-based approach. Since NACK-based scheme cannot handle the unique case where all packets in a communication are lost, we used a last single ACK feedback to make the sender sure that all packets are received successfully. The simulation and experimental results show that the new protocols perform well under various conditions and protocol parameter settings.

Keywords: multi-source multi-sink wireless sensor networks, reliable data transmission, wireless sensor networks.

1 Introduction

A wireless sensor network (WSN) consists of a group of self-organizing, lightweight sensor nodes that are used to cooperatively monitor physical or environmental conditions. Commonly monitored parameters include temperature, sound, humidity, vibration, pressure, motion and so on. Although WSN research was initially motivated by military applications, WSNs are now used in many industrial and public service areas including traffic monitoring, weather conditions monitoring, video

surveillance, industrial automation and healthcare applications [1]. A basic sensor node includes five main components: an embedded processor, memory, radio transceiver, sensor(s), and a power source [2]. Because of the size and cost constraints on sensor nodes, they are limited by energy, bandwidth, memory and other resources. To overcome problems that could come due to these limitations of sensor nodes appropriate design of protocol is needed.

The rest of this section is organized as follows: section 1.1 is an introduction about multi-source multi-destination scheme. Section 1.2 is all about reliable data delivery in a specific wireless sensor network type i.e. multi-source multi-destination WSN. The motivation of this paper is discussed in Section 1.3. Finally this research organization is presented in section 1.4.

1.1 Multi-source Multi-sink Wireless Sensor Networks

In WSNs, sensor nodes sense phenomenon or phenomena and process it. After processing the phenomenon or phenomena that they (multiple sources or sensor nodes) sense from the environment, they have to send the processed data to the sink or sinks. Sinks are base stations at which all the sensed data are collected that have been sent from each sensor nodes. Unlike sensor nodes, sinks are not resource limited such as energy and buffer, so they can be active all the time.

We can categorize wireless sensor networks according to the number of sources and sinks they have. If a wireless sensor network has only one sink and many source nodes that can transmit at the same time, we call it multi-source single-destination wireless sensor network. But in such types of networks if we restrict the network in such a way that one source node can send at a time, we call such type of networks single-source single-sink WSNs as shown in Figure 1. On the other hand, if it has multiple sinks with multiple sensor nodes (sources), we call it multi-source multi-sink wireless sensor network as shown in Figure 2. This network architecture is obviously required when the same WSN is serving multiple applications [19], each running on distinct devices. However, the need for multiple sinks arises also in other situations. For instance, researchers are increasingly investigating the use of actuator nodes in WSNs [3]. Different actuators (sinks) take decision after they analyze or process the information gathered by sensors from the physical world and then perform appropriate actions upon the environment, which allows a user to effectively sense and act from a distance [4]. Moreover, multiple sinks are required when using distributed in-network data mining in order to discover frequent event patterns rather than transmitting raw data streams [5]. As can be seen from the diagram below, the second one (multi-source multi-destination scheme of WSN) is complex. Since this is the case, data collision and data loss while transmitting are likely to occur; this means, such networks have reliability problem and the follow-on data recovery expends much energy of sensor nodes, energy inefficient, unless appropriate protocol is designed to overcome such limitations. So, designing reliable and energy efficient protocol for such WSNs is crucial. In this particular work, only reliable data transmission is considered and we will design an energy efficient protocol which can't affect a 100 % reliable data transmission obtained in here.

1.2 Reliable Data Delivery

Reliable data transmission is an important issue not only in WSNs but also in wireless communication as a whole. Having constrained resources in WSNs makes reliability the very critical and decisive issue. Reliability can be defined as the number of unique data packets received by the receiver divided by the number of data packets queued to be transmitted by the sender.

There exist many applications in Wireless Sensor Networks requiring all data to be transmitted without loss. For example, structure monitoring needs the entire data from all measuring points to build a model and analyze it. Consider a sensor network deployed in a chemical plant to detect harmful gas. It is crucial for sensor nodes to reliably transport every sensor reading back to the sink. Other critical WSN applications such as biological monitoring, health care monitoring, and battlefield surveillance also require high end-to-end reliability [1], [6]. On the other hand, some applications may not require simple 100% guaranteed transmission of data packets [7].

In WSNs, there are plenty of factors that make the system unreliable. Such as node failure due to lack of energy source or malfunctioning of any part of it, as it is a microelectronic device. Moreover, since bandwidth of wireless networks is very narrow, congestion will occur while multiple sensor nodes attempt to transmit at the same time. Problems related to such deficiencies make designing reliable protocol challenging. So, to make sure whether the data sent by the source to the destination delivered without any data loss and as much as possible

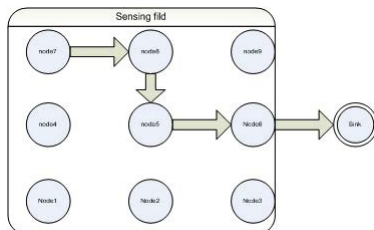


Fig. 1. A sensor network with single sink

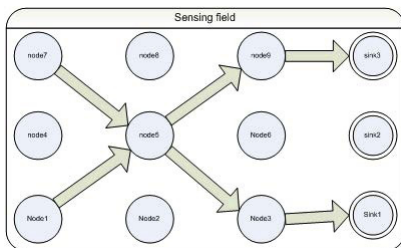


Fig. 2. A sensor network with multiple sinks

to recover from loss, if it is found that there is a data loss, a receiver feedback and sender retransmission mechanism is usually used in WSNs. There are two commonly used receiver feedback mechanisms: the ACK-based approach and the NACK-based approach. In an ACK-based approach, the receiver positively acknowledges receipt of data packets, while in a NACK-based approach, the receiver only returns feedback to the sender if it detects a packet loss. The NACK-based approach incurs less overhead than the ACK-based approach, but a common problem for the NACK-based approach is that it cannot detect when all data packets are lost, since packet loss detection is based on observing gaps in the packet flow.

Reliable data transport in WSNs is a multifaceted problem influenced by the physical, MAC, network, and transport communication layers. Traditional TCP/UDP protocols and end-to-end error detection and recovery mechanism cannot be implemented directly in WSNs because of having insufficient resources. This problem can be overcome by using hop-by-hop error detection and recovery in link layer and by designing appropriate transport control protocol. Some existing transport protocols that are designed for the provision of reliability in WSNs are [9,10,11,12]. Reliability can be provided by designing a reliable routing protocol in routing layer. For example, transmitting multiple copies of a single packet in different paths to the sink increase the probability of delivery of packets without losses [13]. RTS/CTS handshake, MAC layer acknowledgement and randomized slot selection are some of the reliability mechanisms provided by IEEE 802.11 in MAC layer [14]. Using judiciously chosen linear error correcting code, intermediate nodes in a wireless network can directly recover linear combinations of the packets from the observed noisy superposition of transmitted signals [15].

1.3 Motivation

Even though their application areas are very wide, a network of multiple sources and multiple sinks was not considered in many research works. Many researchers' analysis are even considering simple network structures or topologies such as star, tree or linear topologies [8]. In this research work, the reliable data transmission for multi-source multi-sink network is taken as the main objective. In this type of WSNs, when multiple sensor nodes attempt to transmit at the same time, collision of packets would occur. If collision occurs, there will be a data loss. So, reliable protocol that detect a data loss and recover from them is necessary.

1.4 Research Organization

In this section, introduction about WSNs, particularly about multiple sources multiple sink type WSNs, is presented. In section two, recent related works that have been done related to reliable data transmission is discussed. The design of the new protocols is presented in section three. In section four, simulation results and performance evaluation are discussed. At last, the conclusion and possible future work are described.

2 Literature Review

Due to the fact that WSNs are recent technologies [15], research works on them are not matured and well-developed. There are a number of research issues that are not yet addressed. On the other hand, there are a number of research works done on WSNs. Specially, reliable data transmission is a main issue for WSNs. Since it is routine to cover all works that are related to this issue, here are some of the researches that are recent and closely related with this research work. So, in this section, literatures that are related to reliable data transmission protocols are presented.

Reliable data transmission can be provided using receiver feedback mechanisms: ACK-based or NACK-based approaches or by any other mechanisms. In ACK-based approach or commonly known as TCP/IP (if it is implemented in transport layer) the data receiver sends the feedback packet (ACK) if it has successfully received the data packet. When the data packet is sent to any receiver the time stamp will be set for its arrival and if the data packet sender doesn't receive the ACK packet from the receiver within a predefined period, the sender will consider as if the data packet has been lost during the transmission and will resend the unacknowledged packet. In the contrary, the receiver of NACK-based approach will send NACK packet to the data sender if and only if it detects a packet loss.

In most wired and wireless networks that are full of resources such as energy and bandwidth TCP (ACK-based) approach is used. Since the receiver sends ACK for the sender to approve the successful arrival of each and every packet, it exploits the bandwidth and energy of the nodes of a network, especially, when the data packet size is very small. So it is not advisable to use this approach for WSNs, because they usually have insufficient resource. The other drawback of this approach is the high feedback overhead problem. Specially, when the senders are sending for multiple receivers at a time, the acknowledgement that comes from each receiver will create congestion and data loss.

Generally, it is not advisable to employ traditional TCP/IP reliable protocols in WSNs because each sensor node in a WSN has very limited power, bandwidth and storage space and has to cope with a lossy wireless channel. Therefore, the reliable data transport protocols that are widely used in internet i.e. TCP are not suitable for WSNs. Even the modified TCP, hop-by-hop TCP protocol, used in [17] has the problem of traditional TCP which is used in wired and wireless networks that have not deficiencies of the aforementioned resources.

On the other hand, NACK-based approach is more effective than ACK-based, because it only generates an extra packet when data loss occurs. In NACK-based loss detection and recovery scheme, extra packet (NACK) overhead will occur unless it is carefully designed. Another typical NACK problem is the loss of all data packets. In a NACK-based scheme, the receiver can detect and report packet loss only if it is aware of the incoming packet. Thus, a NACK-based scheme cannot handle the unique case where all packets in a communication are lost.

The work done in [1] proposes a new reliable data delivery protocol for general point-to-point data delivery (uni-casting) in WSNs. The proposed protocol adopts a NACK-based hop-by-hop loss detection and recovery scheme using end-to-end sequence numbers. In order to solve the single/last packet problem in the NACK-based approach, a hybrid ACK/NACK scheme is proposed where an ACK-based approach is used as a supplement to the NACK-based approach to solve the single/last packet problem.

However, the wireless sensor network used in the above work is a single-source single-sink type network. In other words, in a given sensor network or sensor field, at a time, there is only one source which can send data only to one sink. But in real world there are many application areas that employ sensor networks in which multiple sources can send to a sink(s) at the same time.

Reliable information forwarding in wireless sensor networks can also be provided by sending multiple copies of a single packet along multiple paths. In [13] ReInForM is proposed to attain desired reliability. It uses redundant copies of a packet to increase its end-to-end probability of data delivery. The degree of redundancy introduced, is controlled using the desired reliability, the local channel error conditions, neighborhood information available at each node. Hence, it relies heavily on the existence of multiple paths from source to sink. On generating a packet, the source node determines the importance of the information it contains and decides the desired reliability for it. It also knows the local channel error. Such a multi-path approach would succeed only if there are sufficiently large number of paths exists from source to sink. Moreover, the lengths of these paths should be close to that of the shortest path. The reliability that can be achieved using multi-path forwarding depends on the expected number of such paths.

When a source node sends multiple copies of a single packet through multiple paths, there would be network congestion. Specially, when multiple sources are allowed to transmit at a time, employing ReInForM for such type WSNs, packet collision would be uncontrolled [13].

In some network protocols like in [16], sending the sensed data to a sink may be receiver initiated. Whenever a sink node wants to process, to take action after processing the data or for any other reason, it sends initiation message for the source node. It provides end-to-end reliability by using end-to-end acknowledgments. In this protocol there are four main problems that are: first, since the data transmission takes place after sink node initiation, the urgent data that has to be transmitted to the sink immediately, will not be transmitted until sink node initiate data transmission. Second, only one source node can send at a time and there should be only a single path to a given sink. So, for the cases that multiple source nodes need to send for a single/multiple sink node(s), it doesn't work. Third, reliability is provided by using end-to-end acknowledgement as a result data loss recovery time is very high. At last this protocol employs only NACK-based approach.

Reliability of the network can be provided a supplement in different layers beside that are reviewed in above. For example, the collision avoidance and

detection method of IEEE 802.11 MAC protocol are loss prevention tools in MAC layer. Whenever any sensor node senses something from the environment, it will not send it immediately to the next hop. Instead, it has to identify whether the medium is busy.

Finally, in this research work, we tried to assure reliable data transmission for highly congested networks by improving problems mentioned above.

3 Protocol Design

This section presents the design procedures for hop-by-hop reliable data delivery and its implementation in NS2 simulator.

3.1 Hop-by-hop Reliable Data Delivery

To deliver reliable data from multiple sources to multiple sinks, our protocol employs hop-by-hop loss detection and recovery scheme. Hop-by-hop loss detection and recovery is achieved by using both negative acknowledgement (NACK-based) and last single positive acknowledgement (ACK-based) approaches.

New Packet Creation in NS2. The new packet added in this research work is only the NACK packet (its packet type and its symbolic name is created in `ns/packet.h`). But the packet type ACK is already created previously. The header format of the ACK and the NACK packets, the basic header fields, and their importance in error detection and recovery are discussed in the following section.

Packet Format. For loss detection and recovery ACK, NACK and DATA packets are used. Data packets are the sensed data from the environment or what the sensor source node wants to send for any sink node that needs the data for further processing/action. All data packets share the same header format. The header of a data packet contains five important fields: `Packet_ID`, `ACK/NACK_Bit`, `Source_ID`, `Receiver_ID`, and `Sender_ID`. Their meaning and their use are presented below.

Packet_ID. This term is used to refer to the identity of packets and it is assumed to be unique for all packets between a particular source/destination pair.

ACK/NACK_Bit. `ACK/NACK_Bit` is stored in the header of a data packet. The possible value of this term is 1 or 0 to indicate whether the sender requests the receiver to respond ACK or NACK for the sent packet respectively. If the packet which is being sent is the last one, the sender will set the value of this bit 1

and for packets that are other than the last packet the value of this bit is set to 0.

Source_ID. Source node is a node at which packets are originally generated. When a single receiver node is receiving packets from multiple source nodes at the same time, Source_ID is used to check in order arrival of packets that come from the same source node, the value of Source_ID is very crucial.

Sender_ID. and **Receiver_ID.** When acknowledgement needed to be sent for the sender of the data packets, the receiver of data packets sets the Sender_ID and Receiver_ID of feedback packets the value of its ID and Sender_ID of the data packets respectively.

Feedback packets (ACK and NACK) share the same header format. The header of these packets contains four important fields: Packet_ID, PacketType, Receiver_ID and LastRecieved_ID.

Packet_ID. It is the similar to that of Packet_ID field of data packets.

PacketType. Used to identify whether data, feedback or control packets are reciving and to take appropriate action.

LastRecieved_ID. when the packet arrival is out of order, the receiver of data packets will set the value of LastRecieved_ID field in the header of NACK packet by the sequence number of the last packet successfully received.

3.2 Protocol Operation

Old Protocol. The link layer protocol that exists before the addition of our protocol

New Protocol. The new link layer reliable protocol

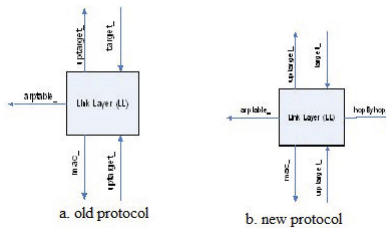


Fig. 3. The architecture of old&new link layer in NS2

Old Protocol Operation. The general architecture of mobile nodes and its description are presented in Figure 4 and section 4.1 respectively. Particularly in this section, the operation of old link layer (LL class in NS2) is graphically presented in Figure 3a.

Sender Operation. As the data packet(s) comes from the upper layer, link layer buffers the packets until their turn to be transmitted is reached and send them to the lower layer according to the default queuing principle used i.e. first in first out (FIFO). The queue length is told during nodes configuration in Tel language as:

```
$ns node-config -ifqLen $val(ifqlen)
```

The life time of these packets in one's node link layer is until it will be transmitted to the lower layer.

If the packet is a routing packet, this layer will consult ARP (address resolution protocol) using a pointer *arptable_*.

Receiver Operation. When the packets, whether they are data or control packets, come from the lower layer, the only responsibility of link layer is to send to the next upper layer that is network layer using a pointer *upertarget_*.

New Protocol Operation. Unlike the old link layer protocol, in our new protocol during both sending and receiving packets, there are additional responsibilities that are detecting a loss and recover from them. The operation of detecting the loss and recovery are explained in sending and receiving operations discussed below. The functions or methods that compute the intended activities are all included in *hop-by-hop.cc* and *hop-by-hop.h* C++ files in NS2. And whenever the packet is arrived in the receiving function (*receive(p,h)*) of link layer (*LL.cc*), functions in the aforementioned files are consulted by using *hopByhop_* pointer.

Sender Operation. The sender node may be a source of data to be transmitted or an intermediate/router node. Whether the sender is a source or an intermediate node, the main responsibilities of the sender node are the following:

- It copies every packet before transmitting. All packets copied in this step will be deleted when ACK is received. When NACK is received all packets with sequence number less than the sequence number of the lost packet will be deleted.
- If the packet being sent is other than the last packet, the sender will fill the ACK/NACK field to 0 otherwise it will fill 1.
- When transmitting the last packet, besides filling ACK/NACK field to 1, the sender sets a time stamp. If ACK packet is not received before time out, the sender will understand that all packets that are after the last negatively acknowledged packet are lost and it will retransmit the whole data packets starting from the last negatively acknowledged one. But if there is no any

negatively acknowledged packet and if ACK is not received, it will transmit the whole packets.

- If NACK packet is received, it extracts the LastReceieved_ID field and it will retransmit packets with sequence number greater than or equal to the value of LastReceieved_ID to the node that sends NACK packet.
- If ACK packet is received, it will understand that all data packets are successfully delivered and it will free the buffer by deleting all packets.

Receiver Operation. The receiver may receive packets in link layer from upper or lower layers. The node is receiving if packets are received in link layer from lower layers and then there will be two main responsibilities of this layer. Firstly, it checks in order arrival of packets. Secondly, it sends feedback packets. Sending feedback packets depends on the result of arrival checking. If there is no loss of packets, the receiving node will send ACK packet when the last packet with ACK/NACK field set to 1 is received. And if there is loss of packets, the receiver node will prepare NACK packet with a sequence number or LastReceieved_ID equal to a sequence number of the last successfully received packet.

4 Simulation Results and Performance Evaluation

The new protocol, hop-by-hop loss detection and recovery, is developed in Link layer of NS2 simulator. As described in the previous portions, the selected sensor field or network is a grid topology in which multiple sensors can send to multiple sinks at the same time without waiting each other and decline in reliability of data transmission. The reliability of data transmission in this network is assured by employing hop-by-hop loss detection and recovery. Hop-by-hop loss detection and recovery is done by using advantages of both ACK and NACK-based approaches.

4.1 Network Simulator 2 (NS2)

Network simulator 2 (version 2) is most commonly known as NS2. It is an event-driven simulator and it is used to simulate both wired and wireless networks. Due to its flexibility, modular nature and open source it is widely used by many researchers [18]. The internal structure/architecture of mobile nodes in NS2 is as shown in Figure 4. When ever needed to add special features or new modules, it is easy and helpful. Its front end is object oriented tk/tcl language or commonly known as OTcl and its back end is c++ language. Our extension or modification is done in c++ part of the simulator.

New Protocol Implementaion in NS2. Our new reliable data transmission protocol is developed in link layer (LL.cc) of NS2. All the features/functions of the new protocol are developed in this file.

The OTcl Simulation Script. This is the part that the activity of the network creation, configuration and parameters settings can be done. As tried to explain in above, NS2 is written in two languages of which OTcl front end language is the one. The basic parameters that are summarized in Table 1 are set in OTcl script.

Table 1. Network configuration and parameter values

MAC type	IEEE 802.11
Buffer Size	500
Number of Nodes	Up to 25
Transport Layer	UDP
Packet Type	CBR
Packet Size	1000 bytes
Inter Sending	0.005ms
Number of Packets	10(for each source)
Topography Area	1000mx1000 m
Distance between nodes	200 m

4.2 Evaluation Methodology

The performance of the new protocol is evaluated from different perspectives. Evaluation metrics and parameters which are used in result analysis are defined. The definitions of evaluation metrics from 1-6 are taken from [1].

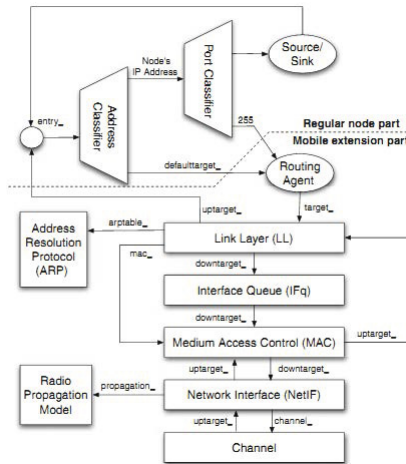


Fig. 4. The architecture of mobile nodes [18]

Evaluation Metrics. The main goal of this research work is designing a reliable data transmission protocol for a multi-source multi-sink wireless sensor networks. In addition to that, extrinsic values such as latency and throughput of a network are also considered.

1. *End-to-end Delay:* The end-to-end delay is measured as the interval between the generation of a data packet at its source and the reception of that packet at the sink. The average end-to-end delay for each source node is calculated as the average end-to-end delay of all data packets generated by that node. The end-to-end delay shows the average amount of time it takes for the network to deliver a data packet from a particular source node to the sink.

2. *Link Delay:* The link delay measures the interval from when a packet is sent from the sender to the time it is received at the next hop. Comparing the link delays is useful for understanding the network congestion level at each link as well as the impact of traffic load on a packet's link delay.

3. *End-to-end Reliability:* The end-to-end reliability for each source node is defined as the number of data packets from the node that are received at the sink divided by the total number of data packets the node generates. The end-to-end reliability reflects the reliability of a given path in the network.

4. *Link Level Reliability:* - The link-level reliability measures the reliability of the link between two adjacent nodes. It is defined as the number of unique data packets received by the receiver divided by the number of data packets queued to be transmitted by the sender at every hop.

5. *End-to-end Throughput:* - The total throughput is measured as the number of unique data packets received at the sink divided by the time interval between when the first data packet is generated and the last packet is received. The achievable total throughput reflects the efficiency of the protocol. The higher the achievable total throughput, the faster source nodes can deliver their data packets to the sink. Both the end-to-end reliability and the end-to-end delay can affect the total throughput.

6. *Link level Throughput:* - The link throughput measures the throughput between two neighbor nodes. Link throughput is calculated as the number of unique data packets received at the receiver divided by the time interval between when the first data packet is generated by the sender and the last packet is received by the receiver.

The reliable protocol developed in this research work is fully efficient in both end-to-end and link level reliability measurements. The network performance regarding throughput and delay is tested: but only its end-to-end values, since, most of the time, these metrics are decisive factors and the network is intended to be the one that its end-to-end delay is as small as expected.

Evaluation Parameters. The performance of the new protocol is tested by changing the following parameters. As the number of sources that can transmit at the same time increase, collision and data loss will increase proportionally. Increasing the number of hops between sinks and sources can also increase the loss of packets. In addition to that when direction of transmission is changed,

particularly for X (cross) transmission direction the probability of occurrence of packet collision is very high. In our simulation, we used a random source generator and we saw the result when neighboring, one hop apart and the like nodes send at the same time. So, the performance of our protocol is evaluated by these worst cases. The parameters used in our simulation are presented as follows.

1. *Number of sources*: - The maximum number of nodes used in our experiment are 25. To show the effect of number of sources in reliable data transmission, we used two and three source nodes in different directions of transmission i.e. when two nodes transmit in parallel and X-direction.

2. *Number of hops*: - the total number of sensor nodes between a sink and source node in the path of transmission plus one is called number of hops. In our experiment from 1 to 4 hops are used when the transmission is in parallel and maximum of 8 hops are use in X direction.

3. *Direction of Transmission (position of sources and sinks)*: - when two or more sensor nodes send the sensed data to the sink which is found in the same row is called parallel transmission. But when two source sensor nodes send the sensed data to sinks which is found in different rows and opposite direction while transmitting, we called it X-transmission(Figure 12). In our experiment both of them were implemented.

4. *Random source*: - In reality, the position of multiple sources in a network is indeterminate. So to make our work realistic, we made the positions of source nodes random i.e. when no, one, two or more nodes are in between.

4.3 Basic Tests of Protocol Performance

The new protocol is tested from different perspectives and in all cases it is 100% reliable and can convey the intended information or data to multiple sinks without any data lose. The quality of the network i.e. throughput and delay are tested by considering different conditions beside the main objective of the network. We have seen the effect of the number of source nodes in a network i.e. scalability test in section 4.3.1. At last, In section 4.3.2, interference and collision test is analyzed.

Scalability Test. The number of sensor nodes in WSNs is proportional with area (sensor field) needed to be controlled. Then when the sensor field is wider the follow-on increment of number of sensor nodes has a direct negative effect on reliability of the data transmission. In other words, when the number of sensor nodes in a given sensor field become many, the high traffic occurrence in a network will be higher than when the number of nodes in a sensor field are less. The follow-on collision among packets will let packets to be corrupted or lost. So the reliable protocol should withstand such type characteristics of WSNs and should recover the lost data (if any). Particularly, our protocol has been tested by increasing the number of sensor nodes up to 25 and by increasing the number of source nodes up to 3 that can possibly transmit at the same time.

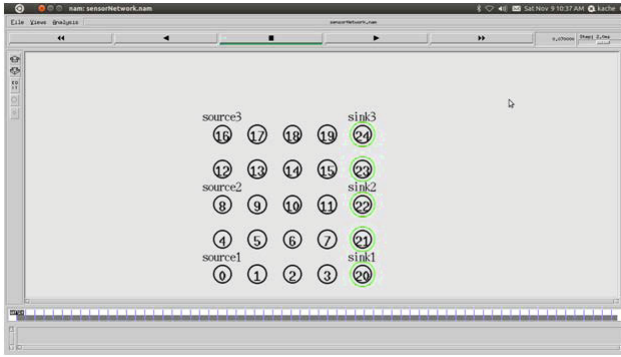


Fig. 5. A multi-source sensor network using NS2

Single Source Single Sink: When a single source is sending to a single sink, no packet loss is observed in our simulation. This is because of there is no network congestion, no packet collision and as a result there is no packet or data loss. But in real world data/packet loss is not limited only if there is collision of packets but it is also may be data packets are corrupted. So to test the performance of the new protocol for this particular case, we dropped packets deliberately that are being transmitted with some time intervals and we call the output of these scenarios *before and after* respectively. Finally, the system has shown that it is reliable and it can recover all the lost packets successfully. Figure 6 & Figure 7, respectively, show that the end-to-end delay and the total throughput of a single source network.

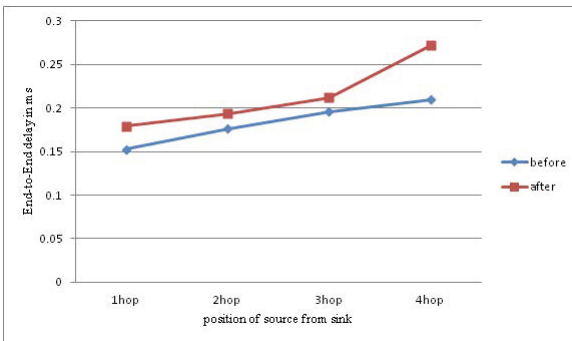


Fig. 6. End-to-end delay in ms

From those network performance figures we can observe that, as the source node is further and further away from the sink the end-to-end delay increases and the corresponding total throughput decreases.

Since end-to-end delay of a packet is the difference between its reception and sending times and this time difference is directly proportional with the distance that the packet covers, its value becomes bigger and bigger as the distance between the source and the sink is getting further and further apart. And as throughput is a total number of successfully delivered packets per the time difference between the reception time of the last packet and the sending time of the first packet, its value decreases as this time difference value getting bigger.

On the other hand, when fault is injected to a network both the end-to-end delay and total throughput values depend on not only the distance between the source and the sink but also the time taken to recover a packet(s). That is all what we can observe from Figure 6 and Figure 7 i.e. the *after* graph.

Multiple Sources Multiple Sinks. In the case of the presence of multiple sources send their data to their corresponding sinks, there is no need to use error model because the atmosphere is exposed to packet collision. However, packet loss due to collision is not always happens when two or more sources send a data. So, if needed to test the protocol when loss is always there, it is possible to use error model as an additional. In this particular work, the protocol is tested when even neighboring nodes and nodes in X-direction are sending. Figure 8 & Figure 9 are for networks with two sources that can send their data at the same time.

By further increasing the number of sources to 3 in the network as can be seen in Figure 5, we have tested the performance of our protocol i.e. its performance of error detection and recovery. And we analyzed that whatsoever the numbers of sources are in a network, the protocol is always 100% reliable. Figure 10 & Figure 11 show that the quality of the network beside reliability.

In general, when two or more sensor source nodes are sending in parallel direction that are not neighbors there may or may not be a collision of data packets. This means that packets may or may not be lost or its occurrence is not always. During our simulation, its observed that, its recovery is also very fast and the overall end-to-end delay is as small as required.



Fig. 7. Total throughput



Fig. 8. End-to-end delay when two sources send to two sinks at a time in ms

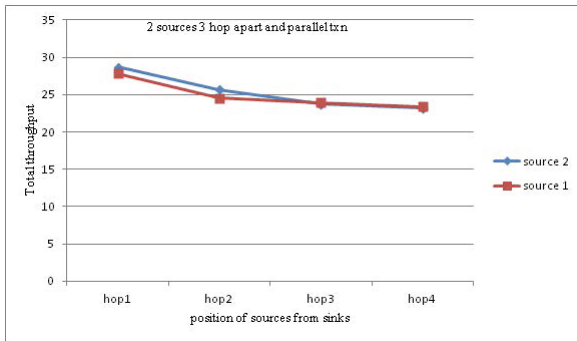


Fig. 9. Total throughput when two sources send to two sinks at a time



Fig. 10. End-to-end delay when three sources send to three sinks at a time in ms

When three sources are sending at the same time sometimes one of the three sources may wait until the other two completes transmitting their data successfully. This job is handled by a collision avoidance principle of MAC 802.11.

Interference and Collision Test. So far, we have shown the effect of number of sources and the quality of our protocol. When we were varying the number of source nodes, the gap among sources were two and three hops. The direction of transmission was only in parallel.

The minimum distanc among sources was a value of two hops but it is still a big distance and less congestion relative to the case when two neighboring nodes are sending. So, the later scenario is taken as one of the worst cases to evaluate reliability, interference and collision tests.

Sources position relative to each other are changed from neighboring up until three hop apart. So, in real world there will not be other than this worst case conditions.

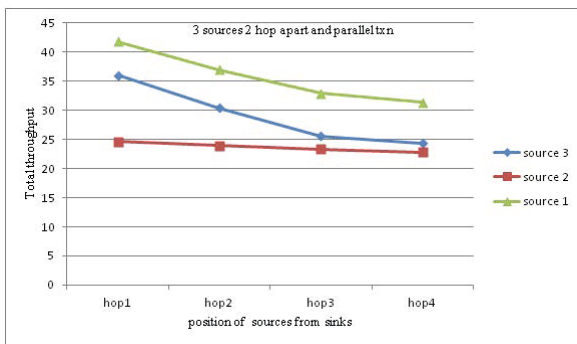


Fig. 11. Total throughput when three sources send to three sinks at a time

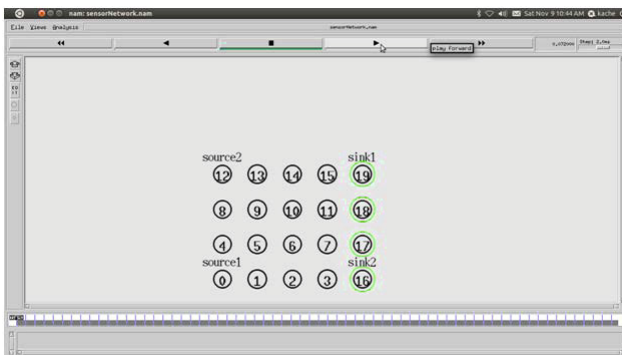


Fig. 12. X-direction transmission

The second and the last worst case is when sources are sending to the X-direction as shown in Figure 12. Source with address of 0 is sending to sink node with address of 19 and source node with address of 12 is sending to sink with address of 16.

Varying Position of Sources. In here we let two neighboring nodes send at a time and during our simulation we observed four different cases for this particular scenario.

If there are other free sensor nodes that are other than the sending nodes and nodes in their path, the sending nodes mostly prefers to use free nodes and try to reduce packet collision that will happen when they are sending through a restricted path.

We restricted sending neighboring nodes not to use alternative path and we simulate the scenario with a two row network. After many repeated simulation results, we could observe three different situations.

- Either one of the sending nodes waits and sends its data after another sending node completes data transmission successfully. In this condition most of the time the frequency of data loss is very small. Or,
- When one of the sending nodes is processing or doing something, the second node will start sending.
- Or, in the worst case, when two sending nodes start sending at the same time, packet collision will be sever and the frequency of packet collision and recovery will be very high. However, it is still fully reliable.

In general, the end-to-end delay is very long time duration as compared with the end-to-end delays that have been seen in the previous senarios and that of the total throughput is very much smaller than the previous values of the same type as shown in Figure 13 and Figure 14 respectively.

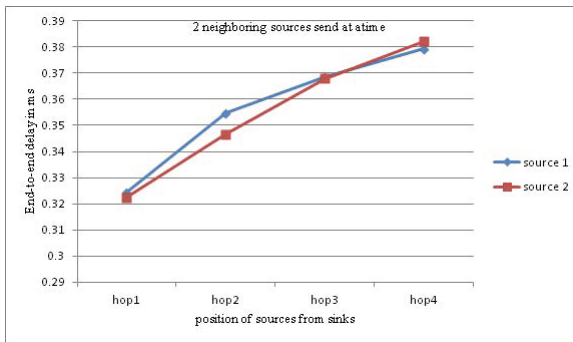


Fig. 13. End-to-end when two neighbouring sources send to two neighboring sinks at a time in ms

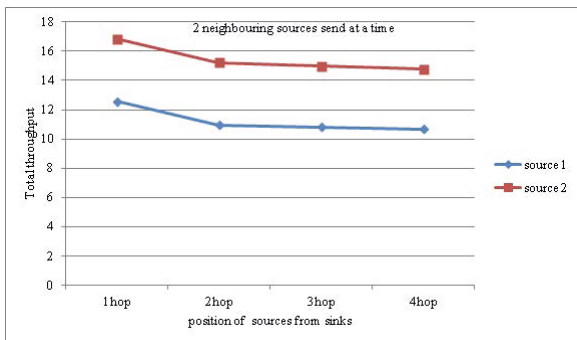


Fig. 14. Total throughput when two neighbouring sources send to two neighboring sinks at a time

Varying Direction of Transmission. As can be seen from Figure 12 two sources are transmitting to the opposite direction (X-direction). In this type of transmission there are some nodes that receive packets from different nodes at the same time. Sensor nodes cannot receive packets from different nodes at the same time. If they try, or forcefully obliged to receive, there will be collision and as a result there will be a data loss. If this is the case, our protocol has to handle such type problems by detecting and retransmitting the lost packets. If there is a loss, it will automatically recover the lost data from the previous hop.

The reason why we don't draw some graphs that show the network performance during X-direction transmission is because, during simulation, it is observed that one of the sending nodes should wait until the other sending node completes transmission successfully. Because there will be one node in common for two paths if they attempt to transmit at the same time. If this is the case, almost all packets will be dropped. The protocol could detect but there will be an all time job of loss and recovery. But thanks to the collision avoidance mechanism of MAC layer, only the first source which receives the CTS message will start sending and the second will be the next after the sending node completes transmitting all data packets.

At last, the reasons why we didn't provide a comparison result with ours are: first we are unable to get works done on multi-source multi-sink WSNs. Second, other reliability approaches are not recommended for WSNs (for example TCP/IP).

5 Conclusion and Future Work

5.1 Conclusion

Wireless sensor networks have scarcity of many resources such as narrow bandwidth, small size and non-rechargeable battery, and small size memory. As a result of that, there are issues that can be raised when we think about WSNs.

Among them reliable data transmission is the very crucial and decisive one. In most research works simple network topologies are considered for experimental analysis, for instance, single-source single-sink WSN. However, real world may be different.

Since WSNs are recent technologies [15] their deployment, energy usage, protocols used and others are not perfectly suited to them. Like any other wireless networks such as ad-hoc networks their bandwidth is also narrow; as a result reliable data transmission was difficult. Due to these deficiencies of WSNs and the result in problem, making data transmission reliable becomes a sensitive issue.

In order to increase the reliability of the network (in our protocol), hop-by-hop data loss detection and recovery is employed. Unlike end-to-end loss detection and recovery which is used in normal internet protocol (like TCP and some others), in every hop the receiving node will check whether the data is successfully receiving. If the receiver detects some data is lost, the recovery process will take place between two adjacent nodes i.e. the sending and the receiving node. Furthermore the NACK-based approach with the last single ACK packet is employed in order to fully detect the data loss. The former approach is the best to detect any data loss but its only drawback is when the whole data is lost. In this particular case the receiver even couldn't know whether packet(s) were sent to it. To avoid this problem the last single ACK packet is used. All these ideas are implemented in NS2 simulator and some features are added to it (the simulator is extended).

The simulation results show that the designed protocol is 100% reliable. In addition to that by setting some protocol parameters the performance of the protocol under various conditions is tested and it performs well. For example, the worst case for reliability of a network with multi-source multi-sink wireless sensor network is when neighboring sources are sending and when two sources are sending in X-direction. However, our protocol is even fully reliable in the aforementioned worst cases.

5.2 Future Work

This research work was only considering uni-cast transmission that is when multiple sources for multiple/single sink(s). But in real world, for a single source node, there could be multiple destinations of outgoing packets (multicasting). Therefore, in the future multicasting in multi-source multi-sink WSNs can be included.

References

1. Yang, B.: Reliable Data Delivery in Wireless Sensor Networks, A thesis submitted to the college of graduate studies and research in partial fulfillment of the requirements for the degree of master of science in the department of Computer Science University of Saskatchewan (May 2010)
2. Kleu, C.: An Ultra-Low Duty Cycle Sleep Scheduling Protocol Stack for Wireless Sensor Networks. Submitted in partial fulfillment of the requirement for the degree Master of Computer Engineering in University of Pretoria (March 2012)

3. Ciciriello, P., Mottola, L., Picco, G.P.: textEfficient Routing from Multiple Sources to Multiple Sinks in Wireless Sensor Networks, Department of electronics and information and department of information and communication technology, Italy
4. Akyildiz, I.F., Kasimoglu, I.H.: Wireless sensor and actor networks: research challenges, School of Electrical and Computer Engineering. Georgia Institute of Technology, USA (May 2004)
5. Romer, K.: Distributed Mining of Spatio-temporal Event Patterns in Sensor Networks. Institute for Pervasive Computing ETH Zurich, Switzerland
6. Kim, S., Fonseca, R., Culler, D.: Reliable Transfer on Wireless Sensor Networks, Electrical Engineering and Computer Sciences, University of California at Berkeley, California
7. Felemban, E., Lee, C., Ekici, E., Boder, R., Vural, S.: Probabilistic QoS guarantee in reliability and timeliness domains in wireless sensor networks. In: Proc. IEEE INFOCOM 2005, Miami, FL, pp. 2646–2657 (March 2005)
8. Khayyat, A., Safwat, A.: Performance evaluation of distributed medium access schemes in multi-hop wireless sensor networks, Tech. Rep., Queen’s University, Kingston, Canada (2008), <http://www.ece.queensu.ca/directory/faculty/Safwat.html>
9. Rocha, F., Grilo, A., Pereira, P.R., Nunes, M.S., Casaca, A.: Performance evaluation of DTSN in wireless sensor networks. In: Cerdà-Alabern, L. (ed.) EuroNGI/EuroFGI 2008. LNCS, vol. 5122, pp. 1–9. Springer, Heidelberg (2008)
10. Wan, C.-Y., Campbell, A., Krishnamurthy, L.: Pump-Slowly, Fetch-Quickly (PSFQ): A Reliable Transport Protocol for Sensor Networks. In: WSN 2002, Atlanta, Georgia (September 2002)
11. Yaghmaee, M.H., Adjeroh, D.: A Reliable Transport Protocol for Wireless Sensor Networks. West Virginia University. In: International Symposium on Telecommunications (2008)
12. Zhang, H., Arora, A., Choi, Y.-R., Gouda, M.: Reliable Bursty Convergecast in Wireless Sensor Networks. In: Proc. ACM Mobihoc 2005, Urbana-Champaign, IL, May 25-28 (2005)
13. Deb, B., Bhatnagar, S., Nath, B.: ReInForM: Reliable Information Forwarding Using Multiple Paths in Sensor Networks, Dept. of Computer Science Rutgers University
14. Stann, F., Heidemann, J.: RMST: Reliable Data Transport in Sensor Networks (May 11, 2003)
15. Nazer, B., Gastpar, M.: Reliable Physical Layer Network Coding, arXiv: 1102.5724v1 [cs.IT], February 28 (2011)
16. Kim, S., Fonseca, R., Dutta, P., Tavakoli, A., Culler, D., Levis, P., Shenker, S., Stoica, I.: Flush: A Reliable Bulk Transport Protocol for Multi-hop Wireless Networks. Computer Science Division in University of California, and Computer Systems Lab in Stanford University
17. Lien, Y.-N.: Hop-by-Hop TCP for Sensor Networks. IJNC 1(1) (April 2009)
18. Issariyakul, T., Hossain, E.: Introduction to Network Simulator NS2, 1st edn & 2nd edn. Springer Science+Business Media, LLC (2012) doi: 10.1007/978-1-4614-1406-3
19. Mutter, T.: Partition-based Network Load Balanced Routing in Large Scale Multi-sink Wireless Sensor Networks, Department of Electrical Engineering, Mathematics and Computer Science, University of Twente (January 5, 2009)

Extraction of Fetal ECG from Abdominal ECG and Heart Rate Variability Analysis

Gizeaddis Lamesgin, Yonas Kassaw, and Dawit Assefa

Center of Biomedical Engineering, Division of Bioinformatics
Addis Ababa Institute of Technology, Addis Ababa University, Addis Ababa, Ethiopia
gizeaddis.lamesgin@ju.edu.et, yonas@bmyedtech.com, dawit.assefa@aaau.edu.et

Abstract. Non-invasive detection of fetal heart activity using abdominal electrocardiogram (AECG) is crucial for diagnosis as well as prognosis of heart defects. Most of the detection schemes proposed in the literature, however, require signals recorded either from both the mother's thoracic and abdominal regions or only the mother's abdomen but with multiple leads. The scheme proposed in the current study uses signals recorded from the mother's abdomen by a single ECG lead for efficient maternal QRS detection, fetal ECG extraction and enhancement, as well as fetal heart rate variability analysis by combining spectral analysis techniques with selected filters. The algorithm has been tested on 20 non-invasively recorded abdominal signals from the MIT/physioNet database and the results found were very promising.

Keywords: ECG, Heart Rate Variability, QRS, Signal Processing, S transform, Wavelets, Hilbert transform.

1 Introduction

Fetal heart defects are among the most common birth defects that may be unnoticed upfront that the baby appears healthy for many years after birth, or can be so severe that its life is in immediate danger. Compared to other procedures that are used to assess the fetal health such as cardiotocography (CTG), electrocardiography (ECG) provides more useful information about the fetal heart conditions such as the fetal heart rate (FHR) with a better predictive value. A typical ECG signal consists of three waveforms: P-wave (the arterial activation), QRS-complex (reflects the ventricular depolarization) and T-wave (refers to the ventricular repolarization) [1]. Fetus cardiac waveform helps physicians to diagnose fetal heart arrhythmia such as Bradycardia, Tachycardia, Congenital heart disease, and Hypoxia. Approximately the FHR ranges from 120bpm to 160bpm [1]. In the case of Bradycardia, the FHR reduces and goes below 120bpm and in Tachycardia the rate can go beyond 180bpm. Changes in the PR and PQ intervals, P-wave, T-wave and ST-segment, and the width of QRS-complex have been associated with oxygen levels. Lack of oxygen for long period can cause permanent damages to the brain and nervous system of the fetus, so early diagnosis helps physicians to have an effective and appropriate intervention.

There are two methods that are widely used in the literature for recording fetal ECG (FECG) signals: non-invasively by placing electrodes on the abdomen of the mother and invasively in which case the electrodes are placed inside the uterus of the mother on the scalp of the fetus during labor. Invasive extraction of FECG is more accurate because of the recording electrode placed on the fetus scalp but has its drawbacks in that it can only be done during delivery. The non-invasive method can be used in all gestation weeks and also during delivery. Eventhough non-invasive abdominal recorded FECG provides significant clinical information about the health condition of the fetus, the signal is often contaminated with large amount of noise making it difficult to accurately detect and extract the FECG. The main sources of noise include the maternal electromyogram (EMG) with a wide frequency range, power line interference (PLI) (50/60Hz frequency), the slow baseline wander (BW), random electronic noise, and the maternal ECG (MECG). In addition, use of several leads for the purpose of detection might make the task too cumbersome. Also, due to the overlap in time and frequency with the MECG and other noise signals, the signal to noise ratio (SNR) of the FECG is relatively low and has been a classical problem in biomedical signal processing. In this regard, several methods have been proposed in the literature for extracting the FECG from composite abdominal signal and the rest of the noise.

In the current study a novel FECG extraction scheme is proposed based on an efficient AECG preprocessing stage followed by a robust FECG denoising technique. Joint time-frequency localized transforms are used for signal visualization and enhancement. The proposed method is applied for extraction of FECG signals from single lead abdominal recorded ECGs with a maximum SNR suitable for effective heart rate variability analysis of the fetus.

2 FECG Detection

In a non-invasive recording method, the abdominal signal comes with the low amplitude FECG and several other signals of high amplitude including the MECG and other noise signals. The large amplitudes of these noise sources hide the transabdominal FECG and a simple high-pass filtering of abdominal signals for FECG extraction cannot be applied due to the over-lapping spectra of the FECG and that of the noise components. Also, an application of a filter may introduce some unwanted phase distortions on the FECG. Moreover, the amplitude of the FECG depends on the electrode configuration and varies among subjects due to body weight and size of the mother as well as positions of the fetus. Thus, it is desirable to eliminate as much noise as possible during recording in order to apply algorithmic analysis for further cleaning of the FECG signal.

Numerous attempts have been made in the literature to detect FECG signals from AECG. For example, an approach to extract the FECG from two ECG signals recorded at the thoracic and abdominal areas of the mother's skin with the help of a hybrid soft computing technique called Adaptive Neuro-Fuzzy Inference System (ANFIS) has been proposed [2]. Another attempt added an

equalizer to the ANFIS to enhance the extraction and increase the SNR [3]. A method that combines independent component analysis (ICA) and projective filtering has also been proposed [4]. The Polynomial Networks technique has also been exploited [5] and the Wavelet transform has been added to this method as a post processing tool to de-noise the extracted FECG [6]. A new mother wavelet called Abdominal Electrocardiogram Mother Wavelet for FECG extraction was suggested to achieve optimal denoising [7]. Another method combines ICA with Wavelets for a more effective de-trending and de-noising [8]. Sequential analysis approach [9] and an artifact reduction procedure in AECG recordings for FECG estimation based on dynamic segmentation of the MECG and subsequent linear prediction of the MECG segments have also been proposed [10].

The above mentioned methods for extracting FECG have, however, their own limitations. Some need priori information for extraction and others need signals recorded from many leads and many of them require recording from the mother's thoracic area. Those methods require electrodes to be scattered on the mother's body and thoracic area making their application inconvenient in non-clinical environments (home care devices, for example). It is therefore desirable to develop a technique that permits extraction of FECG with a minimal abdominal lead setup and independent of any particular recording technique. A single-lead configuration is advantageous for implementation in a mini-apparatus, making it suitable for ambulatory and long-term monitoring. In the current work, we proposed a method of FECG extraction from single lead abdominal signals by applying time-frequency localized transforms and algorithms for filtering, signal enhancement, QRS detection, and fetal heart rate variability (FHRV) analysis.

3 Heart Rate Variability Analysis

Heart Rate Variability (HRV) has long been used as a screening tool for diagnostic purposes. HRV is a physiological phenomenon that reveals the state of having non-consistent R-R durations/intervals over a number of cardiac cycles per unit time and it provides a numeric description of HR fluctuations around a baseline that represents an average HR of a subject. It is established that HRV describes the ability of the heart to adapt and respond to changing circumstances of different types and cases of stimulations and is predominantly dependent on the extrinsic regulation of HR. The aim of HRV analysis is to evaluate the state and conditions of the automatic nervous system (ANS) that is mainly in charge of cardiac activity regulation. A wide variety of HRV analysis tools exist such as statistical methods, geometric methods, and power spectral density analysis.

4 Materials and Methods

The proposed scheme in the current work is comprised of four main steps. First step is signal preprocessing, aiming to remove the BW, PLI, and other noise components. In the second step Wavelets, Hilbert transform and adaptive thresholding based peak detection and filtering is applied to detect the MECG signal.

Third step is composed of extraction and enhancement of the FECG signal and finally FHRV analysis. The AECG signals used for testing were taken from the MIT/PhysioNet non-invasive FECG database. The signals were taken between 38 to 40 weeks of pregnancy from five subjects (RO1, RO4, RO7, RO8, and RO10 as seen in Table (2)) and the duration of each signal is 60sec. The signals were digitized at 1000Hz with 16bit resolution.

4.1 Time-Frequency Localized Transforms

Signal analysis in the time domain is mostly difficult due to many reasons such as the complexity of the signal, the noisy nature and the like. Instead, a frequency domain analysis is commonly used to identify the most distinguished information in the signal. For a stationary signal, where the frequency structure does not change with time, such an analysis may adequately be performed, for example, through the Fourier transform (FT). However, most of our medical signals, including the ECG, are non-stationary with time dependent spectral content. The advent of joint time-frequency analysis techniques have made it possible to efficiently analyze the later signal types. These include the Gabor or Short Time Fourier Transform (STFT), Wigner-Ville distributions, and the Wavelet transform. The fixed resolution of the STFT, the cross line interference in Wigner-Ville are some of the drawbacks. The Wavelet transform was introduced to circumvent those issues. The continuous Wavelet transform of a signal $h(t)$ is given by:

$$W(a, b) = \int_{-\infty}^{\infty} h(t)\psi^*\left(\frac{t-b}{a}\right)/\sqrt{a}dt . \quad (1)$$

where a is the dilation and b is the translation parameters, ψ is the mother Wavelet (Mexican Hat, Gaussian, Daubechies, etc) and $*$ denotes complex conjugation. $h(t)$ can be reconstructed from the Wavelet transform if ψ satisfies the admissibility criteria given by:

$$C = \int_{-\infty}^{\infty} |\psi(w)|^2/|w|dw < \infty . \quad (2)$$

The S transform [11] was introduced to address some of the drawbacks of Wavelets and is given by:

$$S(\tau, f) = \int_{-\infty}^{\infty} h(t)|f|/\sqrt{2\pi}e^{-\frac{(\tau-t)^2 f^2}{2}}e^{-i2\pi f t} dt . \quad (3)$$

The resolution of the S transform is explicitly a function of frequency as opposed to the arbitrary delation parameter used in Wavelets. The S transform not only estimates the local power spectrum but also the local phase spectrum which is absent in wavelets. The localizing window in the S transform is a Gaussian chosen because it uniquely minimizes the quadratic time-frequency moment about a time-frequency point and is symmetric in time and frequency with no side lobes. The S transform is invertible and is directly related to the FT which allows its easy computation using the Fast FT (FFT) through convolutions.

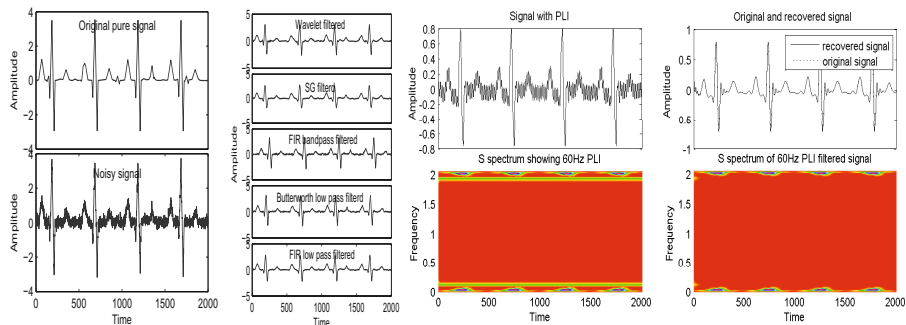


Fig. 1. Column: (1) Original ECG signal (top) and the ECG signal with noise applied (bottom), (2) Performance of different filters applied to remove the noise; (3) An ECG signal with PLI applied (top), its amplitude S spectrum (bottom), (4) The original (with no PLI) and filtered signals plotted together (top), the PLI filtered amplitude S spectrum (bottom) showing the effectiveness of the notch filter.

In this study we make use of both Wavelets and the S transform to analyze our signals and utilize the ability of the Hilbert transform to identify zero-cross points in ECG signals as dominant peaks and apply it in QRS detection.

4.2 Signal Generation and Preprocessing

As a testing stage, a simulated ECG for both the mother and fetus was generated (MATLAB). ECG with a peak voltage of 3.5mV and HR of 89 beats/minute was assumed for the maternal signal. The heart of a fetus beats noticeably faster than that of its mother, with rates ranging from 120 to 160 beats/minute, and the amplitude of the FECG is much weaker than that of the MECG. Accordingly an ECG signal corresponding to HR of 139 beats/minute and a peak voltage of 0.25mV was assumed for the fetus. The measured AECG signal is usually dominated by the maternal heartbeat signal that propagates from the chest cavity to the abdomen. In addition, an uncorrelated Gaussian noise was added to simulate any broadband noises within the measurement. Then filtering was applied to the resulting noisy signal and the recovered signals from different filters were compared to the original signal to evaluate the filters' efficiency.

General framework of morphological ECG preprocessing has two main parts: baseline correction and suppression of noise and the PLI. In abdominal signals the MECG is considered as noise in FECG analysis and vice versa. So appropriate denoising algorithms are required to emphasize one and suppress the other. In the subsections below, we discussed different artifact removal filters and their comparative analysis.

Comparative Analysis of Morphological Noise Filters: Several filters have been tested and compared in this work for their efficacy in denoising signals. These include: window based finite impulse response (FIR) filter (which is often considered superior to both analogue and infinite impulse response (IIR) filters), Savitzky-Golay smoothing, Butterworth filter, and Wavelet transform approach. The mean, SNR, and accuracy were used as parameters for comparison. The mean of the original signal is compared with that of the recovered signal and the error is computed from the difference of the two means. Our results showed that FIR low pass filter, Butterworth filter, and Wavelet based filter have good performance in terms of the mean criteria. Butterworth and Wavelet based filters resulted in higher SNR and showed good accuracy. The performance of different filters applied on an ECG signal with noise is demonstrated in Fig. (1).

Comparative Analysis of BW Removing Methods and PLI Filters: BW is a common phenomenon in biomedical electrical recordings. We compared different BW removing techniques suggested in the literature including: Savitzky-Golay smoothing, Wavelet Approach, Zero phase filtering [12], and Polynomial fitting (parabolic filter) combined with Zero phase filtering. The SNR was used as parameter for comparison. Accordingly, our analysis showed that Savitzky-Golay smoothing approach offered higher SNR.

The PLI, induced by the AC electrical power source, consists of a 50/60Hz artifact together with its harmonics. The typical amplitude of the PLI can be as high as 50% of the peak-to-peak amplitude of the ECG signal. Notch filter is used to remove the PLI due to its ability to pick a specific frequency and attenuate and eliminate it from the input spectrum while leaving the amplitude of the rest of the frequencies relatively unchanged. In our work a 60Hz PLI was added to the synthetic ECG and filtered with IIR notch filter. Last two columns in Fig. (1) present the S amplitude spectrum of an ECG signal with the modeled PLI applied, where the PLI and the rest of the signal are clearly seen and the notch filter applied to remove the PLI performed quit satisfactorily.

4.3 The Proposed Method

Based on the above comparative analysis, we applied Butterworth low pass filter and Wavelet denoising to reduce random noise, Savitzky-golay filter to remove the BW, and Notch filter to extract PLI on the real signals acquired from the MIT/physioNet database as a pre-processing stage and further analysis were carried out. In all tasks the S transform was used for spectral visualization.

Maternal QRS Detection: The QRS complex is the most prominent waveform within the ECG signal with normal duration from 0.06s to 0.1s [13]. The QRS complex consists of three characteristic points within one cardiac cycle denoted as Q, R and S. It reflects the electrical activity within the heart during ventricular activation. Its shape, duration and time of occurrence provide valuable information about the current state of the heart. Because of its specific shape, the QRS complex serves as an entry point for almost all automated ECG analysis algorithms. QRS detection is not a simple task, however, due to

the varying morphologies of normal and abnormal complexes and that the ECG signal experiences different types of disturbances with complex origin. The possibility of maternal QRS (MQRS) and fetal QRS (FQRS) overlaps both in time and frequency makes the detection even more challenging.

Different QRS detection methods are available in the literature and most of the detectors are comprised of two stages: a preprocessor stage to emphasize the QRS complex and a decision stage to threshold the QRS enhanced signal. Some of the algorithms include: derivative-based [14], template matching and morphologic filtering, gene-based design, STFT based [15], and Wavelet based algorithms [16]. A combination of Wavelets, Hilbert transform and adaptive thresholding has also been proposed [17]. Another method operates by feature extraction, event detection, and localization of R peak by counting zero crossings [18].

The methods of QRS detection mentioned above may not be effective for abdominal ECG FQRS/MQRS detection. It is because the FQRS may overlap in time and frequency with the MQRS complex and consequently not possible to separate them using conventional frequency selective filtering or simple thresholding. A method of adaptive thresholding and a way to emphasize the MQRSs is required together with medical knowledge of fetal/maternal heartbeats.

The new algorithm proposed in this study for FQRS detection is comprised of bandpass filtering, Hilbert transform and adaptive thresholding. A bandpass filter is first applied to the signals for noise attenuation while preserving the essential spectral content of the MQRS complex. The frequency components of the QRS complex are different for the mother and the fetus. Therefore the relation between the frequency contents has to be considered. This can be done by considering the QRS durations for both adult and fetus signals and the fact that the QRS frequencies are directly proportional to the QRS durations. The QRS duration cannot exceed 120ms for the adult and 80ms for the fetus [19]. This gives the relation for the QRS duration:

$$QRS_{fetal}/QRS_{adult} = 80ms/120ms = 0.67 . \quad (4)$$

The frequency ratio is therefore $\frac{1}{0.67} = 1.5$. This means that the bandpass filter should be designed to have a frequency range of 18-35Hz for the adult and 27-53Hz for the fetus. So applying a bandpass filter with cutoff frequencies of 18 and 35Hz will emphasize the MQRS complex while reducing the FQRS complexes. The Hilbert transform stage is useful for signal demodulation without knowing the carrier frequency. If we consider a QRS complex as a modulated waveform, the beginning and end of the QRS complex envelope calculated using the Hilbert transform coincide with the QRS onset and offset respectively [20]. In Hilbert transform zero-cross points in ECG signal are formed as dominant peaks in output. Using discrete-time notation, a complex sequence $ECG(k) + iECG_H(k)$ is formed from the real signal sequence $ECG(k)$, where $ECG_H(k)$ is the Hilbert transform of $ECG(k)$. Then, the envelope signal is defined as $ECG_e(k) = \sqrt{ECG^2(k) + ECG_H^2(k)}$. Hilbert transform of the original function represents its harmonic conjugate and the QRS complex corresponds to

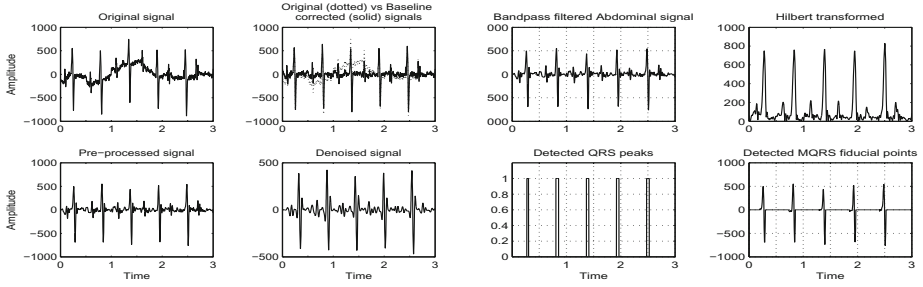


Fig. 2. Pre-processed AECG signal, its Hilbert transform, maternal R-peaks and QRSs

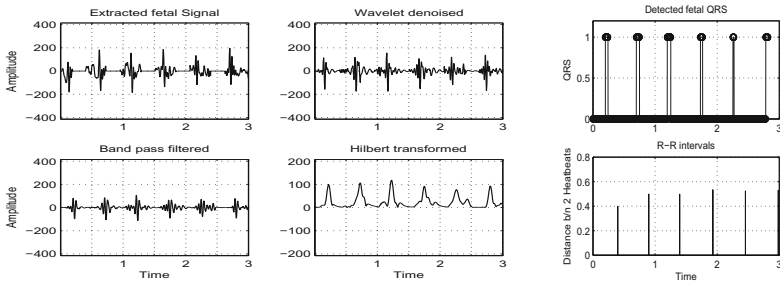


Fig. 3. Extracted and enhanced fetal signal, its Hilbert transform, detected fetal R-peaks, and HRV results

a hump of the envelope signal. Signals with negative amplitude (lead inversion) are seen as positive peaks in the Hilbert transform envelope.

Two important medical rules/truths are applied during the Adaptive thresholding stage: two R-peaks closer than 0.2sec cannot exist in the ECG, and a time interval greater than 2sec cannot exist in the ECG without an R-peak.

Maternal QRS Fiducial Point Detection: In order to remove the MQRS signals it is necessary to detect the QRS onset and offset points. Two intervals of time (search windows) in each cardiac cycle, containing the onset and offset of the QRS respectively, are defined: $[R_p - 70, R_p]$ for QRS onset, and $[R_p, R_p + 80]$ for QRS offset, where R_p is the detected R peak at each cardiac cycle, and the optimal values 70ms and 80ms were selected experimentally. In addition, the MQRS interval is approximately 120 and the long refractory period to the left of the R-peak is also considered. To remove the p-waves with the MQRS the search window was made larger than the maximum QRS duration (120ms).

Fetal ECG Extraction and Enhancement: The MQRS signal is captured within a window as defined above. All samples that fall within this window will be zero padded in the FECG extraction stage. This is done to eliminate or scale down the maternal residual peaks from the abdominal signal. After removing the MQRS signal the remaining signals will be a noise signal, maternal residual

signal and the fetal signal. Wavelet denoising is performed to remove the lower frequency signal components.

Fetal QRS Peak Detection: As discussed above, the FQRS points most dominate in the frequency range of 27-53Hz. So bandpass filtering within this range emphasizes the FQRS signal and suppresses the noise signals. After the bandpass filter, the Hilbert transform is applied so that at the QRS points positive peaked signals will result. The next stage is comprised of an adaptive thresholding and search back algorithm together with the two medical rules discussed earlier. If a peak is not detected within 2sec, the time instant corresponding to the nearest MQRS point is considered as a candidate fetal R-peak assuming that the MQRS is overlapped with the FQRS and rejected during MQRS window removal.

Table 1. Numerical analysis of HRV. Occur time is in seconds.

	Max. HR		Min. HR				Max. HR		Min. HR		
ECG	Value (bpm)	Occur Time	Value (bpm)	Occur Time	Mean HR	ECG	Value (bpm)	Occur Time	Value (bpm)	Occur Time	Mean HR
1	157.48	0.38	119.28	2.45	138.38	6	163.49	1.80	46.19	3.17	104.84
2	150.78	0.40	118.58	2.98	134.68	7	110.70	0.54	48.31	1.81	79.51
3	148.15	0.41	122.95	2.98	135.55	8	155.44	1.40	63.56	0.94	109.50
4	151.13	0.40	123.71	2.98	137.42	9	94.19	0.64	74.63	2.38	84.41
5	158.31	0.94	120.24	0.50	139.28	10	150.00	0.97	67.42	2.88	108.71

FHRV Analysis: The proposed FHRV algorithm in this paper is explained as follows. The intervals between successive normal complexes are determined first. If the peak detected is the first peak, the R-R interval is made to zero to remove first peak duration count and for the rest peak detection points the distance between the two peaks will be equal to the new index of peak detection point minus the old index of peak detection point measured in time until all the peak detection is over. The normal-to-normal (R-R) intervals (i.e. all intervals between adjacent QRS complexes resulting from sinus node depolarizations) or the instantaneous HR is then determined. The analysis also allows to determine the time of ANS intervention and other useful statistical information. Simple time domain variables that are calculated include the mean RR interval, the mean HR, and the instantaneous low and high HRs.

5 Results and Discussion

Figure (2) presents an abdominal signal, preprocessed and Hilbert transformed signals, maternal R-peaks, and detected MQRS onset and offset points. All maternal fiducial signals are correctly identified which was a vital step for fetal signal extraction. Amplitude spectrum of the S transform (not shown) was used for local time-frequency spectral visualization. Figure (3) presents the extracted

and enhanced fetal signal, its Hilbert transform, detected fetal R-peaks, and HRV results. The last column shows plot of the lengths of each successive intervals (in seconds) against time. The longest and shortest intervals show the time instant of dropping and rising of instantaneous HRs respectively. It also indicates the time of the ANS intervention to decrease or increase the HR so as to put it to normal condition. So it is possible to track the ANS activity and the wellbeing of the fetus in general. Table (1) summarizes the highest, lowest, and average HRs and their occurring times.

Sensitivity (SE), specificity and accuracy were calculated in order to evaluate the performance of the proposed method in R-wave detection of the FECG signals. Table (2) lists accuracy measures of the proposed scheme on the real ECG data acquired from MIT/physioNet database. Our proposed method can obtain reasonable estimates of FECG signals for approximately 88% of the abdominal signals in this database, which seems to be a promising result. The 98.38% positive prediction rate shown in the table implies that the algorithm detects almost all existing peaks and the missed peak detection is resulted because of nonexisting peaks in the original abdominal signal. This may occur due to misplacement of leads on the mother's abdomen, the fetus position, or similar other factors.

Table 2. Performance analysis results. SE, +P, and DER are in %.

ECG	Lead	TP	FP	FN	SE	+P	DER	ECG	Lead	TP	FP	FN	SE	+P	DER
R01	1	7	0	0	100	100	0	R08	1	8	0	3	72.72	100	27.28
	2	5	0	2	71.42	100	28.5		2	9	0	2	81.81	100	18.18
	3	5	1	2	71.42	80	28.5		3	9	0	2	81.81	100	18.18
	4	6	1	1	85.71	87.71	28.57		4	9	0	2	81.81	100	18.18
R04	1	7	0	0	100	100	0	R10	1	10	0	0	100	100	0
	2	5	0	2	71.42	100	28.5		2	10	0	0	100	100	0
	3	5	0	2	71.42	100	28.5		3	10	0	0	100	100	0
	4	7	0	0	100	100	0		4	10	0	0	100	100	0
R07	1	5	0	2	71.42	100	28.5	Total Average (%)				88.05	98.39	12.65	
	2	7	0	0	100	100	0								
	3	7	0	0	100	100	0								
	4	7	0	0	100	100	0								

6 Conclusion and Feature Work

Simple, automatic, and effective schemes are developed in this study for use in signal processing of ECG waveforms generated from single lead non-invasive abdominal recordings. The single lead configuration allows FECG monitoring in a non-clinical environment. Minimal abdominal detecting electrodes are used, and the method is simple to operate even by the mothers themselves, and hence can be used in a normal home environment. The hybrid S transform, Wavelets, and Hilbert transform based approach combined with selected filtering schemes have resulted in effective spectral visualization, enhancement, and detection of

useful ECG signals particularly FECG and maternal and fetal QRSs as well as HRV analysis. Experiments with synthetic signals and real FECG data from a database have very well demonstrated the feasibility of the proposed scheme.

Applications of the proposed method on different data sets may be useful for further validation and also to investigate potential clinical implications. Moreover, the method proposed here can be applied only for single fetus signals recorded from the mother's abdomen. To detect twins' signals, the algorithm should be modified in some way. So, our future work will also focus on the development of such modified algorithms to extract twins' fetal signals.

References

1. Chazal, P., de, O.M., Reilly, R.B.: Automatic classification of heartbeats using ECG morphology & heartbeat interval features. *IEEE Trans. Biomed. Eng.* 51(7), 1196–1206 (2004)
2. Saranya, S., Suja, S.P.: A Novel Hybrid Soft Computing Technique for Extracting Fetal ECG from Maternal ECG Signal. *Int. J. Comput. Applicat.* 3(8) (2010)
3. Kulkarni, S.S., Lokhande, S.D.: Detection of Fetal ECG from Abdominal ECG recordings using ANFIS & Equalizer. *IJERT* 2(6) (2013)
4. Kotas, M.: Combined Application of Independent Component Analysis & Projective Filtering to Fetal ECG Extraction. *Biocybernetics & Biomed. Eng.* 28(1), 75–93 (2008)
5. Assaleh, K., Al-Nashash, H.: A novel technique for the extraction of fetal ECG using polynomial networks. *IEEE Trans. Biomed. Eng.* 52(6), 1148–1152 (2005)
6. Ahmadi, M.: Fetal ECG Signal Enhancement. MSc Dissertation, American University of Sharjah, Sharjah, UAE (2008)
7. Suarez, F., Gransky, I.: A new mother wavelet for Fetal ECG to achieve optimal denoising & compressing results. *Int. J. Acad. Res. Applied Sci.* 2(2), 1–39 (2013)
8. Zhou, Z., Yang, K.: Fetal Electrocardiogram Extraction & Performance Analysis. *J. Comput.* 7(11) (2012)
9. Martens, S.M.M., Rabotti, C., Mischi, M., Sluijter, R.J.: A robust fetal ECG detection method for abdominal recordings. *Physiol. Meas.* 28, 373–388 (2007)
10. Vullings, R., Peters, C., Mischi, M., Sluijter, R., Oei, G., Bergmans, J.: Artifact reduction in maternal abdominal ECG recordings for fetal ECG estimation. In: 29th Int. Conf. IEEE EMBS, Lyon, France, August 23–26 (2007)
11. Stockwell, R.G., Mansinha, L., Lowe, R.: Localization of the Complex Spectrum: The S Transform. *IEEE Trans. Sig. Process.* 44(4) (1996)
12. Sonali, P.: Patial, Different Techniques of Baseline Wandering Removal - A Review. *Int. J. Enhanced Research Sci. Techno. & Eng.* 2(5), 37–43 (2013)
13. Johneff, V.: Complex Valued Wavelet Analysis for QRS Detection in ECG signals. Technical Uni. of Sofia
14. Pan, J., Tompkins, W.J.: A Real-Time QRS Detection Algorithm. *IEEE Trans. Biomed. Eng.* 32(3), 230–236 (1985)
15. Uchaipichat, N., Inban, S.: QRS detection using the short time fourier transform. *IJCA Special Issue on Comput. Aided Soft Computing Techniques for Imaging and Biomed. Applica.* (2012)
16. Vagadiya, B.D., Dasi, S., Bhatt, V., Doshi, K.: QRS detection using the Wavelete transform. *IJERT* 2(4) (2013)

17. Rabbani, H., Mahjoob, M.P., Farahabadi, E., Farahabadi, A.: R Peak detection in Electrocardiogram Signal based on an optimal Combination of Wavelet Transform, Hilbert Transform & Adaptive Thresholding. *J. Med. Signals Sens.* 1(2) (2011)
18. Kohler, B.-U., Hennig, C., Orglmeister, R.: QRS detection by zero-crossing counts. *Progress in Biomed. Res.* 8(3) (2003)
19. Sigurdardottir, I.E.: R-wave detection algorithms using adult & fetal ECG signals. MSc Dissertation, Chalmers Uni. of Techno., Gothenburg, Sweden (2013)
20. Manriquez, A.I., Zhang, Q.: An Algorithm for Robust Detection of QRS Onset & Offset in ECG Signals. *Comput. in Cardiology* 35, 857–860 (2008)

Addition of Static Aspects to the Intuitive Mapping of UML Activity Diagram to CPN

Jan Czopik, Michael Alexander Košinár, Jakub Štolfa, and Svatopluk Štolfa

VSB - Technical University of Ostrava,
17. listopadu 15, 708 33 Ostrava,
Czech Republic
{jan.czopik.st,michal.kosinar,jakub.stolfa,
svatopluk.stolfa}@vsb.cz

Abstract. Software process is a core of every software company and even in cases the processes are not documented they are still there. The cases when the process must be documented raise the need of powerful process framework, methodology, and tools that are able to catch every process aspect. Our current research focuses on creation of formal software process framework that combines mathematically precise approaches like OWL and Petri Nets with semi-formal techniques based on UML that make the framework easy to understand and use. Our recent results discussed the tools covering the modeling of dynamic process aspects with Colored Petri Nets transformed from UML activity diagrams. This paper presents the results by incorporating static aspects in the process of SP modeling and its simulations including resources, artifacts, and workers.

Keywords: software process, formal methods, OWL, UML, formalization, Activity Diagram, CPN.

1 Introduction

Traditional approach to process modeling is to model the process based on a structure, behavior and functional description of a company; i.e. modeling from centralized perspective collecting information about activities, conditions, control paths, objects, roles and artifacts along with internal rules and constraints.

Traditional specifications approaches work well for simple business processes; however it is not sufficient for business processes that could suffer from big complexity or changing environment where the main requirement is to have precise, stable and reusable models easily used in simulation methods [1, 7, 12].

With rising need of adequate adaptability and flexibility in business processes many research groups have been investigating various methodologies, adaptive process techniques, etc. [10, 17]. Another solutions involve utilization of formal systems, knowledge bases and mathematical models (e.g., Discrete-Event Simulation, System Dynamics, etc.). Implementation of formal rules and facts, the process model and its execution may be more adaptive to unexpected situation and events.

On the contrary conventional semi-formal approaches may lead to misunderstandings and errors and are not easy to use in process simulations without further processing due to their lack of formal power and formal systems are not easy to use in practice because of their complexity.

This paper provides an extension to our previous effort introduced in [5], in terms of simulating consumption of resources, generation of artifacts and utilization of workers in the process thus in conjunction with the intuitive mapping framework providing a complex framework for simulation, analysis and verification of software processes.

The paper is organized as follows: Section 2 covers the state of the art in the field of our research; Section 3 introduces the basic theories and proposal of the methodology with a simple software process example hierarchy and analysis including the software tool utilizing the approach; Sections 4 and 5 conclude the benefits of the approach and a summary of the work and outlines future work.

2 State of the Art

Nowadays there are many modeling techniques for process modeling as is mentioned in Vergidis's paper [20]. On the other hand, the software process is quite specific [18] and it has been characterized as "the most complex endeavor humankind has ever attempted" [4].

By converting software process models modeled in UML2 Activity diagrams to their Colored Petri Nets (CPN) counterparts we get advantages of well-formed formal language. Some of these advantages might include: *verification of the model*, *performance evaluation*, *state space analysis* and *simulation* [8].

For modeling, simulation and analysis we can use CPN tools that provide means for syntax checking, efficient simulation of either timed or untimed nets and their analysis.

Other researchers had already proposed methods for UML formalization to CPN. Some of them focused on static part of the software (business) process formalizing Class diagrams and Collaboration diagrams like in the paper of Du [6], others focused on the dynamic part by formalizing Activity diagrams [9, 14, 19] and Interaction overview diagrams [16]. We are focused on formalizing dynamic part of the software process (UML activity diagrams). Our work extends previous effort of researches mentioned above to eliminate any loopholes and to make the formalization as simple as possible. The problem we found in work of Staines [19] and Jung [9] was lack of precise definition of the colorset declaration, variable declaration and arc expression/guard definition. Our mapping framework provides *complete* set of conversion rules with clear arc expression definitions, guard definitions and explanation of colorset declarations. The strongpoint of our framework is its simplicity with maintaining completeness and simple extensibility. Because we are using *focus token* approach to pass the token of execution between activities, it is easy to extend the rules to accompany advanced structures like roles, actors and resources by using OWL for modeling hierarchy and relationships between those entities. OWL has been successfully used for modeling of static aspects of the process in [3, 11, 13, 17]. The extension of additional colorsets will be presented in this paper. The conversion can be still completely automated without need of any user action no matter of the newly added colorsets.

3 Proposed Approach of Extended Intuitive Mapping from Semi-formal to Formal Model

Modeling with formalization framework doesn't differ that much from normal modeling with UML Activity diagrams. The main difference is additional post-modeling step consisting of automated translation of the Activity diagram to CPN.

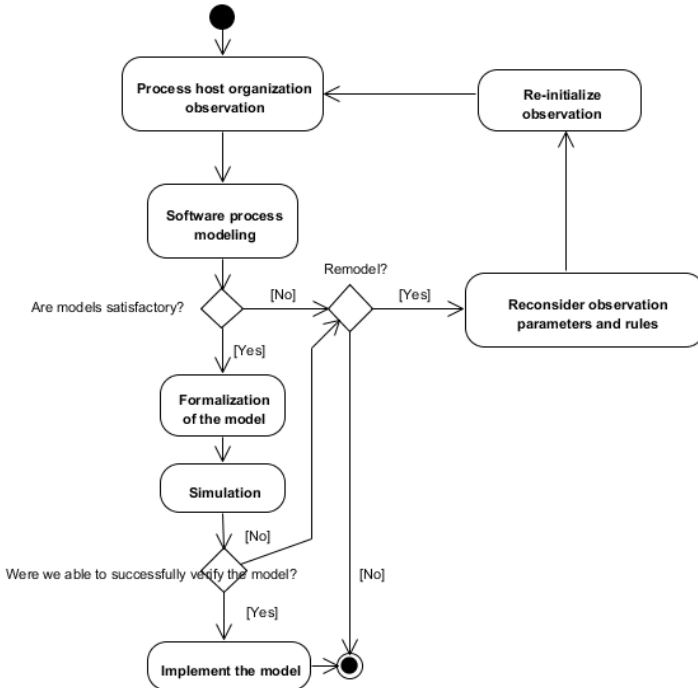


Fig. 1. Software Process formalization framework application

The formalization activity in Fig. 1 is defined by framework that consists of easy to read and understandable intuitive mapping rules.

At first, we define code generation of *colorsets*, variables and values. As a core, we will use *colorset* of type *bool* that will be used to pass focus from one activity to another. Next, we will examine all *decision* nodes to find out number of conditions used. Based on that step we will declare different enumerated *colorset* for each number of conditions used by decision node. If decision node has two conditions, we will generate enumerated *colorset* with two values, if decision node has three conditions, we will generate enumerated *colorset* with three values etc.

Let's consider *resources*, *artifacts* and *workers*. We can look at this three "categories" as if they were superclasses or stereotypes which can be further subclassed or refined. What it means exactly can be seen on Fig. 2. Each activity in the process requires a worker which will perform this activity in a certain role.

Activity (respectively worker performing this activity) may or may not produce an artifact as an output of the activity [15]. In the paper [2], Aalst is speaking of resource as an executor of an activity and by his definition, resource can be either a *machine* or a *person*. We would like to separate resource category into two, not differentiating between people and machines, but separating them by their activeness. Making workers *active* participants and resources *passive* participants of the process. Workers will execute activities utilizing/consuming resources and creating artifacts. The hierarchy will be modeled in OWL and the OWL model will be used in the formalization process (the additional colorsets will be generated from OWL model).

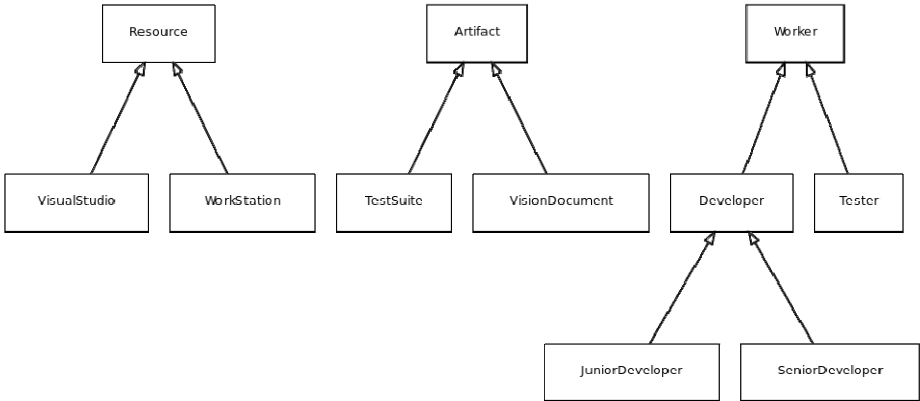


Fig. 2. Example of roles, artifacts and workers hierarchy

The subclasses of *resource* class and subclasses of *artifact* class are considered as specializations (concrete instances) of their superclasses. Whether the subclasses of *worker* class are looked upon as *roles*, more than concrete instances. Three colorsets will be generated: Resources, Artifacts, Workers. All will be enumeration colorsets. The concrete named identifiers will be generated based on class diagram associated to the activity diagram. The associated class diagram contains the hierarchy (in form of a graph) of the resources, artifacts and workers. Names of the classes will be simplified to their alphanumeric counterparts to match the constraints of CPN ML language. The generation of colorsets mentioned above counts on uniqueness of different resources, artifacts and workers.

As for the variables, we declare one variable d (as decide) for each *Conditions colorset* and it will be used for conditional routing. The simple d variable with no suffix is used for control flow conditional expressions and for *precondition* and *postcondition* of an activity. The last variable called b will be used in exception handling. No variables are declared for resources, artifacts and workers since we will generate the arc expression directly based on the concrete resource, artifact or worker (role) value needed. We are also going to declare two value constants T and F each representing multiset of one *true* (*false* respectively) token.

In terms of CPN ML language it looks like the example on the Fig. 3 (roles, artifacts and workers were taken from Fig. 2 as an example).

```
colset Focus = bool;
colset Conditions2 = with First | Second;
colset Conditions3 = with First | Second | Third;
colset Resources = with Workstation | VisualStudio;
colset Artifacts = with VisionDocument | TestSuite;
    colset workers = with Developer | JuniorDeveloper | SeniorDeveloper
                    | Tester;

var b : Focus;
var d : Focus;
var d2 : Conditions2;
var d3 : Conditions3;
val T = 1`true;
val F = 1`false;
```

Fig. 3. CPN ML declarations of colorsets, variables and constants used for translation

3.1 Extended Mapping Rules Used for Formalization

Fig. 4 shows the transformation rule for the transformation of an activity extended by additional information to corresponding CPN elements.

The resulting net is extended by three special places which represent buffers/accumulators for *resources*, *artifacts* and workers under *roles*. The *resource pool place* can have both incoming and outgoing arcs since the activity is borrowing the resource for limited amount of time and then returning it making it available for other activities. In case of resource, that can be used only once (does not return to resource pool after used) the connection from *resource pool place* to an activity is only one-directional (resource is lost after execution of the activity). The *artifact pool place* can have only incoming arcs since we are only creating artifacts in activities. If the created artifact is required as resource later in the process a special arc (dotted arc on Fig. 4) is drawn from the activity creating the artifact to *resource pool place* making it available as resource (but the artifact has to inherit also from the Resource superclass, not just the Artifact superclass). The *role pool place* has both incoming and outgoing arcs thus simulating utilization of workers (including conflicts). All three of these special places are singleton places meaning every activity extended by resources, artifacts or roles is connected to the same place.

An incoming arc is added to the first transition of an activity for either resource and/or worker if needed (if present in meta-information). The arc inscription is generated based on the given meta-information (object flow, swimlanes).

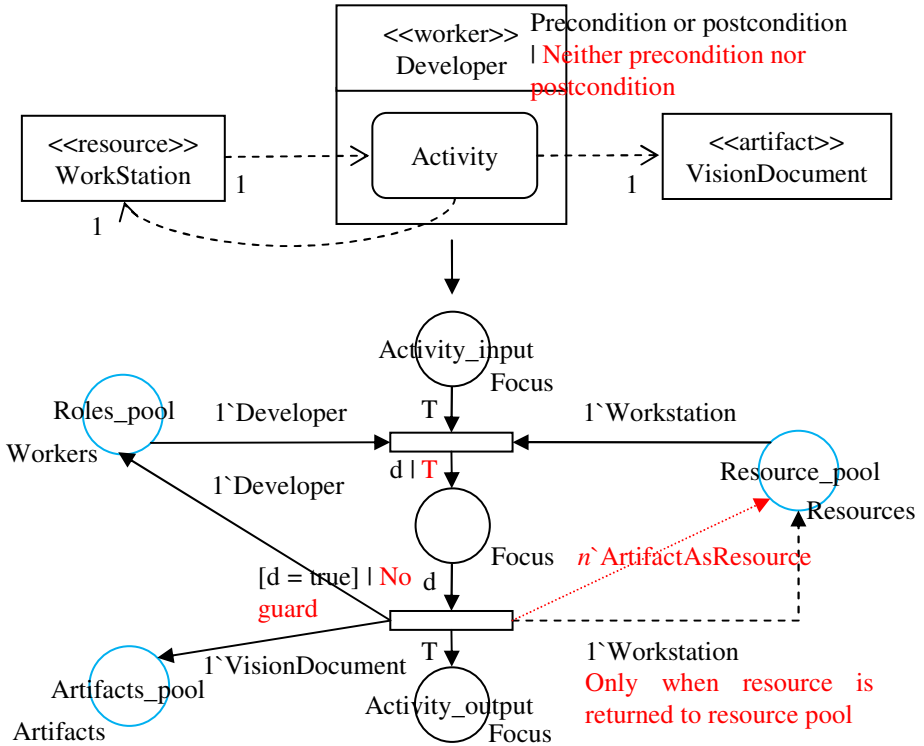


Fig. 4. Extended mapping rule of an Activity

For the *resources* the object associated with the activity with given multiplicity provides information about what type of resource is needed and how many instances of that resource are needed. We can easily convert the object association to valid CPN ML expression (see Fig. 4).

If the activity is in a swimlane, the swimlane (role) is used for generation of arc inscription of the incoming arc. But unlike in the case of *resources*, we consider the *workers* hierarchy to represent *roles* in the process. And because of that, we might want to restrict execution of an activity to either very specific role (for example *SeniorDeveloper*, see Fig. 2) or to allow group of roles to execute the activity (for example *Developer* role, in this case doesn't matter if the activity will be executed by *JuniorDeveloper* or *SeniorDeveloper*, both of them have the rights to execute the activity, see Fig. 2). Bearing this in mind, the generation of arc inscription for *roles* (*workers*) is a little bit trickier. This situation is modeled in Fig. 5.

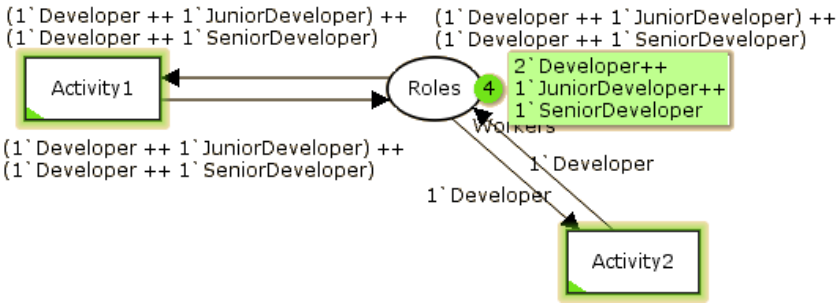


Fig. 5. Arc inscriptions for roles

Fig. 5 models a situation where we have two developers available, one *JuniorDeveloper* and one *SeniorDeveloper*. We have two activities that are parallel to each other, *Activity1* and *Activity2* (they are represented as single transition for simplicity, but in reality they are more complex, see Fig. 2). From the activity diagram and its meta-information (swimlanes) we know that *Activity1* requires two workers, one *JuniorDeveloper* and one *SeniorDeveloper*. *Activity2* requires only one worker and it can be either *JuniorDeveloper* or *SeniorDeveloper* (because of the generalization to *Developer* role). The two activities are in a conflict (even though they are parallel, they can't be executed at the same time) because if *Activity1* is executed, it allocates one *JuniorDeveloper* and one *SeniorDeveloper* for some period of time necessary to complete this task. While those two developers are working on *Activity1* no other developer required for *Activity2* is available (since we have only two developers available at any given time). You can picture it as a small team of two developers and the team has two user stories to implement (*Activity1*, *Activity2*). As you can see, the arc inscriptions of incoming arcs are different than in the *resources case* mentioned above. For each leaf of the inheritance tree, we have to also include in the inscription its superclasses all the way back to the root (the element which inherits from the *Worker* class). In our example, we have only one level of additional inheritance from *Worker* specializations (the *JuniorDeveloper* and *SeniorDeveloper*). So for each *JuniorDeveloper* or *SeniorDeveloper* we use in any inscription related to *roles* (the arc inscriptions and initial marking inscription of the *role pool place*) we have to add (via ++ multiset operator) its superclasses (in our case only one, the *Developer* class). If we would require two *JuniorDevelopers* for some activity, we would also have to include two *Developers* in the arc inscription. Applying that rule to the Fig. 5, if the *Activity1* is executed, one *JuniorDeveloper* and one *SeniorDeveloper* is allocated for the activity (accompanied with their superclass, one *Developer* for each required subclass instance) thus leaving the roles place empty. While the first activity is in progress, the second activity can't be executed because of the lack of needed developers. The other way around, if we execute *Activity2* one *Developer* is allocated (we don't know which one because the *JuniorDeveloper* and *SeniorDeveloper* tokens stays in roles place) thus blocking execution of the *Activity1*. In both scenarios of the conflict the model behaves as expected thus proving itself as a sound solution.

As mentioned, in case of the *artifacts*, only arcs from an activity to *artifact pool place* are added. The arc inscriptions are generated from object associations as it is in the *resources case*.

We finish the CPN model by putting certain amount of tokens (based on estimated or real numbers of resources and workers) to resource pool place and to roles pool place.

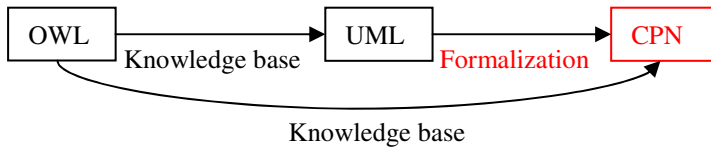


Fig. 6. Overall architecture

Let's close this chapter by presenting the overall architecture of the proposed solution on Fig. 6. The architecture consist of three parts, where the OWL is used as knowledge base for semi-formal UML model and as a knowledge base for the process of formalization to formal CPN model of the business process. The knowledge base consists of definition of domain specific entities via meta-modeling and definition of concrete resources, workers and artifacts along with their constraints towards the process. UML is used for modeling dynamic aspect of the process which is then converted to formal CPN model with help of OWL database.

4 Benefits of the Approach

Proposed software process transformation approach could be a powerful tool that may help adapting the quickly changing needs of software development market and enrich the process building with reusable formal rules. The benefits of the approach are:

- Complex software process methodology based on combination of well-known techniques – UML and CPN. Readability and ease of use of the framework are preserved by utilization of semi-formal models combined with strong semantics of Petri nets;
- Assuming the UML models are serialized as XML we can formalize the diagrams to XML model which is supported by CPN tools using simple XSLT template;
- Powerful simulation and analysis methods provided by CPN tools;
- Approach would improve process planning, analyze and evaluate tasks quantitatively, assess costs and benefits of applying new tools and technologies on a project, train and support project staff and improve the communication with the customer. Now extended with enhanced simulation of reality because of the resource conflict modeling and statistics, artifact generation statistics and worker conflict modeling and utilization statistics;

5 Conclusion and Future Work

In this paper, we have presented that the software process modeling using semi-formal techniques (represented by UML) could be successfully supported with powerful formal modeling language. We have proposed set of intuitive mapping rules for conversion between *UML2* Activity diagram and *CPN*. We have shown the undisputable benefits of modeling with *CPN* featuring state space analysis which provides powerful method of analyzing reachability of each activity, simulation which can be used in a step-by-step manner to debug to find a root cause of possible deadlocks in modeled process and of course means for validation and verification of the model itself.

Some features of Activity diagrams such as the data part and *swimlanes* were left out because we need a good balance between readability and completeness so we have to choose what is really essential for this approach. By incorporating data component to *CPN* in form of additional *colorsets* and inscriptions the net model would grow rapidly.

This method is trying to support the process by combining both approaches to get rid of ambiguity of modeling with semi-formal methods and discomfort and rigidity of modeling with formal methods.

Obviously, further research and studies are required to cover all problems as the domain of software process management and requirements specification are quite complex and current incomplete methodology is far from being possible to solve the discussed issues.

Our future work will be dedicated to the improvement and formalization of this methodology. Improvements will be especially on the processes, proper description of formal procedures and creation of *UML2* Activity diagram to *CPN* compiler that would become a headstone for other independent projects. To avoid any procedural flaws and misunderstandings, our task will be to create a proper formal procedure (description) of the methodology.

Acknowledgment. The research was supported by the internal grant agency of VSB Technical University of Ostrava, Czech Republic, project no. SP2014/157 "Knowledge modeling, simulation and design of processes". Michael A. Košinár is supported as a *Grand aided student of Municipality of Ostrava, Czech Republic*.

References

1. Aalst, W.M.P., van der Hee, K.M., van Houben, G.J.: Modeling Workflow Management Systems with High-Level Petri-Nets. In: De Michelis, G., Ellis, C., Memmi, G. (eds.) Proceedings of the Second Workshop on Computer-Supported Cooperative Work, Petri nets and Related Formalisms, pp. 31–50 (1994)
2. Aalst, W.M.P.: The Application of Petri Nets to Workflow Management. The Journal of Circuits, Systems and Computers 8(1), 21–66 (1998)

3. Allemang, D., Hendler, J.: Semantic Web for the Working Ontologist: Effective Modeling in RDFS and OWL (2008)
4. Brooks, F.P.: No Silver Bullet - Essence and Accidents of Software Engineering (reprinted from information processing 86, 1986). *Computer* 20(4), 10–19 (1987)
5. Czopik, J., Kosinar, A.M., Stolfa, J., Stolfa, S.: Formalization of Software Process Using Intuitive Mapping of UML Activity Diagram to CPN. Paper presented at the Proceedings of the 5th International conference on Innovations in Bio-Inspired Computing and Applications, Ostrava (2014)
6. Du, Z., Yang, Y., Xu, J., Wang, J.: Mapping UML Models to Colored Petri Nets Models based on Edged Graph Grammar (2011)
7. Jennings, N.R., Faratin, P., Norman, T.J., O'Brien, P., Odgers, B.: Autonomous agents for business process management. *International Journal of Applied Artificial Intelligence* 14(2), 145–189 (2000)
8. Jensen, K., Kristensen, L.M.: Coloured Petri Nets: modelling and validation of concurrent systems. Dordrecht: Springer. ISBN 978-3-642-00283-0 (2009)
9. Jung, K., Joo, S.: Transformation of an activity model into a Colored Petri Net model (2010)
10. Kammer, P.J., Bolcer, G.A., Taylor, R.N., Hitomi, A.S., Bergman, M.: Techniques for supporting dynamic and adaptive workflow. *Computer Supported Cooperative Work* 9 (3-4), 269–292 (2000)
11. Kaufmann, M., Silver, G.A., Lacy, L.W., Miller, J.A.: Ontology based representations of simulation models following the process interaction world view. Paper presented at the Proceedings of the 38th conference on Winter simulation, Monterey, California (2006)
12. Klein, M., Dellarocas, C.: A knowledge-based approach to handling exceptions in workflow systems. *Computer Supported Cooperative Work* 9(3-4), 399–412 (2000)
13. Košinár, M., Štolfa, J., Štolfa, S.: Knowledge Support for Software Processes. In: Proceedings of the 24th European-Japanese Conference on Information Modeling and Knowledge Bases
14. Kowalski, T.: Net Verifier of Discrete Event System models expressed by UML Activity Diagrams (2006)
15. Kruchten, P.: The rational unified process: an introduction. 3rd ed., xviii, 310 s. Addison-Wesley, Upper Saddle River (2004) ISBN 03-211-9770-4
16. Luati, A., Jerad, C., Barkaoui, K.: On CPN-based Verification of Hierarchical Formalization of UML 2 Interaction Overview Diagrams (2013)
17. Narendra, N.C.: Flexible support and management of adaptive workflow processes. *Information Systems Frontiers* 6(3), 247–262 (2004)
18. Raffo, D.M.: Modeling software processes quantitatively and assessing the impact of potential process changes on process performance. Carnegie Mellon University (1996)
19. Staines, T.S.: Intuitive Mapping of UML 2 Activity Diagrams into Fundamental Modeling Concept Petri Net Diagrams and Colored Petri Nets (2008)
20. Vergidis, K., Tiwari, A., Majeed, B.: Business Process Analysis and Optimization: Beyond Reengineering. *IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews* 38(1), 69–82 (2008)

Flash Assisted Segmented Bloom Filter for Deduplication

Girum Dagnaw¹, Amare Teferi², and Eshetie Berhan³

¹ Addis Ababa Universtiy, ICT-COE,
Addis Ababa, Ethiopia
girumdagnaw@gmail.com

² Afmobi,
Addis Ababa, Ethiopia
amarechina@gmail.com

³ Addis Ababa Institute of Technology
Addis Ababa, Ethiopia
eshetie_ethio@yahoo.com

Abstract. It is known that large amount of unique fingerprints are generated during deduplication for storage on the cloud. This large number of fingerprints is often tackled using bloom filters. Usually bloom filters are implemented using Random Access Memory (RAM) which is limited and expensive. This leads to higher false positive probability which decreases the deduplication efficiency. In this paper, a Flash Assisted Segmented Bloom Filter for Deduplication (FASBF) is proposed which implements its bloom filter (BF) on solid state drive where only part of the whole bloom filter will be kept in RAM while the full bloom filter is on solid state drive. This improves duplicate lookup in three ways. First the size of the bloom filter can be sufficiently large. Second, more number of hash functions can be used. And last, there will be more RAM space for fingerprint cache. This approach is evaluated using a prototype implemented on a study oriented network backup system. The result shows that this approach saves a considerable amount of memory space while satisfying an underlying 100MB/s backup throughput.

Keywords: Segmented bloom filters, Solid state drives, Deduplication.

1 Introduction

Backup service providers use a technique called deduplication which is a data compression technique which eliminates duplicate data at different granularities. In deduplication, unique files or chunks are identified and stored.

Given the probability that similar byte patterns will appear hundreds, if not thousands, of times deduplication process will result in significantly less data being stored or transmitted, when compared with the data which should have been stored or processed had it not been deduplicated.

In deduplication, having smaller chunk size improves deduplication ratio, but it creates duplicate-lookup disk bottleneck problem. This problem arises as a result of keeping chunk indexes in the RAM. As data volume - which the index represents -

grows, the index also grows and will reach to a point where it will be very big for the RAM to accommodate as a whole. This, traditionally, is solved by putting part of the index on disk and fingerprint lookups will often need to be done on the on-disk index.

Performance wise this traditional approach of putting the indexes on disk is unthinkable due to the frequent accesses to the indexes for lookups and updates.

The other method is to use Bloom filters [1], which are space efficient probabilistic data structures which tell if an element is NOT a member of a set with probability 1. Because of its bit representation, bloom filters occupy very less space and hence can handle huge amount of data in a very short time.

The use of bloom filters solves the problem of duplicate-lookup disk bottleneck to some degree. But when the number of items inserted to bloom filter grows large most of the bits in it are set and hence most of the membership queries result in false positives. This in turn will result in frequent on-disk index lookups.

With a quest for better performance studies like [2][3][4] have implemented flash memory to have dynamically growing bloom filters which can go beyond the size of available RAM space.

In this paper, a Flash Assisted Segmented Bloom Filter for Deduplication (FASBF) is proposed which implements its bloom filter (BF) on solid state drive where only part of the whole bloom filter will be kept in RAM while the full bloom filter is on solid state drive. This improves duplicate lookup in three ways. First the size of the bloom filter can be sufficiently large. Second, more number of hash functions can be used. And last, there will be more RAM space for fingerprint cache.

Contribution of this paper is: Bloom Filter space is segmented (Segmented Bloom Filter) and these segments are grouped into a bloom filter array (Segmented Bloom Filter Array) where an array of these Segmented Bloom Filter Arrays is stored partially in RAM and as a whole in SSD (Solid State Drive). The SSD is also used to store full file fingerprint HBM(Hush Bucket Matrix) and partial HBM of chunk fingerprints.

The rest of the paper is organized as follows. Section II gives brief introduction of standard bloom filters and design of two bloom filters which are implemented on flash memory. The architecture and detailed design of FASBF will be presented in section III. Section IV describes the environment FASBF is evaluated on and the result of the evaluation. And finally section V gives conclusion.

2 Literature Review

2.1 Standard Bloom Filters

A bloom filter is a probabilistic data structure that is used to test whether an element is a member of a set. It is implemented using bit arrays, with m bits, which have all bits initially set to zero. k different hash functions will each be used to map or hash some set element to one of the m array positions.

One big problem with standard bloom filters is false probability, which is a condition where all array positions being checked for an element are 1 even when that element does not exist.

In order to reduce the false positive probability there has to be a balance between the potential number of elements, the size of the bloom filter and the number of hash functions.

Given bloom filter bitwise array m and the number of hash functions k we can compute the expected false positive probability f with:

$$f = \left(1 - \left(1 - \frac{1}{m}\right)^{kn}\right)^k \simeq \left(1 - e^{-\frac{kn}{m}}\right)^k \quad (1)$$

As in Broder et al [1], given bloom filter's bitwise m and total number of elements n , the optimal number of hash functions k can be expressed as:

$$k = \left(\frac{m}{n}\right) \times \ln 2 \quad (2)$$

If we have lower bound ϵ to the false positive rate f , i.e., $f \leq \epsilon$, we can deduce that the bitwise m should satisfy

$$m \geq n \frac{\log_2(1/\epsilon)}{\ln 2} = n \log_2 e \cdot \log_2(1/\epsilon) \quad (3)$$

2.2 Other Bloom Filter Designs on Flash

Canim et al. [2] designed a Bloom Filter for flash memory which pre-allocates a single large space on flash as bloom filter. This bloom filter space will be divided into n number of equal spaces or sub-bloom-filters and every of these sub-bloom-filters will have a corresponding but smaller sized partition in RAM which will be checked first while searching for keys and at the same time is used as buffer when inserting into the bloom filter.

The main drawback of this design arises from the size gap between the in-RAM buffer space and the sub-bloom-filters in flash. Frequent block erase and page writes are required because the in-RAM buffer space becomes full fast due to its small size.

In Lu et al. [3] a forest-structured bloom filter is designed to solve the drawbacks of [2] and optimize the process of both key insertion and lookup on flash memory while it can also grow dynamically.

When a root layer bloom filter in RAM gets full a new layer of bloom filter (considered as chilled bloom filter) will be created in flash. The root-layer bloom filter will be written to flash in its entirety leaving the space it occupied in RAM for buffer space used to delay insertions. This buffer space will be partitioned in a similar fashion as that of [2] and at any time, only sub-BFs at one layer of the forest will be buffered in RAM.

As the size of data grows the bloom filter design of [3] will face two problems. First, all lookups for keys which does not exist in the set will have flash reads equal to the height of the forest. This will reduce the performance of the application using the bloom filter. Second, after a few layers are added to the flash there will be more number of blocks which needs to be buffered in RAM. This results in less buffer space per block, which in turn causes more flash write operations, leading to performance degradation and a shorter life time of the flash memory, Debnath B. et al [5].

3 Architecture and Design of FASBF

This experiment was done in a study oriented network backup system environment. This storage system has Backup Server (BS), Storage Proxies (SP) and Metadata Servers (MDS) at the front end and the backend consists of clustered storage nodes which themselves can have one or more Storage Components (SCs).

Upon arrival of backup jobs from clients, the BS splits metadata from files and sends them to the MDS. Then Rabin fingerprinting algorithm, which implements content-defined and variable-sized chunking, is used in SP in order to divide the files received from BS. SP also computes the fingerprints of both the files and the chunks. The file recipes which are necessary for archival will be generated. These file recipes and chunk contents generated by SP are distributed among multiple SCs in their backup sequences.

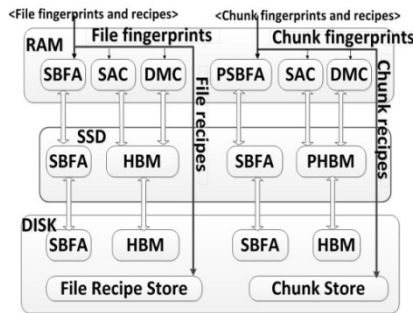


Fig. 1. Architecture of FASBF

As shown in Fig. 1, the HBM (Hash Bucket Matrix) is organized in two ways. The HBM for the file fingerprints will be stored both in SSD and disk in its entirety. Because of its large size only part of the HBM for chunk fingerprints will be stored in SSD and so the whole part will be stored in disk.

The bloom filters, organized as Segmented Bloom Filter Array (SBFA), will be stored in RAM, SSD and disk. The design of FASBF required the use of segmented bloom filters for the sake of dividing its space and storing its partial content in RAM.

3.1 Segmented Bloom Filter Array in RAM and SSD

FASBF implements bloom filters in two phases. While the second phase involves SSD, the first phase is done totally on the RAM.

In phase one (shown in Algorithm 1-(a)), the bloom filter will be inserted into and searched like any other segmented bloom filter. After deciding on the size of the bloom filter m and the number of hash functions k using the equations given in equations 2 and 3, the required bloom filter space is allocated in RAM. This bloom filter space is divided into k equal spaces and assigned to the k different hash functions, h_1 through h_k . When this bloom filter (current bloom filter) reaches its capacity it is copied to SSD and is reset. This marks start of the second phase.

In phase two (shown in Algorithm 1-(b)), the first k/p consecutive segments of all the individual bloom filters (PSBFA) will be in RAM (as shown in the upper part of

Fig. 2) along with the current bloom filter¹. Here p is the fraction of segments from the individual SBFAs in SSD chosen to create the in-RAM PSBFA.

Setting higher values for p will save more RAM space but the false positive probability for the lookups on the in-RAM PSBFA will be higher resulting in frequent accesses to the in-SSD FSBFA. Choosing p to be 1 will mean the in-RAM PSBFA will be the exact copy of the FSBFA.

Phase one:

1. Decide size of m and k : // m is current bloom filter
2. Allocate m and divide m in to k segments;
3. Initialize all m bits: $m = 0$;
4. **while** (m is not full)
5. hash(fingerprint):
6. **for** ($j=0; j < k$)
7. use hash function j : $hj(\text{fingerprint})$ in to segment j ;
8. **end**
9. copy m to SSD;

(a)

Phase Two:

1. Decide value of p ; // p is fraction of the k segments to keep in RAM
2. copy first k/p segments of m into PSBFA;
3. reset m : $m=0$;
4. do
5. **while** m is not full
6. hash(fingerprint):
7. **for** ($j=0; j < k$)
8. use hash function j : $hj(\text{fingerprint})$ in to segment j ;
9. end
10. copy m to SSD;
11. append first k/p segments of m into PSBFA;
12. reset m : $m=0$;
13. **while** (there is fingerprint to be hashed)

(b)

Algorithm 1. Hashing in to PSBFA

¹ Current bloom filter is a segmented bloom filter array which is currently in RAM and is being inserted into and has similar structure to that of any the rows on the bottom part of Fig. 2.

Fingerprint insertions are made only into the current bloom filter. Hash function k_1 will compute a bit position in its assigned segment h_1 and set that position to 1, and hash function k_2 will compute and set a bit position in h_2 , and so on. When the current bloom filter reaches its capacity, its first k/p segments will be appended to the PSBFA and the entire SBFA will be appended to the in-SSD FSBFA. This insures that the rows in the PSBFA have one-to-one correspondence with the rows in the FSBFA which is important during fingerprint queries.

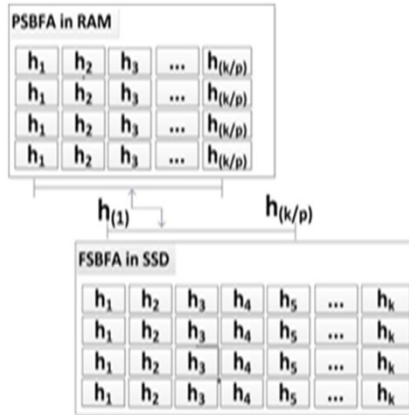


Fig. 2. Partial Segmented Bloom Filter Array (PSBFA) in RAM and Full Segmented Bloom Filter Array (FSBFA) in SSD

Whenever the bloom filters are to be searched for fingerprints, first the current bloom filter will be looked up and if it is not found there the PSBFA will be searched. When the lookups on this PSBFA return true the remaining part of the corresponding SBFA from SSD will be brought to RAM and checked. The design of FASBF makes it possible to identify the row number of the row(s) in the PSBFA which returned true for the lookup. This row number(s) will be used to read the corresponding row(s) of the FSBFA in the SSD.

3.2 Hash Buckets in SSD and HDD

In FASBF the Rabin fingerprint algorithm is used to determine variable sized chunks of files or objects. An MD5 hash value of these chunks or objects, their offsets and length comprise an entry of the buckets.

FASBF supports file level and chunk level deduplication, and maintains two independent sets of HB for each of them. FASBF stores partial copies of the chunk fingerprint buckets and full copies of file fingerprint buckets in SSD as shown in Fig. 1. When the size of the file fingerprint grows large, it is copied to disk and only part of it will be kept in SSD making it similar to that of the fingerprint bucket.

Though the chunk duplicate locality in network storages is less than that of DDS's [7] which is done on disk-to-disk backup, it still can be exploited for better performance. This is achieved by maintaining a tanker for every independent SBFA.

While fingerprints are hashed into the current bloom filter array their most significant bits will be used to identify the hash bucket they should be put into. This insures that chunk fingerprints which belong to a file or a backup job are inserted into one tanker. It also makes it easy to identify which bucket to look into during lookups.

3.3 FASBF Work Flow

Because of the limited available space only the chunk level deduplication process is explained here.

When a fingerprint arrives for checking, first the current bloom filter is checked. If there is a positive response and the current tanker² also returns true, it means the chunk is a duplicate so is not stored. Negative response means we need to further check the fingerprint in the in-RAM partial bloom filter (PSBFA). If this query returns negative, the fingerprint is definitely unique and needs to be processed accordingly. If the query on the PSBFA returns true, the chunk fingerprint needs to be checked in the corresponding in-SSD FSBFA³. A negative response from the query on the in-SSD FSBFA means a unique chunk fingerprint and a positive response means the corresponding tanker need to be checked for confirmation.

```

1. Check current BF for fp; // fp=md5(chunk);
2. if FOUND
3.   check current tanker;
4.   if FOUND
5.     fp is duplicate;
6.   end
7. else
8.   check PSBFA;
9.   if NOT FOUND
10.    hash fp to current BF; // fp is unique;
11.  else
12.    check FSBFA;
13.    if NOT FOUND
14.      hash fp to current BF; // fp is unique;
15.    else
16.      fp is duplicate;
17.    end
18.  end
19. end

```

Algorithm 2. Chunk level deduplication

² Current tanker is a tanker which resides in RAM and corresponds to the current bloom filter.

³ Even though there are two steps involved in querying the BF a hash value h_j by a hash function $h(j)$ is computed only once at the beginning.

4 Prototype Implementation and Evaluation

4.1 Prototype Implementation

The prototype used to evaluate FASBF is developed using Visual C++ programming language in the Visual studio 2010 express environment.

The prototype was then run on windows server environment. Hardware setup of this server constitute an Intel(R) Xeon(R) CPU (2 CPUs each with 4 cores and at 2.13GHz), 8 GB RAM, 2GB network interface card. A 64GB solid state drive with a sequential read/write speed of 260/100MB/s respectively is used.

4.2 Data Used for the Study

Data from 20 users representing wide range of professions is used. Users were backing up their data for one month. They backup their data incrementally on a daily basis and do full back up on every seventh day of their previous full backup (as a mandatory) or any day (optional). At the end of this period a total of 11.089TB of data constituting a total of 39.69 million files was backed up.

4.3 Result and Evaluation

The performance of FASBF is evaluated using three performance measures. While its throughput and memory consumption are compared against the traditional approach which uses only RAM for its bloom filter, its efficiency is evaluated by comparing results for 4kb and 8kb average chunk sizes.

Zero Length Chunks

As a result of file systems representing file holes with null characters (zero length strings), the Rabin fingerprinting algorithm generated a multitude of chunks which just contain a zero length string which are common even among dissimilar files.

In order to speed up the deduplication process the MD5 of these zero length strings is deduplicated directly without checking in the bloom filter and hence the HBMs.

From the set of the data collected a total of 12005360 zero length chunks were observed in the 8KB average chunk size set up while 17279830 zero length chunks were found in the 4KB average chunk size setup.

Deduplication Efficiency

The Rabin fingerprint algorithm was run twice with different chunk size settings on the same data. In the first setup the Rabin fingerprinting module was set to run with a 1KB, 4KB and 32KB minimum, average and maximum chunk sizes respectively. And the second set up was 2KB, 8KB and 64 KB minimum, average and maximum chunk sizes respectively.

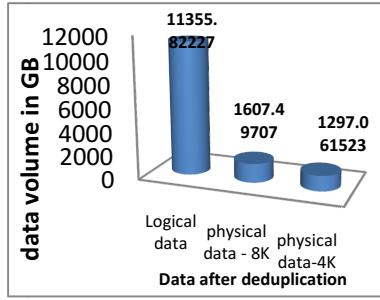


Fig. 3. Deduplication efficiency for 4KB and 8KB average chunk sizes

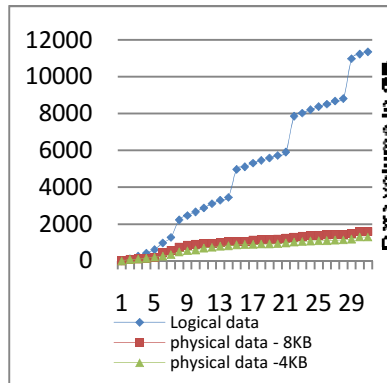


Fig. 4. Deduplication efficiency for 8KB and 4KB average chunk size

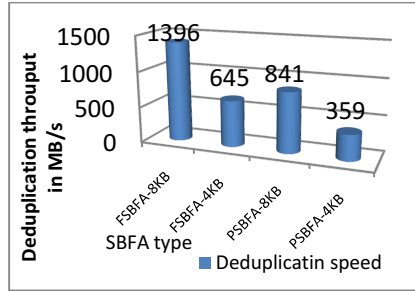
In the first setup, a total of 2.77×230 fingerprints were generated. Among these fingerprints 187,708,985 were unique for a total size of 1.26TB. The second set up resulted in 1.38×230 fingerprints among which $1.7 \times 229 (0.17 \times 230)$ were unique for a total size of 1.57TB. As shown in Fig. 3, the deduplication ratio achieved in the first set up is much better than the second. But, as the number of fingerprints increased, so did the memory requirement for the bloom filter and the storage space required for the fingerprints.

Deduplication Throughput

Deduplication throughput of FASBF is evaluated using the two data sets collected using the two different sets of average chunk size configurations. For both of these setups, first the data sets are deduplicated in the traditional way by putting the whole of the bloom filter array in RAM (FSBFA) in its entirety. Second, the bloom filter array was put partially in RAM and fully in SSD (PSBFA). Table 1 gives summary of these configurations.

Table 1. Segmented Bloom Filter Array test configuration

2K-8K-64K (MIN-AVG-MAX)		1K-4K-32K (MIN-AVG-MAX)	
Full BF in RAM (FSBFA)	Partial BF in RAM (PSBFA)	Full BF in RAM (FSBFA)	Partial BF in RAM (PSBFA)

**Fig. 5.** Deduplication speed in MB/s

As can be observed in **Chyba! Nenalezen zdroj odkazů.**, the deduplication speed is maximum for the 8KB average chunk size and full bloom filter in RAM configuration and it is the least for the 4KB average chunk size and partial bloom filter in RAM configuration.

As the four test cases in the chart are generated from the same data, it can be seen that the size of the chunks play great role in deciding the throughput of deduplication. In both cases - where SSD is used and not used - the deduplication speed is almost halved when the average chunk size is reduced from 8KB to 4KB. The impact of using SSD on the speed of deduplication is even greater since it reduced the deduplication throughput by approximately 51% for the 8KB average chunk size by a little over 75% for the 4KB average chunk size configuration.

Memory Consumption

SSDs are not used to support file level deduplication because the memory consumption for the file level bloom filters is relatively very small.

In the 1KB-4KB-32KB chunk size configuration there were 2839 million chunk fingerprints (assuming 4KB average chunk size) and this number halved to 1419 million for the second configuration. In both cases the false positive probability is kept below $\frac{1}{2} \cdot 20$.

When not using the SSD, the first configuration required 10.132GB of RAM and the second 5.132GB of RAM. Using the SSD, the RAM requirement was dependent on the number of bloom filter segments kept in RAM. In this evaluation two-third ($\frac{2}{3}$) of the bloom filter segments were put in RAM. Hence, the memory requirement was 6.8GB and 3.46GB for the 4KB and 8KB average chunk size configurations respectively.

Table 2 gives summary of the memory consumption.

Table 2. Segmented bloom filter memory consumption summary

(MIN-AVG-MAX) chunk size configuration	No of fingerprints In millions		Memory consumption (in GB)			
			File level (Size in MB)	Chunk level		Total (With SSD)
	File level	Chunk level		Without SSD	With SSD	
1K-4K-32K	39.69	2839	136.52	10	6.66	10.132 (6.8)
2K-8K-64K	39.69	1419	136.52	5	3.33	5.132(3.46)

5 Conclusion

Chunking a stream of data into very small sized segments gives the advantage of having less data to store after deduplication and less data to transfer over the network. But it reduces the deduplication speed to a much degree and consumes a considerable amount of memory.

Another important factor which determines the deduplication speed is where the bloom filter is put. By putting the bloom filter partially in RAM and entirely in SSD, a great amount of RAM space is saved for further use. But, considerable amount of speed compromise is required. **Chyba! Nenalezen zdroj odkazů.** shows that when SSD is used the deduplication speed drops to well below half of the speed when it is not used. When this is combined with a value of 4KB for average chunk size the deduplication speed falls further to less than 25%.

Because it greatly reduces the deduplication throughput the use of SSD in areas where there is abundant network bandwidth might not be feasible. But in areas where there is limitation on the availability of network bandwidth it will be tempting to consider its use.

A good option to support backup on the network in relatively limited bandwidth network areas is to transfer as small amount of data as possible. This can be achieved by removing as much duplicate data as possible at the source. This in turn can be achieved by splitting data to very small chunks like an average size of 4KB or even less.

In order to be in the right range of false positive probability and thus reduce the disk access needed to verify these false positives, SSDs can be implemented. As shown in the evaluation result, a rough 360MB/s deduplication can be achieved with SSD.

References

1. Broder, A., Mitzenmacher, M.: Network Applications of Bloom Filters: A Survey (2003)
2. Canim, M., Mihalia, G.A., Bhattacharjee, B., Lang, C.A., Ross, K.A.: Buffered bloom filters on solid state storage (2010)
3. Debnath, B., Du, D.H.C., Lu, G.: A Forest-structured Bloom Filter with Flash Memory. In: Mass Storage Systems and Technologies, MSST (2011)
4. Mokbel, M.F., Lilja, D.J., Du, D., Debnath, B.: Deferred updates for flash-based storage. In: Mass Storage Systems and Technologies (MSST), Washington, DC, USA (2010)
5. Jiang, H., Zhou, K., Feng, D., Wei, J.: MAD2: A Scalable High-Throughput Exact Deduplication Approach for Network Backup Services. In: 26th IEEE MSST, Incline Village, NV, USA (May 2010)
6. Li, K., Zhu, B.: Avoiding the Disk Bottleneck in the Data Domain Deduplication File System. In: 6th USENIX Conference on File and Storage Technologies
7. Policroniades, C., Pratt, I.: Alternatives for Detecting Redundancy in Storage Systems Data. In: Proceedings of the 2004 USENIX Annual Technical Conference, Boston, MA, USA (June 2004)
8. Chazelle, B., Kilian, J., Rubinfeld, R., Tal, A.: The bloomier filter: an efficient data structure for static support lookup tables. In: Proceedings of the Fifteenth Annual ACM-SIAM Symposium on Discrete Algorithms, Philadelphia, USA, pp. 30–39 (2004)
9. Meister, D., Brinkmann, A.: dedupv1: Improving deduplication throughput using solid state drives (SSD). In: MSST 2010 Proceedings of the 2010 IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST), Washington, DC, USA, pp. 1–6 (2010)
10. Debnath, B., Sengupta, S., Li, J.: ChunkStash: Speeding up Inline Storage Deduplication using Flash Memory. In: 2010 USENIX Annual Technical Conference (ATC) (June 2010)
11. Debnath, B., Sengupta, S., Li, J., Lilja, D.J., Du, D.: BloomFlash, D.: Bloom Filter on Flash-based Storage. In: International Conference on Distributed Computing Systems (ICDCS), Minneapolis, USA (2011)
12. Tarkoma, S., Rothenberg, C.E., Lagerspetz, E.: Theory and Practice of Bloom Filters for Distributed Systems. *IEEE Communications Surveys & Tutorials* 14(1), 131–155 (2012)
13. Rothenberg, C.E., Macapuna, C.A.B., Verdi, F.L., Magalhães, M.F.: The deletable Bloom filter: a new member of the Bloom family. *IEEE Communications Letters* 14(6), 557–559 (2010)
14. Xia, W., Jiang, H., Feng, D., Hua, Y.: Silo, A similarity-locality based near-exact deduplication scheme with low ram overhead and high throughput. In: USENIX Conference on USENIX Annual Technical Conference, Berkeley, CA, USA

A Train Run Simulation with the Aid of the EMTP – ATP Programme

Maroš Ďurica

Department of Electrical Engineering, Faculty of Electrical Engineering and Computer Science
VŠB – Technical University of Ostrava, Ostrava, Czech Republic
d.maroseznam.cz

Abstract. Demandingness of the energy calculation by means of the rail vehicle momentary output time integration method lies in a need for detailed input data as a tachogram - a time-speed relation and a corresponding passed distance which can be obtained by a calculation or from a train entrepreneur. A locomotive load value, comprising wagons and resistance components of a track, determines a required tractive power value on a wheel circumference. Each traction power value is assigned to a traction motor current value in traction characteristics, which can be further used for a locomotive output calculation. To make the calculation more accurate a model has been created in the EMTP – ATP software, its advantage is to simulate a dynamic behavior of a model which is affected by more influences at the same time. The simulation time can be shortened in an exact ratio to the real time and thus an evaluation and implementation of changes of the model parameters can be performed more quickly.

Keywords: rail transportation, mathematical model, vehicle driving, motors, impedance.

1 Introduction

Railways play a significant role in the World transport network. For assessing the economic costs, a result of an energy consumption calculation in the given segment of the railway track is a significant factor. The train run energy consumption calculation can be performed as follows:

- Time integration based on the vehicle input. Relation of the driving vehicle output and speed on time needs to be known for the calculation.
- Determination of consumption based on traction work and loss components.

The aim of my observation is to improve the method of the electrical energy calculation with the help of the method of the locomotive output direct integration using the EMTP – ATP programme model. By means of the EMTP – ATP programme the progression of voltage and current in traction motors during various periods of a run has been simulated. The main issue in this observation is how much the given train

model comes close to the actual conditions. The influence of the line feeder resistance and the locomotive load value has been included in the model.

A run of a train with a 182 type locomotive has been simulated in the Přerov - Česká Třebová track segment (track no.270 according to the Czech Railways nomenclature), to be specific, between the stations with traction substations Grygov, Červenka, as shown in Fig. 3, where a placement of the traction substations corresponding with the actual condition can be seen [5].

2 Train Data

The 182 type locomotive output control is a resistance one, traction motors are direct-current six-pole series electric motors sequenced in three ways:

- Series (6 motors in a series, $U = 500$ V),
- 1. Series-parallel (2 legs of 3 series-connected motors, $U = 1000$ V),
- 2. Series-parallel (3 legs of 2 series-connected motors, $U = 1500$ V).

For modeling, resistance of six motors is considered with current and resistance depending on speed according to the characteristic shown in Figure 5 [4].

The locomotive comprises 8mH smoothing choke connected in a series with a RL component replacing a series traction motor. The motor resistance varies in dependence on speed, 5.95 Ω for 40 km/h speed and 6.66 Ω for 30 km/h speed, inductivity is of 11mH value. For the energy calculation, I have used data from the progression of voltage and current designated V30, V40 on the RL component to which a RC filter has been connected in parallel, serving for attenuation of transient phenomena when switching-over the locomotive speed.

The locomotive has been loaded with eleven Faccs type freight wagons of 657 t total weight, as shown in Figure 1 [2, 7].



Fig. 1. 182 type locomotive and Faccs type freight wagon

3 Data of Traction Mains and Traction Substations

The basis of the traction substation is a three-winding traction transformer 23/ 2 x 2.5 kV with a winding in Yyd, as it is shown in Figure 2. In the same Figure there is also a corresponding model of a rectifying block in the EMTP – ATP comprising a BCTRAN transformer with corresponding values of a traction oil transformer of 5.3 MVA output which is made by Power – Energo company [3].

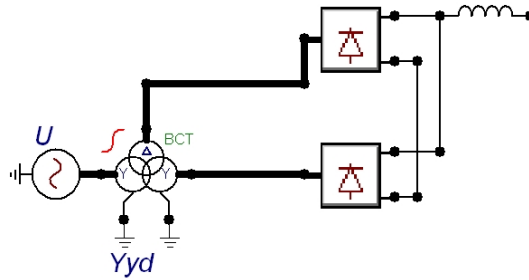


Fig. 2. The traction transformer and twelve-pulse rectifier model in the Emtip – Atp

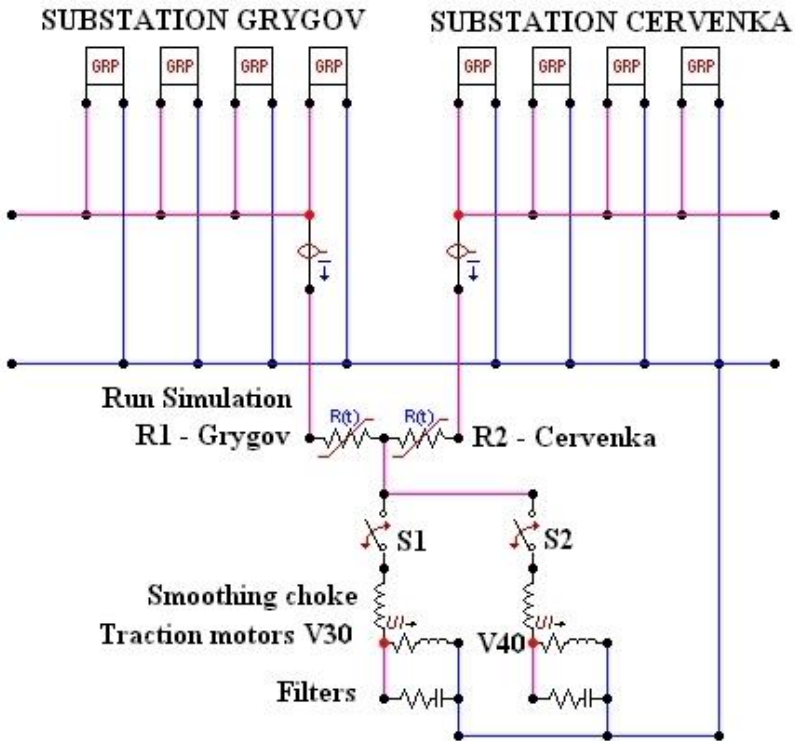


Fig. 3. 182 type locomotive traction motors in the series-parallel connection and a circuit diagram for the train run simulation in the Grygov – Červenka segment in the EMTP – ATP

To save a space in the scheme in Figure 3, the block with the transformer and the rectifier is depicted as a square with a designation GRP. Each of the substations consists of four blocks connected in parallel, connected to the trolley wire through a choke of 4 mH value, in the simulation it helps to smooth a progression of the rectified voltage and restrains the rate of the short-circuit current rise. The two-side power supply of the traction mains is realized by a Cu trolley wire 150 mm² and a Cu suspension cable 120 mm² together with an AlFe line feeder 240 mm² with 0.06 Ω/km specific resistance.

4 The Simulation of the 182 Type Locomotive Run in the Grygov – Červenka Track Segment

The train run according to the time-speed tachogram begins with starting up to 30 km/h speed in the Grygov station, goes-on at a constant speed for a minute and accelerates again up to 40 km/h speed which is maintained till the arrival to the Červenka station where the train begins to brake to stop. According to time the passed distance and the corresponding train tractive power has been calculated according to the particular phases of the run and the reduced gradient value, as shown in Figure 4.

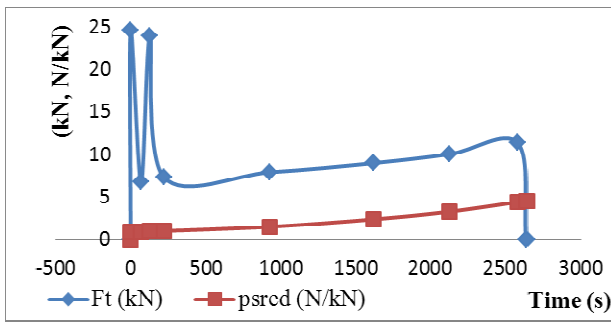


Fig. 4. Required tractive power on a wheel circumference and during the constant speed run, where it is increased depending on the track reduced gradient

The train run real time 2640 s is too long for the simulation in the EMTP – ATP programme, therefore it has been speeded-up a thousand times:

$$\frac{\text{Real time}}{1000} = \text{Simulation time (s, -, s)}. \tag{1}$$

For the 27.86 km distance between the stations the line resistance value is 1.67172 Ω. The train run simulation according to the progression of voltage and current V30 to V40 between the Grygov – Červenka stations is solved by two time-dependent resistors with a VA characteristic. The R1 resistor Grygov resistance value is changing depending on the distance in the segment (the given distance in km x 0.06 Ω), while the R2 Červenka resistance value is decreasing, so that the sum of the resistances in each time step is equal to the line resistance 1.67172 Ω between the stations.

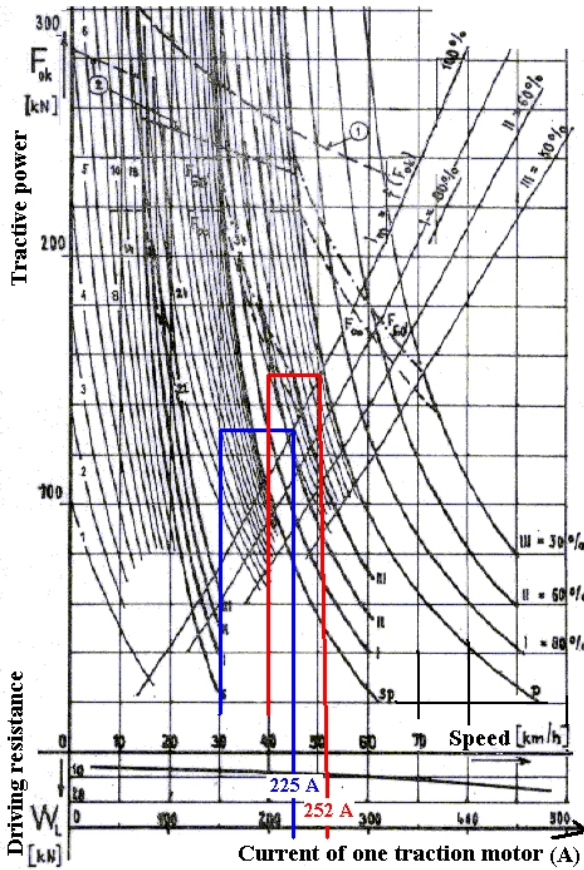


Fig. 5. 182 type locomotive traction characteristic, dependence of 1 motor current for speed 30 (blue line) and 40 km/h (red line)

Determination of the motor resistance values for the simulation has been performed according to the locomotive total current, whereas the current of one traction motor needed to be known [6]. Motor groups sequencing for 30 km/h speed is series-parallel with pre-sequencing a resistance step no.25, where one motor current is 225 A (cor-responding tractive power is 135 kN) and the locomotive total current is determined according to the formula below:

$$I_{30} = \frac{6 \cdot 225 A}{3} = 450 A . \tag{2}$$

A calculation of the resistance for the simulation:

$$R_{30} = \frac{3000 V}{450 A} = 6.66 \Omega . \tag{3}$$

Motor groups sequencing for 40 km/h speed is series-parallel with a shunt section no.2, where one motor current is 252 A (corresponding tractive power is 150 kN) and the locomotive total current is determined according to the formula below:

$$I_{40} = \frac{6 \cdot 252 \text{ A}}{3} = 504 \text{ A} . \quad (4)$$

A calculation of the resistance for the simulation:

$$R_{40} = \frac{3000\text{V}}{504\text{A}} = 5,95 \Omega . \quad (5)$$

5 Calculation of Electrical Energy Consumption of a Freight Train and the Grygov and Červenka Substations

A calculation of electrical energy consumption for the run and the Červenka and Grygov substations has been performed for the run at 40 km/h speed by means of the locomotive momentary output time integration method [1]. A period of a calculation step was used: 0.0001, to which a simulation time step 0.05 s corresponds according to the output data file View LIS file 0.05 s = 50 s real time step = 0.01388 hour. The output progressions of current and voltage of the traction motors at 30 and 40 km/h speed are shown in Figure 6.

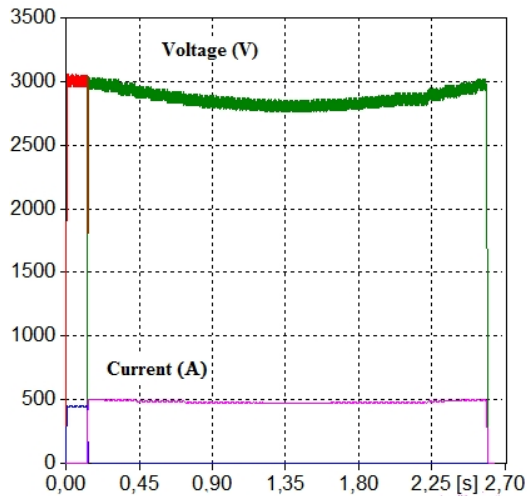


Fig. 6. The locomotive voltage and current progression at 30 – 40 km/h speed in various time periods during the run

The electrical energy consumption has been calculated for 40 km/h even speed, for this run phase the traction motor output has been calculated and integration for the simulation time step $t = 0.01388 \text{ h}$ has been performed.

$$f = Power = \frac{V_{40} \cdot I_{40}}{1000} \text{ (kW)} . \quad (6)$$

For E182 energy consumption calculation the integration has been performed using the Mathcad software, the result is a matrix of values of electrical energy consumption for the given time step. The result of the sum of all elements of the matrix is the total electrical energy consumption:

$$\sum E182 = \int f(t) \text{ (kWh)} = 914,5 \text{ kWh} . \quad (7)$$

In the railway transport the specific energy consumption (W.h / t.km) is determined which is the total consumption rate (W.h) related to t.km (ton-kilometre). At energy consumption 914.5 kW, the specific energy consumption is (designated as Wmsee):

$$Wmsee = \frac{T_{40} \cdot \sum E182 \cdot 1000}{M \cdot L_{40}} = 29 \text{ Wh/tkm} . \quad (8)$$

where **T40** is the run time at 40 km/h even speed – 0.6389 h,

M is a weight of the locomotive and wagons – 777 tons,

L40 is a length of the passed distance at 40 km/h even speed – 26.07 km.

It is necessary to add the specific energy consumption (designated as **Wpp**) of auxiliary drives as compressors, ventilators, control circuits; in literature, for 182 type locomotive the value **Ppp** = 52.2 kW is mentioned as the accessory circuits input.

$$Wpp = \frac{T_{40} \cdot P_{pp} \cdot 1000}{M \cdot L_{40}} = 1,6 \text{ Wh/tkm} . \quad (9)$$

The specific energy consumption for the locomotive collector at 40 km/h even speed **in summer** is $Wmsee + Wpp = 30.48 \text{ W.h / t.km}$.

For the Grygov and Červenka substations energy consumption calculation (designated as EGrygov and ECervenka) the integration of the output has been performed which is the result of a product of V1, V2 voltage and I1, I2 current according to the values of the voltage and current progression for 40 km/h even speed run for the simulation time step of 0.01388 h in the Mathcad software.

$$EGrygov, ECervenka = \int f(t) \text{ (kWh)} . \quad (10)$$

The result is a matrix of the EGrygov, ECervenka values, the sum of all the matrix elements gives the total electrical energy consumption of the substation:

$$\sum EGrygov = 486,6 \text{ kWh} . \quad (11)$$

The sum of the ECervenka values matrix is less by approximately 13kWh:

$$\sum ECervenka = 473,7 \text{ kWh} . \quad (12)$$

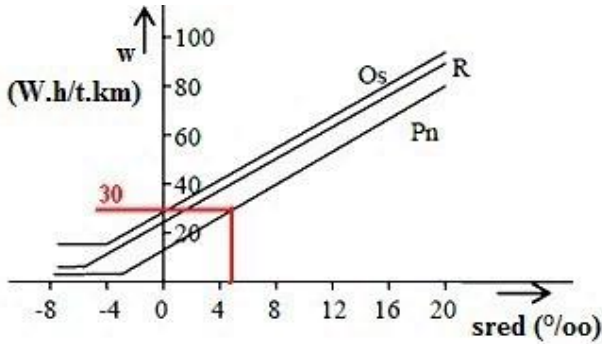


Fig. 7. Nomogram for determination of orientational specific energy consumption for typical trains

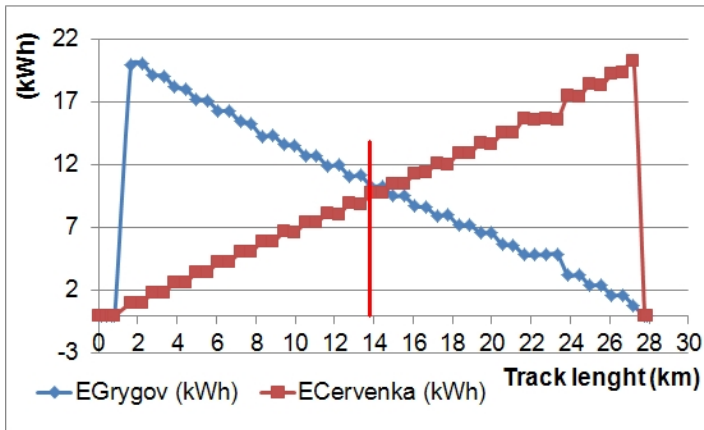


Fig. 8. Rate of the energy consumption rise and drop for the Grygov and Červenka traction substations during the train run at 40 km/h even speed with regard to the track segment length. The center of the supplied track segment is marked with the red line.

6 Conclusion

In the paper a problem of the specific energy consumption calculation has been solved for a selected train – 182 type locomotive with Faccs freight wagons when running at 40 km/h even speed in the Grygov – Červenka track segment. The accuracy of the calculation has been reduced due to not including the output of the traction motors at the train start up in the calculation, however, considering the start up phase lasting only few minutes in comparison to the total train run, omitting this phase of the run from the calculation has not caused a relevant deviation from the correct result. In this context, a greater error in the calculation can be caused by an inaccurate specification of the run time at the even speed, in no.8 formula indicated as T40.

A corresponding simulation in the EMTP – ATP programme has been performed, the resulting voltage and current progression values have been used for the specific energy consumption calculation by input integration for the given time.

The specific energy consumption result for summer is 30.48 W.h/t.km, which corresponds to a Pn through freight train consumption at a track reduced gradient up to 4.5 N/kN, as shown in the nomogram on Figure 7.

Further, electrical energy consumption calculation has been performed for the entire period of the run at 40 km/h even speed between the Grygov and Červenka substations, whereas a difference in the consumption of the particular substations is by 12.9 kWh higher for the Grygov substation which is caused by an unequal length of the distance passed at 40 km/h even speed with regard to the center – 13.93 km of the two-side supplied segment, as shown in Figure 8. Thus, the energy consumption of the substations of the track segment with a two-side power supply depends on a position of the train starting or stopping in the segment.

References

1. Paleček, J.: Electroenergetics in Transportation, study materials for lectures (Book style), Ostrava, VŠB – Technical University of Ostrava (Winter Semester, 2010)
2. Encyclopedia of locomotives. Locomotive 182 (2004c), <http://www.atlaslokomotiv.net/loko-182/> (retrieved November 18, 2013)
3. Power Energo. Traction transformer (2011c), <http://www.power-energo.cz/p114-soubory-ke-stazeni-specialni-transformatory.html> (retrieved October 14, 2013)
4. Šíroký, J.: Traction characteristics (2007c), http://homen.vsb.cz/~s1i95/mvd/vozidla_e (retrieved 2013 October 19, 2013)
5. Railway tracks history. Tracks data (2011c), <http://historie-trati.wz.cz> (retrieved September 5, 2013)
6. A mathematical model for a train run. Simulation of a train run (2008c), <http://diplom.utc.sk/wan/1881.pdf> (retrieved May 1, 2014)
7. Company RYKO PLUS. Wagon type Faccs 38 m³ (2010c), <http://www.rykoplus.cz/doc/typovy-list-faccs-38-m3.pdf/> (retrieved October 1, 2013)

Comparative Assessment of Temperature Based ANN and Angstrom Type Models for Predicting Global Solar Radiation

Darlington Ihunanyachukwu Egeonu, Howard Okezie Njoku,
Patrick Nwosa Okolo, and Samuel Ogbonna Enibe

Department of Mechanical Engineering,
University of Nigeria,
Nsukka, Enugu State, Nigeria
darlington.egeonu@unn.edu.ng

Abstract. In this study, temperature based artificial neural network (ANN) models and Angstrom type models for predicting global solar radiation were developed for selected locations in Nigeria. The ANN models were standard multi-layered feed forward, back-propagation neural networks trained with the Levenberg Marquardt algorithm using seventeen years data collected from Nigerian Meteorological Agency (*NIMET*), Abuja, Nigeria and tested with twenty-two years monthly averaged data downloaded from National Aeronautical Space Administration (*NASA*) online database. The network inputs were latitude, longitude, elevation, month, maximum ambient temperature (T_{max}) and minimum ambient temperature (T_{min}), while monthly average global solar radiation was the network output. The Angstrom type empirical models correlated global solar radiation with minimum and maximum ambient temperatures. The performance of the models were evaluated using statistical performance indicators, namely $RMSE$, MBE , R^2 and rank score. The coefficients of determination (R^2) of the ANN models were always greater than 99% for all the selected locations while the highest coefficient of determination for the empirical models was 89%. The temperature-based ANN models were thus shown to deliver superior and more reliable outcomes in comparison with the empirical models.

Keywords: Solar Radiation Prediction, Artificial Neural Network Models, Empirical Models.

1 Introduction

Many countries of the world are confronted with current or impending energy crises. In this scenario, solar energy plays a vital role as a renewable energy because of its unpolluted nature and its reliability in tropical countries like Nigeria. In any energy conversion system, the knowledge of global solar radiation is germane for optimal system design and performance prediction, and the ability to understand and quantify its value accurately is important in the initial

conceptualization and modeling of solar energy dependent processes in the earth-ocean-atmosphere system. In most regions of the world, the efficient application of solar energy seems inevitable because of the widespread sunshine availability. Coupled with problems of access to energy supply, solar energy could serve as an alternative for both thermal and photovoltaic applications.

The solar radiation data of any location can be obtained by measurement or by estimation. The traditional way of obtaining these data is to install pyranometers at the required locations. This is however not easy since the equipment are usually expensive and require high technical knowhow for both installation and maintenance. As a result, in most developing nations, these equipment are located at few places /stations with some being poorly installed and others manned by unskilled personnel thus resulting in poor data recording and near total absence of accurate solar radiation data in some cases [1, 2, 3].

In order to overcome the challenges posed by shortage of solar radiation data, estimation methods based on statistical correlations have been resorted to. Solar radiation is related to many meteorological and geographical parameters some of which are easily measurable. Being the primary driver of all terrestrial processes, it is therefore cost-effective in most instances to develop methods to estimate the global solar radiation using these parameters [4, 5, 6]. Some of these parameters include sunshine duration, air temperature, relative humidity, cloud cover, precipitation etc.

2 Background

Models for predicting global solar radiation on a horizontal surface have been proposed which differ in the number of meteorological parameters used in correlating global solar radiation, accuracy and applicability. These solar radiation prediction models can be categorized into empirical models, radiative transfer models and ANN models. Empirical models are sets of equations that correlate global solar radiation with a few other measurable meteorological parameters [7]. The coefficients of the empirical models are location-dependent thus limiting their application. The radiative transfer models entail complex modelling of solar radiation using geographical and meteorological parameters [7].

ANN models employ artificial intelligence techniques and are data driven. Essentially, ANNs are used to learn the behaviour of a system and subsequently used to simulate and predict this behaviour [8]. They have proved their efficiency and superiority to other kinds of models through their ability to use input parameters which have no specified relationship and still produce reasonably accurate predictions. Apart from modelling solar radiation, ANNs have been used in broad range of applications including pattern recognition and classification [9], function approximation and prediction [10], identification and control, optimization and diagnostics [11].

The single artificial neuron, simple network models and example sigmoid functions are well discussed in the ANN literature (e.g [8, 11, 12, 13]). An ANN is a network consisting of connected artificial neurons, including inputs, outputs and

hidden layers. The input and output layers contain neurons of equal number to the number of input and output parameters, respectively. The number of hidden layers depends on the training algorithm. Basically, the mathematical equation of an artificial neuron can be represented below:

$$y(x) = g \left(\sum_{i=0}^n w_i x_i \right) \quad (1)$$

where $y(x)$ is one output axon, x is a neuron with n input dendrites ($x_0 \dots x_n$), w is a weight and g is an activation sigmoid function. g should be a simple threshold function returning 1 or 0 (see [13]).

Angstrom [14] was the first scientist to suggest a linear relationship between global solar radiation and sunshine duration. The model he proposed has been used in practical applications for many years to estimate the daily, monthly and annual global solar radiation (H) from comparatively simple measurements of sunshine duration (S). Later, Prescott [15] expressed Angstrom equation in a more convenient form as

$$\frac{H}{H_o} = a + b \frac{S}{S_o} \quad (2)$$

where H and S are daily global solar radiation received on a horizontal surface at the ground level and the sunshine duration respectively; H_o and S_o are extraterrestrial radiation and daylength respectively; the parameters a and b are regression coefficients.

Many investigators have used Angstrom type model to predict monthly average global solar radiation using minimum and maximum ambient temperatures as input parameters [16, 17, 18]. Okundamiya and Nzeako [18] proposed a two-parameter-temperature based linear regression model for estimating global solar radiation on a horizontal surface. This model was applied to selected cities in Nigeria namely Lagos, Benin, Nsukka, Abuja, Yola and Katsina. The performance of the developed model was evaluated using statistical performance indicators namely root mean square error, mean bias error and t-statistic. The good performance of the model in the selected cities re-emphasizes the fact that temperature is an important parameter that should be used in solar radiation modeling in Nigeria.

Reddy and Ranjan [11] obtained estimates of solar radiation at locations in India using ANN models and compared the estimates with those of some correlation models. ANNs for estimating monthly mean daily and hourly values of global solar radiation were created. The ANNs were trained and tested using data from 13 stations spread over India. The ANN solar radiation estimates were in good agreement with the actual values, the maximum mean absolute relative deviation of predicted hourly global solar radiation being 4.07%.

Mubiru and Banda [19] used ANN to estimate monthly average daily global solar irradiation on a horizontal surface at four locations in Uganda based on weather station data (sunshine duration, maximum temperature, and cloud cover) and location parameters of (latitude, longitude, and altitude). Results

showed good agreement between the estimated and actual values of global solar radiation. A correlation coefficient of 0.974 was obtained with MBE of $0.059MJ/m^2$ and $RMSE$ of $0.385MJ/m^2$.

In Nigeria, ANN models have also been applied to the prediction of meteorological parameters for selected cities [20, 21, 23, 22]. Using sunshine hours, maximum temperature, relative humidity and cloud cover as network inputs, they developed feed-forward backpropagation neural networks to analyze and predict solar radiation. The results of the ANN model prediction were found to be more accurate than those of Angstrom-type models.

In the papers reviewed, most developed ANN models used all available meteorological and climatological parameters as inputs but not all these input variables are readily available in weather stations in Nigeria. But minimum and maximum ambient temperatures are readily available in weather stations and are easily measurable hence a purely temperature based artificial neural network model which will be handy for predicting monthly average global solar radiation will be studied and compared with temperature-based Angstrom type empirical model.

3 Analysis

3.1 Data

Five locations, representing the range of climatic conditions obtainable in southern part of Nigeria were selected for this study: Calabar (4.95°N , 8.33°E , 246m above sea level (ASL)), Enugu (6.50°N , 7.50°E , 183m ASL), Owerri (5.48°N , 7.00°E , 127m ASL), Port Harcourt (4.67°N , 7.17°E , 83m ASL) and Warri (5.50°N , 5.68°E , 21m ASL). Two sets of data comprising monthly average daily global solar radiation ($MJ/m^2/day$), minimum ambient temperature ($^\circ\text{C}$), maximum ambient temperature ($^\circ\text{C}$), were used in this study. The first set were seventeen years (1991-2007) measured data from ground stations at the locations obtained from the Nigerian Meteorological Agency (NIMET), Federal Ministry of Aviation, Abuja, Nigeria. This data set was used for training the networks. The second data set were twenty-two years monthly averaged data for the locations obtained from the NASA-SSE online database, and used for testing the networks. The data sets were normalized before being used for network development to obtain zero means and unit standard deviations.

3.2 The Empirical Model

An Angstrom-type regression model based on minimum and maximum ambient temperatures was used in this study as follows:

$$\frac{\bar{H}}{H_o} = a + b \frac{\bar{T}_{min}}{\bar{T}_{max}} \quad (3)$$

where \bar{T}_{min} and \bar{T}_{max} are monthly average minimum and maximum ambient temperatures respectively, a and b are empirical coefficients which are obtained by statistical regression.

The monthly mean extraterrestrial radiation on a horizontal surface, \bar{H}_o used in equation 3 is obtained by

$$\bar{H}_o = \frac{1}{D} \sum_{d_m=1}^D H_o(d_m) \quad (4)$$

where d_m is day of the month and D is total number of days in the month. H_o is the total daily extraterrestrial radiation on a horizontal surface given by (see [24])

$$H_o = \frac{24 * 3600}{\pi} I_{sc} \left(1 + 0.033 \cos \frac{360d}{365} \right) \left(\frac{\pi \omega_s}{180} \sin \delta \sin \phi + \cos \delta \cos \phi \sin \omega_s \right) \quad (5)$$

where I_{sc} is the solar constant, taken to be $1357kWh/m^2$, ϕ is the latitude and δ is declination (see [24]):

$$\delta = 23.45^\circ \sin \left(\frac{360(284 + d)}{365} \right) \quad (6)$$

and ω_s is the solar hour angle (see [24]):

$$\omega_s = \cos^{-1}(-\tan \delta \tan \phi) \quad (7)$$

3.3 ANN Model Development

A multilayer feedforward, back-propagation neural network was created and initialized using Matlab ANN tool box. Since it has been established that any network which has 'tansig' transfer function in the hidden layer and 'pureln' transfer function in the output layer can arbitrarily approximate any function [25, 26], this was the configuration selected. The Levenberg Marquardt algorithm was also selected as the training function. The artificial neural network model has month, longitude, latitude, elevation, maximum ambient temperature and minimum ambient temperature as input parameters while monthly average global solar radiation is the network output. The network architecture of the model is shown in Figure 1.

The training process required a set of examples of proper network behavior, network inputs and target outputs. The network was trained for 2000 epochs and a network architecture with two hidden layers of ten neurons in each hidden layer were used. During training which was performed with measured (NIMET) data, and based on Levenberg Marquardt algorithm, the weights and biases of the network were iteratively adjusted to minimize the network performance function (i.e. mean square error).

Next, the performance of the models developed were tested by exposing them to new data. The online monthly averaged daily climatological data obtained from the NASA database [27] was used for testing the networks.

The performance of the developed models were evaluated quantitatively to verify the existence of underlying trends in performance. Statistical measures

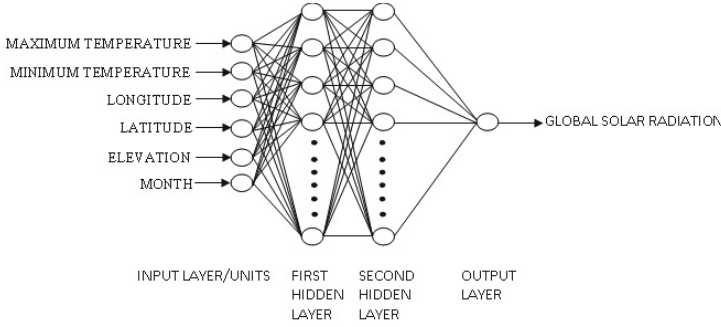


Fig. 1. Network architecture of ANN Model

including the coefficient of determination (R^2), the root mean square error ($RMSE$) and the mean bias error (MBE) were used. (R^2) is the proportion of variability in a data set that is accounted for by a statistical model, where the variability is measured quantitatively as the sum of square deviations. This is represented notationally in equation (8).

$$R^2 = 1 - \frac{\sum_{i=1}^n (\bar{H}_p - \bar{H}_m)^2}{\sum_{i=1}^n (\bar{H}_m)^2} \quad (8)$$

$RMSE$ is a measure of the variation of predicted values around the measured data and it provides information on the short term performance. The lower the $RMSE$, the more accurate is the estimation. $RMSE$ is shown mathematically by equation (9).

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\bar{H}_p - \bar{H}_m)^2} \quad (9)$$

MBE is an indication of the average deviation of the predicted values from the corresponding measured data and can provide information on long term performance of the models; the lower the MBE , the better is the long term model prediction. A positive MBE value indicates the amount of overestimation in the predicted global solar radiation and vice versa. MBE is represented by equation (10). \bar{H}_m is the measured monthly average daily global solar radiation; \bar{H}_p is the predicted monthly average daily global solar radiation; n is the number of observations.

$$MBE = \frac{1}{n} \sum_{i=1}^n (\bar{H}_p - \bar{H}_m) \quad (10)$$

The best performing model was determined using a ranking method proposed by Mubiru *et al.*[28]. Following their procedure, a rank score was computed as the sum of both normalized MBE and the normalized $RMSE$. Both were obtained by dividing with the mean as shown in equation (11). The model with the lowest rank score received the highest ranking.

$$\text{Rank Score} = \frac{\text{abs}(MBE)}{\text{Mean}} + \frac{RMSE}{\text{Mean}} \quad (11)$$

4 Results and Discussion

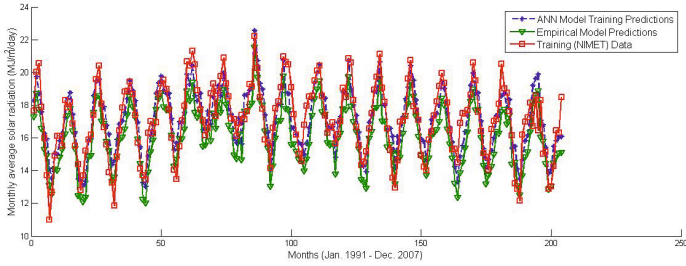
The output of the developed ANN and empirical models implemented using the measured (NIMET) dataset are shown in Figures (2a-3b). Visual observation shows that both models give good agreement with measured data. This is confirmed by the values of the accuracy measures of the predictions of both models.

The $RMSE$, MBE , and R^2 of the predictions of the ANN models, developed (trained) with measured (NIMET) data and testing of the ANN models with (NASA-SSE) data, are listed in Table 1. Very high R^2 values ranging from 0.9950 (for Port Harcourt) to 0.9964 (for Owerri) were recorded, while a very low MBE values ranging from $1.1955e-5$ (for Enugu) to $2.3125e-6$ (for Warri) were obtained for the training predictions. The results are the characteristics of the very high prediction accuracies associated with the ANN models.

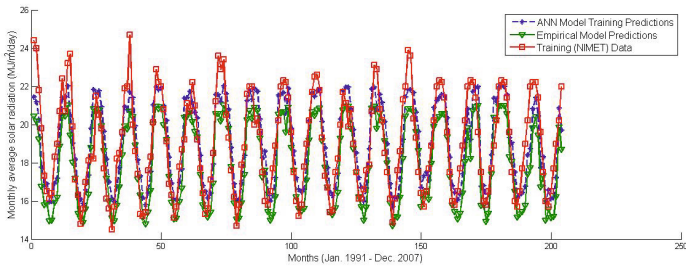
The range of R^2 and MBE for the testing predictions respectively were 0.9965 (for Enugu) to 0.9987 (for Owerri) and 0.0060 (for Owerri) to -0.2674 (for Enugu).

The $RMSE$, MBE and R^2 values of the temperature-based Angstrom-type model, developed using measured (NIMET) data and implemented (tested) on NASA-SSE data, for the locations considered, are presented in Table 2. The R^2 values of the developed models ranged from 0.8167 (for Port Harcourt) to 0.8851 (for Enugu), while the MBE values ranged from 0.5612 (for Warri) to 1.4204 (for Enugu) showing these models to be of acceptable accuracy. When these models were however implemented on the on the NASA-SSE data, the R^2 and MBE respectively, deteriorated markedly to ranges of 0.7212 (for Port Harcourt) to 0.7665 (for Enugu) and 0.9430 (for Owerri) to -2.1167 (for Enugu). This extent of deterioration was not recorded with the ANN models.

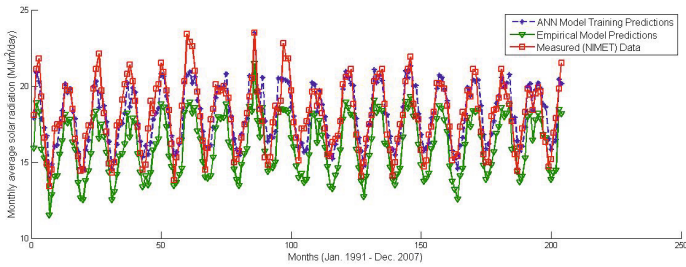
Values of monthly average solar radiation obtained from NASA-SSE, together with ANN and empirical model predicted values are plotted in Figures(3a-4b), for the locations considered. The plots show that ANN models consistently give predictions of closer agreement with the test data. This is confirmed by the rank scores computed by equation 11 for these models, which are presented in Table 3. Since lower rank score translates to better predictions, it is seen that better predictions were obtained with ANN models for all the locations considered. This was so for both the training and test cases – the ANN model were found to give more reliable and accurate predictions of the monthly average global solar radiation than the empirical models.



(a) Calabar

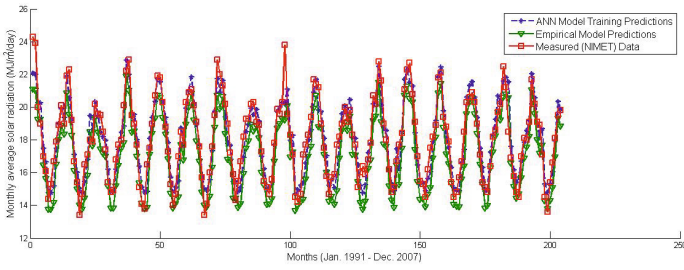


(b) Enugu

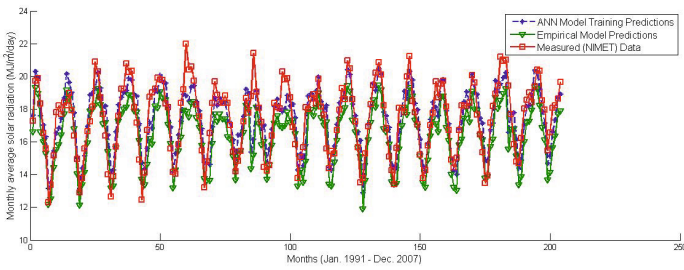


(c) Owerri

Fig. 2. Plots of Predicted and Measured Monthly Average Global Solar Radiation Based Using NIMET Data for (a) Calabar (b) Enugu (c) Owerri (d) Port Harcourt and (e) Warri



(a) Port Harcourt

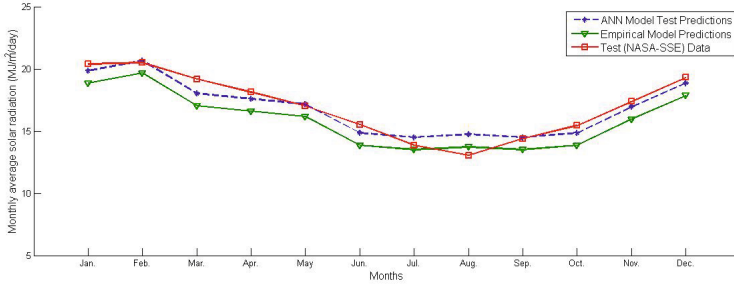


(b) Warri

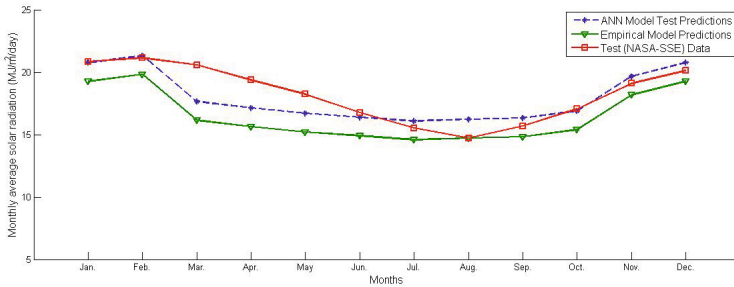
Fig. 2. (continued)

Table 1. Results for the ANN Model

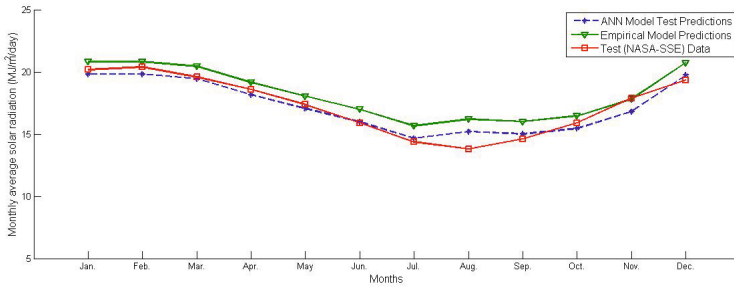
Location	Data Set	RMSE	MBE	R^2
Calabar	Training	1.0935	-1.245e-6	0.9962
	Testing	0.7325	-0.1359	0.9982
Enugu	Training	1.2844	1.1955e-6	0.9951
	Testing	1.1278	-0.2674	0.9965
Owerri	Training	1.0941	6.7099e-4	0.9964
	Testing	0.6391	0.0060	0.9987
Port Harcourt	Training	1.2927	-1.7317e-4	0.9950
	Testing	0.6579	-0.0120	0.9984
Warri	Training	1.1272	2.3125e-6	0.9959
	Testing	0.7910	0.0083	0.9978



(a) Calabar

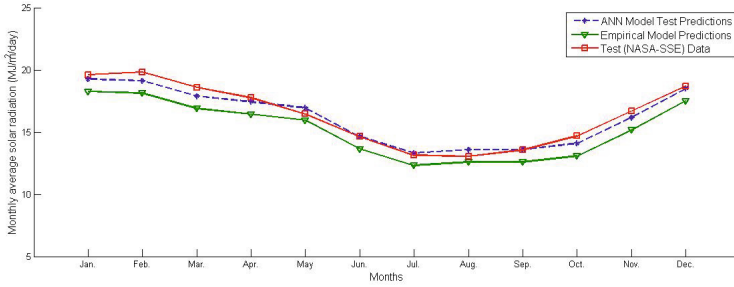


(b) Enugu

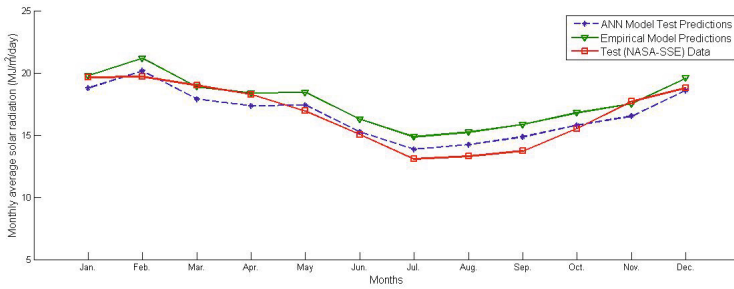


(c) Owerri

Fig. 3. Plots of Predicted and Measured Monthly Average Global Solar Radiation Based Using NASA - SSE Data for (a) Calabar (b) Enugu (c) Owerri (d) Port Harcourt and (e) Warri



(a) Port Harcourt



(b) Warri

Fig. 3. (continued)

Table 2. Results for the Empirical Model

Location	Model	Data Set	RMSE	MBE	R^2
Calabar	1.8304 – 1.8036(T_{min}/T_{max})	NIMET	1.5136	-0.8512	0.8267
		NASA-SSE	2.1134	-1.5214	0.7598
Enugu	1.7417 – 1.7156(T_{min}/T_{max})	NIMET	1.8492	-1.4204	0.8851
		NASA-SSE	1.9997	-2.1167	0.7665
Owerri	1.6719 – 1.6080(T_{min}/T_{max})	NIMET	1.9141	0.5761	0.8493
		NASA-SSE	2.4472	0.9430	0.7317
Port Harcourt	1.6191 – 1.5335(T_{min}/T_{max})	NIMET	1.5911	-0.6117	0.8167
		NASA-SSE	2.1140	-1.1619	0.7212
Warri	1.9203 – 1.9193(T_{min}/T_{max})	NIMET	1.9116	0.5612	0.8508
		NASA-SSE	2.1719	1.0121	0.7408

Table 3. RANK SCORES

Location	Data Set	ANN Model	Emperical Model
Calabar	Training	0.063888	0.1438
	Testing	0.051039	0.2210
Enugu	Training	0.067536	0.1658
	Testing	0.076403	0.2087
Owerri	Training	0.060568	0.1425
	Testing	0.037230	0.1941
Port Harcourt	Training	0.071588	0.1184
	Testing	0.040869	0.1761
Warri	Training	0.064570	0.1470
	Testing	0.047813	0.1892

5 Conclusions

Temperature-based artificial neural network models and empirical models for predicting monthly average global solar radiation were studied in this work. The ANN model has minimum ambient temperature, maximum ambient temperature, month, longitude, latitude, elevation as input parameters while monthly average global solar radiation is the network output. The developed empirical models were Angstrom-type model which correlated global solar radiation with minimum and maximum ambient temperatures. The performance of the models were evaluated based on dimensionless error statistics called rank scores and coefficients of determination. Using the rank score approach suggested by Mubiru *et al.* [28], temperature based ANN model performed better than the temperature based empirical model for the locations considered. This reaffirms the fact that minimum and maximum ambient temperatures are the most important weather parameters for predicting monthly average global solar radiation. It is therefore suggested that temperature based ANN models should be the choice models for predicting monthly average global solar radiation in locations where directly measured data is lacking since maximum and minimum temperature data are fairly easy to obtain.

Acknowledgements. The authors are grateful to Mrs Kelechi Ojide of the Department of Statistics, Federal University Ndufu-Alike Ikwo, Ebonyi State, Nigerian, for graciously providing the NIMET data used in this study.

References

- [1] Ezekwe, C.I., Ezeilo, C.O.: Measured Solar Radiation in a Nigerian Environment Compared with Predicted Data. *Solar Energy* 26, 181–186 (1982)
- [2] Aidan, J., Yadima, A., Ododo, J.C.: Modeling of Un-available Solar Radiation Using Some Climatological Parameters. *Nigerian J Solar Energy* 15(1), 18–26 (2005)

- [3] Augustine, C., Nnabuchi, M.N.: Correlation Between Sunshine Hours and Global Solar Radiation in Warri, Nigeria. *The Pacific J. Sci. Tech.* 10(4), 574–579 (2009)
- [4] Kassem, A., Aboukarima, A., El Ashmawy, N.: Development of Neural Network Model to Estimate Hourly Total and Diffuse Solar Radiation on Horizontal Surface at Alexandria City (Egypt). *Journal of Applied Sciences Research* 5(11), 2006–2016 (2009)
- [5] Falayi, E.O., Adepitan, J.O., Rabi, A.B.: Empirical Models for the Correlation of Global Solar Radiation with Meteorological Data for Iseyin, Nigeria. *International Journal of Physical Sciences* 3(9), 210–216 (2008)
- [6] El-Sebaai, A., Trabea, A.: Estimation of Global Solar Radiation on Horizontal Surfaces Over Egypt. *Egypt. J. Solids* 28(1), 163–175 (2005)
- [7] Donatelli, M., Bellocchi, G., Fontana, F.: Software to Estimate Daily Radiation Data from Commonly Available Meteorological Variables. *Agric. For. Meteorol.* 18(3), 363–367 (2003)
- [8] Kalogirou, S.A.: Artificial Neural Networks in Renewable Energy Systems Application: A Review. *Renewable Sustainable Energy Rev.* 5, 373–401 (2001)
- [9] Knutti, R., Stocker, T.F., Joos, F., Plattner, G.K.: Probabilistic Climate Change Predictions Using Neural Networks. *Climate Dynamics* 21, 257–272 (2003)
- [10] Mohandes, M., Rehman, S., Halawani, T.O.: Estimation of Global Solar Radiation Using Artificial Neural Network. *Renewable Energy* 14(1), 79–84 (1998)
- [11] Reddy, K.S., Ranjan, M.: Solar Resource Estimation Using Artificial Neural Networks and Comparison with other Correlation Models. *Energy Conversion and Management* 44(15), 2519–2530 (2003)
- [12] Haykin, S.: *Neural Networks and Learning Machines*, 3rd edn. Pearson Education, Inc., New Jersey (2009)
- [13] Nissen, S.: Implementation of a Fast Artificial Neural Network. Department of Computer Science, University of Copenhagen, PDF Notes, 29 p. (2003)
- [14] Angstrom, A.: Solar and terrestrial radiation. *Quart. Jour. Roy. Meteorol. Soc.* 50, 121–125 (1924)
- [15] Prescott, J.A.: Evaporation from a water surface in relation to solar radiation. *Trans. R. Soc. Sci. Australia* 64, 114–125 (1940)
- [16] Garcia, J.V.: *Principios Fisicos de la Climatologia*, Ediciones. UNALM (Universidad Nacional Agraria La Molina: Lima, Peru) (1994)
- [17] Hargreaves, G., Samani, Z.: Estimating Potential Evaporation. *Journal of Irrigation and Drainage Engineering* 108, 225–230 (1982)
- [18] Okundamiya, M.S., Nzeako, A.N.: Empirical Model for Estimating Global Solar Radiation on Horizontal Surfaces for Selected Cities in the Six Geopolitical Zones in Nigeria. *Research Journal of Applied Sciences, Engineering and Technology* 2(8), 805–812 (2010)
- [19] Mubiru, J., Banda, E.J.K.B.: Estimation of Global Solar Radiation Using Artificial Neural Network. *Solar Energy* 82(2), 181–187 (2008)
- [20] Chukwu, S.C., Nwachukwu, A.N.: Analysis of Some Meteorological Parameters Using Artificial Neural Network Method for Markurdi, Nigeria. *African Journal of Environmental Science and Technology* 6(3), 182–188 (2012)
- [21] Abdulazeez, M.A.: Artificial Neural Network Estimation of Global Solar Radiation Using Meteorological Parameters in Gusau, Nigeria. *Archiv. Appl. Sci. Res.* 3(2), 586–595 (2011)
- [22] Egeonu, D.I., Njoku, H.O., Enibe, S.O.: Performance of Global Solar Radiation ANN Models with Different Input Parameter Combination. Submitted to: *International Journal of Energy Research* (2014)

- [23] Fadare, D.A.: Modelling of Solar Energy Potential in Nigeria Using an Artificial Neural Network Model. *Applied Energy* (2009), doi:10.1016/j.apenergy.2008.12.005
- [24] Duffie, J.A., Beckman, W.A.: *Solar Engineering of Thermal Processes*. Wiley & Sons Inc., New York (2006)
- [25] Howard, D., Beale, M.: *Neural Network Toolbox for Use with MATLAB, User's Guide, Version 4*, pp. 133–205. The Math Works, Inc. Product (2000)
- [26] Martin, T.H., Howard, B.D.: *Neural Network Design*. PWS Publishing Company, United States of America (1996)
- [27] NASA: <https://eosweb.larc.nasa.gov/> (accessed: December 10, 2012)
- [28] Mubiru, J., Banda, E.J.K.B., Ujanga, F.D., Senyonga, T.: Assessing the Performance of Global Solar Radiation Empirical Formulations in Kampala, Uganda. *Theoretical and Applied Climatology* 87(1-4), 179–181 (2007)

Position Control and Tracking of Ball and Plate System Using Fuzzy Sliding Mode Controller

Andinet Negash and Nagendra P. Singh

Addis Ababa Institute of Technology, Addis Ababa University,
King George IV St., P.O.Box 385, Addis Ababa, Ethiopia
{mail2andinet, npd.singh}@gmail.com
<http://www.aait.edu.et>

Abstract. Sliding mode control, one of the tools available to design robust controllers, is introduced in the outer loop of a double-loop feedback control of the Ball and Plate (B&P) system. Fuzzy Logic is used to attenuate the chattering introduced by Sliding mode control. Genetic algorithm is implemented to determine the parameters of the fuzzy system in an optimal manner. Linear algebraic method is used to design an inner loop angle controller by solving a set of Diophantine equations. The mathematical model of the B&P system, obtained from Euler-Lagrange Equations of Motion, of the proposed controller is evaluated through simulation studies. Simulation results show that the ball could be stabilized anywhere on the plate in 3.5 seconds and it could also track a circular trajectory of 0.4 m radius at 0.8 rad/s in 10 seconds without significant chattering.

Keywords: Fuzzy Sliding Mode Control, Linear Algebraic Method, Ball and Plate control, Fuzzy Logic Control, genetic Algorithm.

1 Introduction

The B&P apparatus is a two dimensional electromechanical device which can be further categorized as a nonlinear, multivariable and unstable system [1]. Furthermore, the system is under-actuated as it possesses more degrees of freedom than the number of available actuators [2]. For a more effective control of the Ball and Plate system, a double feedback loop structure is utilized [1]. The inner loop is designed as an actuator (angular) position controller for the plate inclination while the outer loop is implemented so as to control the balls (linear) position on the plate.

The design of the inner control loop is based on the concept of linear algebraic method. An overall transfer function is chosen that minimizes the Integral of Time Multiplied by Absolute Error (ITAE) and a Two-Parameter configuration is used for the implementation of the compensators that are obtained from the solution of a Diophantine equation.

Because of the existence of uncertainty due to friction, measurement time delays and parameter uncertainties, practical applications require nonlinear control methods to be adopted in the design of the outer control loop [1]. In order to fulfill this requirement, Sliding Mode Control is implemented. The robustness in Sliding Mode Control is achieved as a result of high frequency switching in the control signal which results in chattering effects. Therefore, fuzzy logic is introduced to tune the switching gain based on the distance of the system trajectory from the sliding surface. Genetic algorithm is employed to determine the parameters of the fuzzy membership functions in an optimal manner.

1.1 Mathematical Model

In the modeling of the Ball and Plate, the method of Lagrange shall be used. To begin with, let us define the Lagrangian function as follows [3].

$$L(q_i, \dot{q}_i, t) = T(\dot{q}_i, t) - V(q_i, t) \quad (1)$$

Then, the Euler-Lagrange equations are given by:-

$$\frac{d}{dt} \left(\frac{\partial L(q_i, \dot{q}_i, t)}{\partial \dot{q}_i} \right) - \frac{\partial L(q_i, \dot{q}_i, t)}{\partial q_i} = F_{q_i}, 1 \leq i \leq n \quad (2)$$

where, $T(\dot{q}_i, t)$ and $V(q_i, t)$ are respectively the kinetic and potential energies with respect to the inertial axes, n is the degree of freedom, q_1 through q_n are the generalized co-ordinates, F_{q_1} through F_{q_n} are the generalized forces.

Define a right-handed coordinate as shown in Fig. 1 and then choosing generalized co-ordinates as: $q_1 = x$, $q_2 = y$, $q_3 = \theta_x$ and $q_4 = \theta_y$ [2], and taking small angle approximations, the Lagrangian becomes:

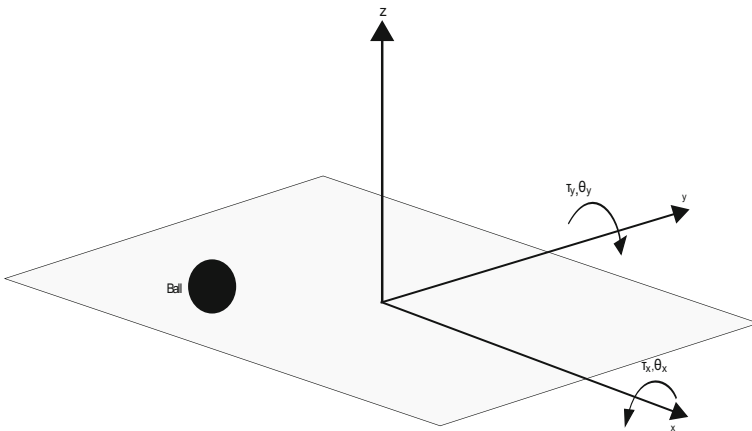


Fig. 1. Schematic diagram of the Ball and Plate system

$$L(q_i, \dot{q}_i) = \frac{1}{2}m \left\{ \dot{x}^2 + \dot{y}^2 + [x\dot{\theta}_x + y\dot{\theta}_y]^2 \right\} + \frac{1}{2} \frac{J_B}{R^2} (\dot{x}^2 + \dot{y}^2) + \frac{1}{2} J_B (\dot{\theta}_x^2 + \dot{\theta}_y^2) + \frac{1}{2} J_{P_x} \dot{\theta}_x^2 + \frac{1}{2} J_{P_y} \dot{\theta}_y^2 + mg(x \sin \theta_x + y \sin \theta_y)$$

Substituting this expression into (1) and evaluating the derivatives and rearranging terms, we obtain four differential equations as follows:-

$$\left(m + \frac{J_B}{R^2} \right) \ddot{x} - mx (\dot{\theta}_x)^2 - my \dot{\theta}_x \dot{\theta}_y = mg \sin \theta_x \quad (3)$$

$$\left(m + \frac{J_B}{R^2} \right) \ddot{y} - my (\dot{\theta}_y)^2 - mx \dot{\theta}_x \dot{\theta}_y = mg \sin \theta_y \quad (4)$$

$$(mx^2 + J_B + J_{P_x}) \ddot{\theta}_x + 2mx\dot{x}\dot{\theta}_x + mxy\ddot{\theta}_y + m(\dot{x}y + x\dot{y})\dot{\theta}_y = \tau_{\theta_x} - mgx \cos \theta_x \quad (5)$$

$$(my^2 + J_B + J_{P_y}) \ddot{\theta}_y + 2my\dot{y}\dot{\theta}_y + mxy\ddot{\theta}_x + m(\dot{x}y + x\dot{y})\dot{\theta}_x = \tau_{\theta_y} - mgy \cos \theta_y \quad (6)$$

In these equations, $m(\text{kg})$ is the mass of the ball; $R(\text{m})$ is the radius of the ball; $x(\text{m})$, $\dot{x}(\text{m/s})$, and $\ddot{x}(\text{m/s}^2)$ are the ball's position, velocity and acceleration respectively along X-axis; $y(\text{m})$, $\dot{y}(\text{m/s})$, and $\ddot{y}(\text{m/s}^2)$ are the ball's position, velocity and acceleration respectively along Y-axis; $\theta_x(\text{rad})$, $\dot{\theta}_x(\text{rad/sec})$ are respectively the plate's deflection angle and angular velocity about X-axis. $\theta_y(\text{rad})$, $\dot{\theta}_y(\text{rad/sec})$ are respectively the plate's deflection angle and angular velocity about Y-axis; $\tau_{\theta_x}(\text{Nm})$ and $\tau_{\theta_y}(\text{Nm})$ are respectively the torque exerted on the plate along X-axis and Y-axis.

We note that the mass and the moment of inertia of the ball are negligible compared to the moment of inertia of the plate [4]. If one assumes that the angle do not change much, $\pm 30^\circ$, the sine function can be replaced by its argument [5]. In equations (3) and (4), the velocities $\dot{\theta}_x$ and $\dot{\theta}_y$ are small and hence have negligible effects when squared or multiplied together [6],[7]. Furthermore, the system inputs are considered as θ_x and θ_y and not the torque moments τ_{θ_x} and τ_{θ_y} . Due to this reason, (5) and (6) could be dropped out in the consecutive study of the system. As a result of these assumptions, linearized, simplified and uncoupled ordinary differential equations are obtained. Substituting for the moment of inertia of the ball, $J_B = \frac{2}{5}(mR^2)$ we get:

$$\frac{7}{5} \ddot{x} = g \theta_x \quad (7)$$

$$\frac{7}{5} \ddot{y} = g \theta_y \quad (8)$$

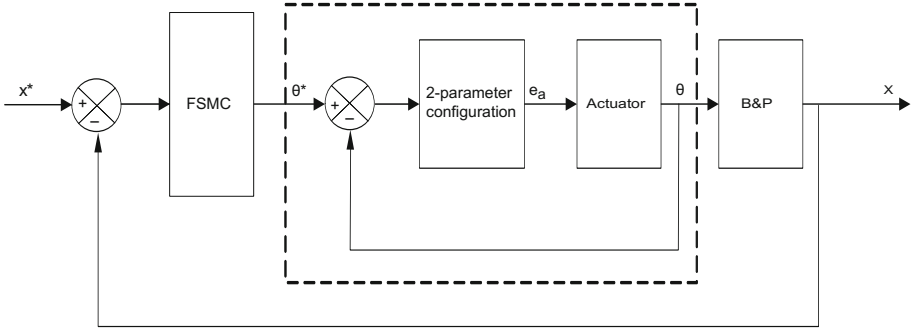


Fig. 2. The proposed double feedback loop for the B & P system

2 Controller Design

One of the most effective control schemes for the Ball and Plate system is the double feedback loop structure, a loop within a loop [1]. Refer to Fig. 2.

2.1 Selecting Parameters of the Actuator

For the B & P application consider the actuator to be a permanent magnet DC motor. The transfer function between θ_L (angular position of the output gear) and e_a (angular position at the output gear) becomes:

$$\frac{\theta_L}{e_a} = \frac{K_t \frac{N_1}{N_2}}{[J_{eq}L_a s^3 + (J_{eq}R_a + D_{eq}L_a) s^2 + (D_{eq}R_a + K_t K_a) s]} \quad (9)$$

where, $J_{eq} = J_a + J_L \left(\frac{N_1}{N_2}\right)^2$ and $D_{eq} = D_a + D_L \left(\frac{N_1}{N_2}\right)^2$

The load torque could be approximated from (5) and the values of the parameters of HUMUSOFT CE151 given in Table 1. The required angular speed of $n_L = 20$ rev/min is estimated from animated videos using VRML. The moment of inertia of the load is obtained from Table 1 as: $J_L = J_P + J_B \approx 0.5 \text{ Kg m}^2$.

Based on the values of τ_L , J_L and n_L , we select a DC motor which will fulfill these requirements- in fact, constraint of maximum and average torques are also

Table 1. Parameters of the Ball and Plate System

No.	Description	Parameter	Value	Unit
1	Mass of the Ball	m	0.11	Kg
2	Radius of the Ball	R	0.02	m
3	Plate Dimension (square)	$l \times w$	1	m^2
4	Moment of Inertia of the Plate	$J_{P_{x,y}}$	0.5	Kg m^2
5	Moment of Inertia of the Ball	J_B	1.76×10^{-5}	Kg m^2
6	Maximum Speed	v	4	mms^{-1}

incorporated. Substituting values for parameters of the selected motor into (9) we obtain:

$$\frac{\theta_L}{e_a} = \frac{4.87 \times 10^{-4}}{s} \tag{10}$$

2.2 Design of the Inner-Loop Controller

Here, the two-parameter configuration is used to implement the overall transfer function $G_o(s)$ and to calculate the desired compensators. Let us choose, $G_o(s)$, the ITAE optimal overall transfer function with zero position error as [8]:

$$G_o(s) = \frac{\omega_o^2}{s^2 + 1.4\omega_0s + \omega_o^2} \tag{11}$$

In order to determine the value of ω_o , we require that a response due to a step input should settle in 0.4 sec. Furthermore, we want the actuating signal $U(s)$ due to a step input not to exceed the rated motor voltage. Through iterative simulation, we find that the value $\omega_o = 20$ rad/sec gives the desired steady state response with an additional gain of 253 being provided using a pre-amplifier to limit the step response. $G_o(s)$ is implemented using the two-parameter feedback configuration shown in Fig. 3, where $L(s)$, $M(s)$ and $A(s)$ are polynomials that describe the compensator and p is a disturbance input.

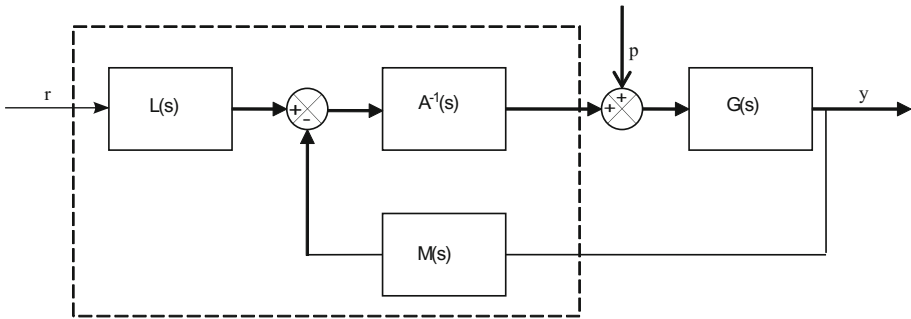


Fig. 3. The two-parameter feedback configuration

Solving the Diophantine equation, we find that $A(s) = 27.877 + s$, $M(s) = 3, 252 + s$ and $L(s) = 3, 252$.

2.3 Design of the Outer-Loop Controller

The sliding surface for a SISO second order plant is designed as:

$$s = \dot{e} + \lambda e \tag{12}$$

For the stability of the system trajectory confined to the sliding surface, λ should be positive. Hence, we choose $\lambda = 1.25$ so that the response due to a unit step input settles in about 4 seconds. We obtain the control signal based on sliding mode control as:

$$u = -b_o^{-1} \{-\ddot{x}_d + \lambda \dot{e} + k \operatorname{sgn}(s)\} \quad (13)$$

Under computer simulation, the value of k that will not result in too large settling time is if it is assigned the value 2.

2.4 Implementation of Fuzzy Logic and Genetic Algorithm

In non-ideal conditions, the system trajectories chatter rather than slide along the sliding surface. As chattering is not a desired phenomenon, we shall fuzzify the relationship between the gain k and the distance s using fuzzy logic. Three fuzzy sets named S, M, L respectively for small, medium and large are chosen to fuzzify the sliding surface. Similarly the switching gain k is fuzzified into three fuzzy sets named S, M, L respectively for small, medium and large.

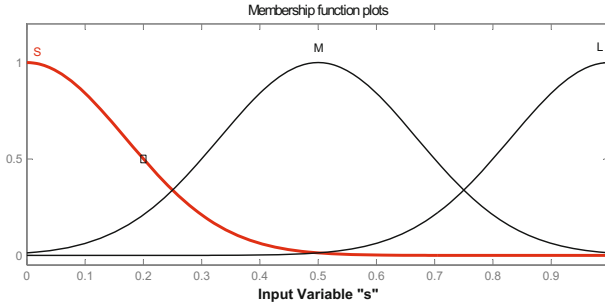


Fig. 4. Approximate shapes of the fuzzy sets for the normalized variable s

In our discussion *shape* shall mean the location of each fuzzy set along the universe of discourse, the width of the membership functions and the amount of overlap with other fuzzy sets. Next, we shall find the optimal shape of the membership functions using genetic algorithm. Since we have 12 variables to be encoded into the chromosome (2 for each fuzzy set, i.e. the mean and the variance) and if we represent each variable using 6 bits, then we have an individual population of length 72 bits. 100 populations are randomly generated each obtained by concatenating a randomly generated bit string of length 72. These are next scaled into the basic ranges for the variables $s(0 - 0.4)$ and $k(0 - 2)$. The design parameters used in the simulation are given in Table 2.

Table 2. Design parameters for the fuzzy sets and genetic algorithm

No.	Design Parameter	Value
1	Population size (Number of individual)	100
2	Membership function	Gaussian : m, σ
3	Number of input fuzzy sets (s)	3
4	Number of output fuzzy sets(k)	3
5	Total number of variable to be determined	2x6=12
6	Length of bits to encode each variable	6
7	Chromosome length	12x6=72
8	Optimization function	minimize (e^2)
9	Reproduction rate	0.15
10	Crossover rate	0.8
11	Mutation rate	0.05
12	Fuzzy inference rule	max

3 Simulation Studies

The 3-D model for the B&P system is designed using V-Realm Builder. The model is used in the study of static and dynamic position tracking of the B&P system.

3.1 Static Position Tracking

The evolution of the sliding surface and the angle θ_x with time as the ball executes free translation to reach and stay at the set-point $(-0.3, 0.4)$ with Fuzzy Sliding Mode Controller (FSMC) is shown in Fig. 5. The variable nature of the control θ_x which depends on the sign of s is evident.

Simulation results with Sliding Mode Controller (SMC) for the same setting is shown in Fig. 6. It is seen that there is too much chattering in the control signal θ_x even after the error dynamics reaches within 2% of its steady state value.

3.2 Trajectory Tracking

In Fig. 7, the frequency is set at 0.2 rad/sec and the ball executes a complete motion along the reference circle in approximately 34 seconds. At an angular frequency 0.8 rad/sec the ball is able to execute a complete motion along the reference circle in 10 seconds but with more tracking error.

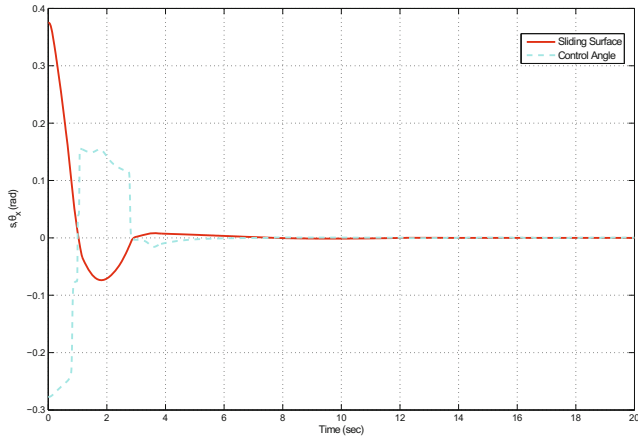


Fig. 5. Plot of s and θ_x with the designed FSM controller

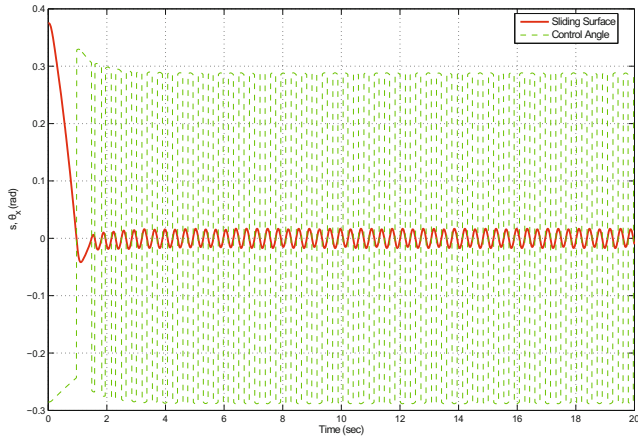


Fig. 6. Plot of s and θ_x with SM controller

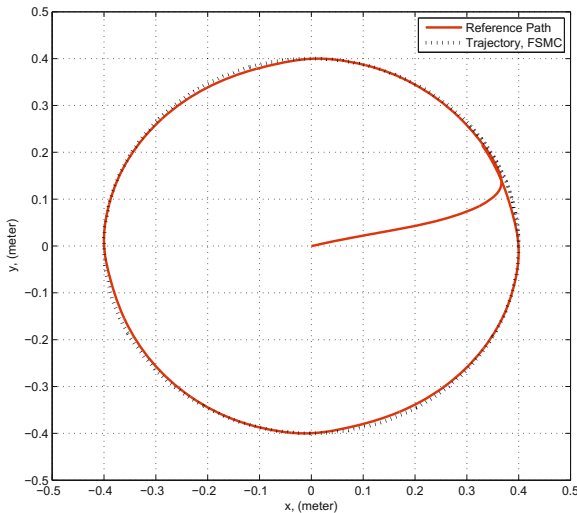


Fig. 7. Circular trajectory tracking performance at 0.2 rad/sec

References

1. Liu, H., Liang, Y.: Trajectory tracking sliding mode control of ball and plate system. In: 2nd International Asia Conference on Informatics in Control, Automation and Robotics, vol. 3, pp. 142–145. IEEE Press, NJ (2010)
2. Wang, H., Tian, Y., Sui, Z., Zhang, X., Ding, C.: Tracking Control of Ball and Plate System with a Double Feedback Loop Structure. In: International Conference on Mechatronics and Automation, ICMA 2007, pp. 1114–1119. IEEE Press (2007)
3. Wellstead, P.E.: Introduction to Physical System Modeling. Academic Press Ltd, London (1979)
4. Fan, X., Zhang, N., Teng, S.: Trajectory planning and tracking of ball and plate system using hierarchical fuzzy control scheme. *Fuzzy Sets and Systems* 144(2), 297–312 (2004)
5. Wang, W.: Control of a Ball and Beam System, http://data.mecheng.adelaide.edu.au/robotics/projects/2007/BallBeam/Wei_Final_Thesis.pdf
6. Knuplei, A., Chowdhur, A., SveEko, R.: Modeling and Control design for the ball and plate system. In: International Conference on Industrial Technology, pp. 216–221. IEEE Press, Shanghai (2009)
7. Liu, D., Tian, Y., Duan, H.: Ball and Plate Control System based on sliding mode control with uncertain items observe compensation. In: International Conference of Intelligent Computing and Intelligent Systems, vol. 2, pp. 216–221. IEEE Press (2003)
8. Chen, C.-T.: Analog and Digital Control System Design: transfer-Function, State-Space and Algebraic Methods. Saunders College Publishing (2006)

9. Duan, H., Tian, Y., Wang, G.: Trajectory Tracking Control of Ball and Plate System Based on Auto-Disturbance Rejection Controller. In: Asian Control Conference, pp. 471–476. IEEE Press (2009)
10. Moarref, M., Saadat, M., Vossoughi, G.: Mechatronic Design and Position Control of a Novel Ball and Plate System. In: Mediterranean Conference of Control and Automation, pp. 1071–1076. IEEE Press (2008)
11. Ashrafzadeh, F., Nowicki, E.P., Boozarjomehry, R., Salmon, J.C.: Optimal synthesis of fuzzy sliding mode controllers [for induction motors]. In: Thirty-First IAS Annual Meeting on Industry Applications Conference, vol. 3, pp. 1741–1745. IEEE Press (1996)
12. Hung, J.Y., Hung, J.C.: Chatter reduction in variable structure control. In: 20th International Conference on Industrial Electronics, Control and Instrumentation, vol. 3, pp. 1914–1918. IEEE Press (1994)
13. Rojko, A., Jezernik, K.: Adaptive Fuzzy Sliding Mode Control of Robot Manipulator. In: The 15th IFAC World Congress, vol. 15, pp. 1077–1077. International Federation of Automatic Control (2002)
14. Hung, J.Y., Gao, W., Hung, J.C.: Variable Structure Control: A Survey. IEEE Transactions on Industrial Electronics 15, 2–22 (1993)
15. Andrews, G., Colasuonno, C., Herrmann, A.: Ball on Plate Balancing System, Rensselaer Polytechnic Institute (2004), <http://cats-fs.rpi.edu/>
16. Slotine, J.-J.E., Li, W.: Nonlinear Applied Control. Prentice-Hall, Inc. (1991)
17. Kung, C.-C., Liao, C.-C.: Fuzzy-sliding mode controller design for tracking control of nonlinear system. In: American Control Conference, vol. 1, pp. 180–184. IEEE (1994)
18. Song, F., Smith, S.M.: A comparison of sliding mode fuzzy controller and fuzzy sliding mode controller. In: 19th International Conference of the North American Fuzzy Information Processing Society, pp. 480–484. IEEE (2000)
19. Ross, T.J.: Fuzzy Logic with Engineering Applications, Second Edition. John Wiley & Sons, Ltd. (2004)
20. Sabanovic, A., Fridman, L.M., Spurgeon, S.K.: Variable Structure Systems: From Principles to Implementation. The Institution of Engineering and Technology (2004)

A Simulation of Project Completion Probability Using Different Probability Distribution Functions

Erimas Tesfaye, Kidist Girma, Eshetie Berhan, and Birhanu Beshah

Addis Ababa University, Addis Ababa institute of Technology, School of Mechanical and Industrial Engineering, Addis Ababa, Ethiopia
ermiastes@gmail.com

Abstract. Estimation of time in project that involves several activities requires expert's knowledge to give an accurate estimation of project duration. In most project, Program Evaluation and Review Technique (PERT) are used to estimate completion time of project with the basic assumption of normality. This assumption is accepted by many scholars without assessing the errors in the result. Thus, this paper examines the validity of this assumption for activities with different probability distribution functions. The effect on the project's completion time is addressed by considering a single probability distribution for the entire activities of a network and mixed probability distribution functions within the network, using a computer based simulation called ARENA in Villa house construction project. The findings show that the project completion time for activities that follow different probability distribution function do not follow normality.

Keywords: Project completion time (PCT), Probability distribution function (PDF), project duration.

1 Introduction

In projects, the variability of time estimates for an activity is assumed to follow beta distribution [1, 5] and the Project Completion Time (PCT) is assumed to follow normal distribution. Different scholars have attempt to prove the validity of this assumption. Elmaghraby [4] in his work on Project Planning and Control by Network Models recommended the use of uniform probability density function for the activity durations. Hamdy A. Taha [8] rely on the central limit theorem to postulate that the completion time can be portrayed using a normal distribution as a function of the cumulative mean and variance of all the activities within the longest path. However, Dong-Eun Lee [3] stated that a normal distribution should not always be assumed in PCT if one wants to get a reliable result. If one assumes that PCTs are normally distributed, PERT may lead to an approximately 10 to 30% more optimistic PCT than when activity durations are generated assuming Triangular, Uniform, Exponential and Weibull functions. Whereas, Looetsma [9] proposed the use of the gamma probability density function. Moreover, another assumption is that all the activities of a network have same Probability Distribution Function (PDF). Reliable estimation of project completion period occupies a central role in the decision making process. In the past,

scheduling methods by relaxing some of the restrictive assumptions of Program Evaluation and Review Technique (PERT) were commonly practiced. However, [3] hypothesized that in practice, it is possible that different activities exhibit different characteristics that necessitates the use of different PDFs

Several studies have been conducted on PCT of a project. Based on classification of the literature [2], these studies can be categorized into three methods: (1) exact methods, (2) approximation methods, and (3) simulation methods. Nowadays, simulation based scheduling methods is a well-accepted technique, which gives project planners flexibility in determining project completion time. Recently, simulation based scheduling methods has been developed for the prediction of project completion period in use of Probability Distribution Functions (PDF) [3]. With the advent of this method, the reliability of estimating project completion period has been enhanced; enabling more realistic decisions to be taken.

While this technique is serving as a practical tool for scheduling project activities, it is important to investigate its reliability as the final result is dependent on the assumption that we take. Commonly the activities of a network are assumed to have a normal distribution. Thirty years of experiences on simulation based scheduling techniques relied on normality assumptions of project completion time (PCT) for the simulation output [6, 7, 10].

It is well recognized that for a certain activity with different estimation time, evaluation of project completion time is commonly carried out using Program Evaluation and Review Technique (PERT) using beta distribution. Recently, the effect of using different probability distribution function on project completion times in simulation based scheduling has been studied. These studies basically assume that, all activities of a network have the same probability distribution function. However, projects with high degree of scatter and complexity, it is possible that different activities in the same network might have different probability distribution function.

Therefore, this paper presents a simulation of PCT that has several activities and when the individual activities undergoes different PDF to validate the normality assumption of PCT. This research attempts to address the validity of normality assumption and conform the level of accuracy on the normality assumption on project completion time.

2 Methodology

The concept of the research were based on two approaches, the first approach focused on the theoretical simulation of a project completion time based on a hypothetical project, with randomly generated data with different PDFs for the identical activities of the network. Intentionally, the data were generated without altering the mean completion just simply by changing the PDFs. Keeping the mean completion time same and differing the PDFs the results of having wide and narrow scattered activity time were analyzed. To back the results obtained using theoretically generated data, from

practical point of view it is important to investigate the PCT of an actual project having wide scatter and different distribution characteristics. Since construction projects are expected to have such characteristics, a project on a villa house was considered. Data that are necessary for simulation were collected from nine construction companies. Every planner estimate the duration of each activity based on the company's trend. Finally, data were analyzed using assumed PDFs and the best fit PDF were plotted using ARENA Vr. 14 input analyzer.

3 Simulation of Project Completion Time (PCT)

The estimation of a project time could be relatively simple as the nature of project is of the same simple nature. However, when doing time estimations for major projects that involve several activities, it requires a little more intellect in order to give an accurate estimation of the time it would take to complete the project.

For projects with high degree of scatter and complexity, it is possible that different activities in the same network might have different PDF. This theoretical simulation analysis tries to shows the characteristics of different PDFs of individual activity time on the total PCT.

The network shown in Figure 1 is introduced through true random number generation to demonstrate this probability distribution functions. It consists of an Activity On-Arc (A-O-A) network with 12 activities. Each activity is assigned a most likely duration based on random numbers generated. Two classes of random number are generated to show both wide and narrow scatter of activity duration. To find the probability distribution function that best describes the distribution of the project completion time, analyses were made by using normal, beta, exponential and Weibull PDFs.

3.1 Simulation of Project Using Randomly Generated Data

According to the traditional PERT technique the probability of a certain project meeting a specific schedule time can be described as follows:

$$Z = \frac{x - \mu}{\sigma}$$

Here, Z is the number of standard deviations of the due date or target date (x) lies from the mean or expected date.

First with the normality assumption the expected completion time is determined using generated activity network time.

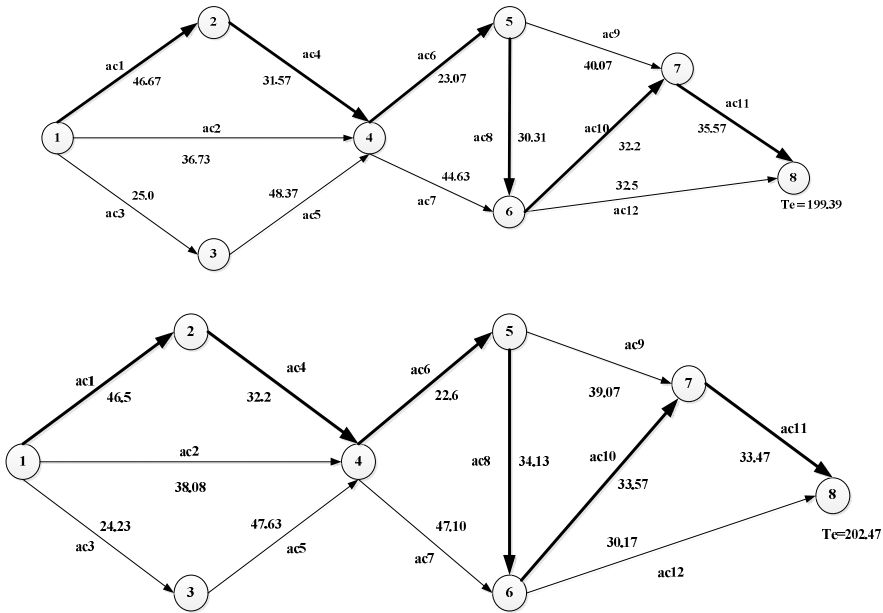


Fig. 1. Activity Network with Narrow and Wide Scatter

The normal expected time (T_e) which is equal to the sum of normal expected times of activities on critical path. i.e. $t_1, t_2, t_3, \dots, t_k$ are the expected times of critical path activities, then

$$T_e = \sum t_i, i = 1, 2, \dots, k$$

Thus, critical path for both narrow and wide scatter activity is the bold line as shown in figure 1 with the expected time T_e of 199.39 and 202.47 days respectively.

The project completion time with 90% confidence is 250.09 days for wide scatter and 270.98 days for the narrow scatter activities as show in Table 1. However, the probability of completion, assuming a normal distribution, may not hold for every activity. Proceeding with examining the behavior of the exponential and beta distributions functions on the project completion time, on the same activity network, the expected time T_e , is 173.39 and 245.33 days for wide scatter activities and 216.56 and 245.33 days for narrow scatter activities respectively.

Table 1. Expected Time for Normal Probability Distribution Function

Activity	Expected Time (T_e)				Probability of Project Completion with 90 % Confidence Interval	
	Wide Scatter		Narrow Scatter		Wide	Narrow
	Mean	Standard Deviation	Mean	Standard Deviation		
ac1	46.5	6.83	46.67	2.43	55.28	62.40
ac2	38.03	6.53	36.73	1.78	-	-
ac3	24.23	6.08	25	2.46	-	-
ac4	32.2	8.86	31.57	2.15	43.59	42.72
ac5	47.63	7.22	48.37	3.19	-	-
ac6	22.6	5.98	23.07	2.56	30.28	32.20
ac7	47.1	5.04	44.63	3.03	-	-
ac8	34.13	3.70	30.17	2.64	38.88	41.41
ac9	39.07	7.33	40.07	2.08	-	-
ac10	33.57	6.76	32.2	2.47	42.26	43.85
ac11	33.47	4.93	35.57	2.69	39.81	48.40
ac12	30.17	5.99	32.5	2.29	-	-
Expected Time (T_e)					250.09	270.98

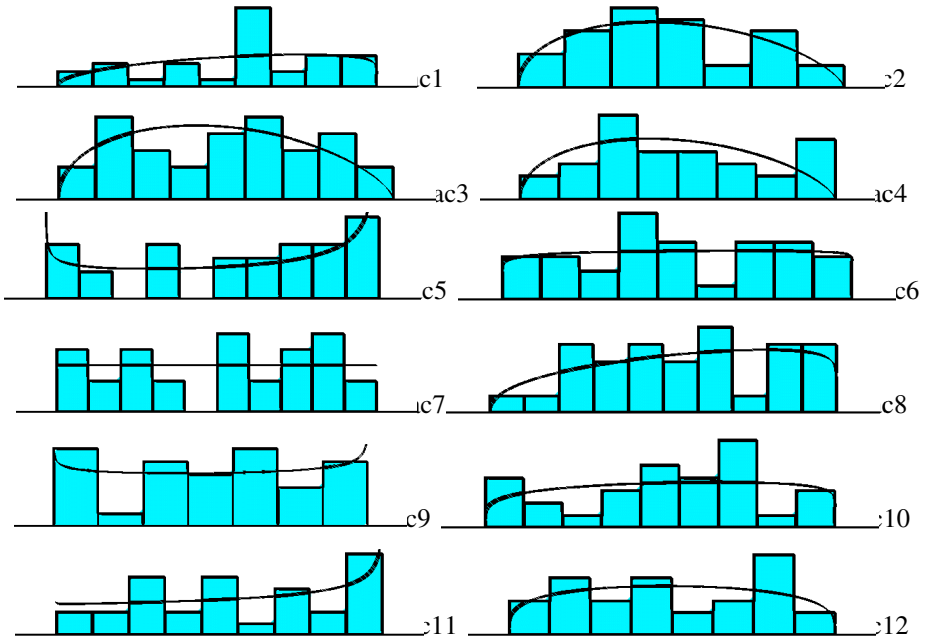


Fig. 2. Narrow Scattered Best Fit Chart

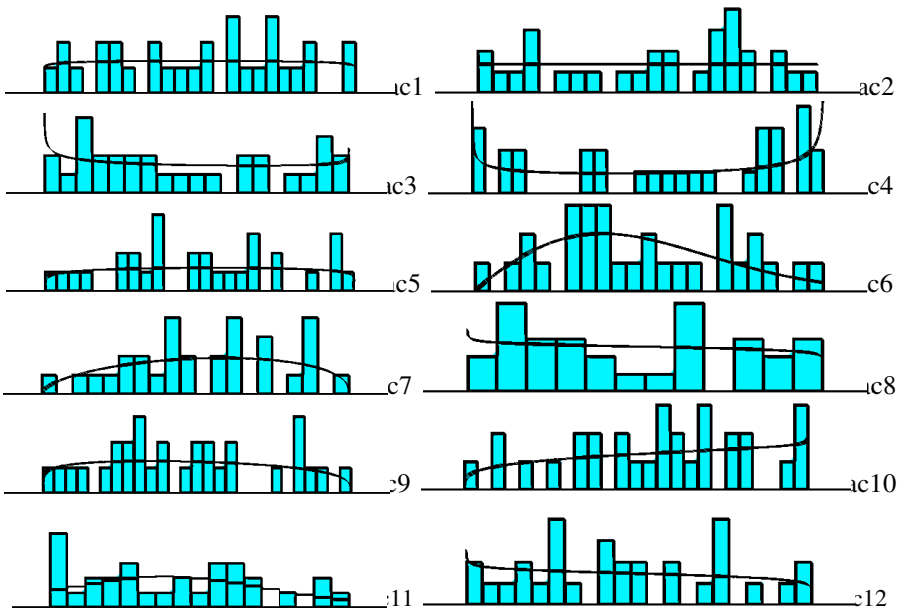


Fig. 3. Narrow Scattered Activity Best Fit Chart

Consequently, it is necessary to see the best fit PDFs for each activity to check the normality assumption. The best fit PDFs is developed by using ARENA input analyzer to find the completion time distribution. Furthermore, the program simulates the project completion time based on the best fit PDFs expressions.

As shown in Figure 2 and 3 the best fit PDFs for both wide and narrow scattered activity network varied from the normality assumption where most activities show beta PDF and there are some with Weibull, Poisson and uniform PDF which makes the individual activity's in the network having different PDF's. Therefore, the next step is to check the normality assumption of the total PCT of the network.

By simulating this best fit probability distribution expression on each activity the project completion time was computed. The computed results for different PDF's are shown in Figure 4 and 5.

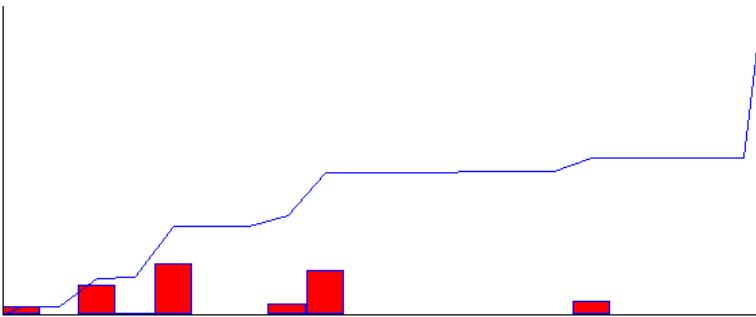


Fig. 4. Completion Time Distribution for narrow Scatter

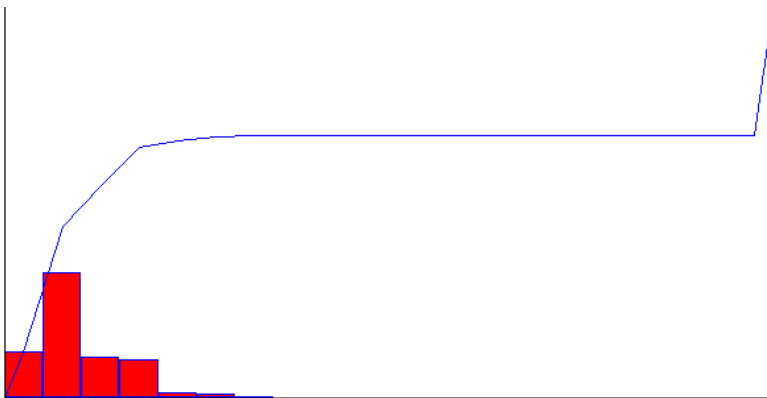


Fig. 5. Completion Time Distribution for Wide Scatter

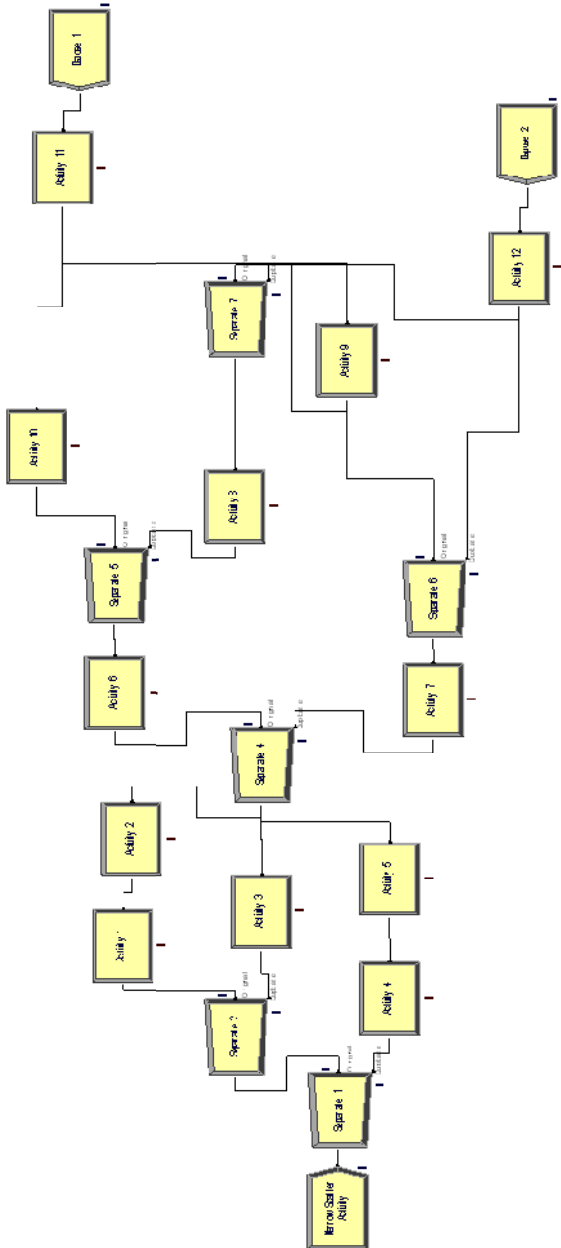


Fig. 6. Completion Time Distribution Simulation Network for Wide and Narrow Scatter

3.2 Simulation of Project Using the Actual Data

The probability of the project completion time assuming the normality, beta, exponential and best fit PDF were tried to be simulated on the theoretical simulation part using a randomly generated activity network. As it is explained above, there is a deviation from the normality assumption where each activity follows different PDF. Therefore, this case study tries to address ascertains that the normality assumption, which frequently used in construction simulation studies.

This paper considers a villa house construction projects from nine construction companies. For the purpose of simulation, only excavation and earth work, concrete work for sub and super structure, masonry and block work activities were considered. As shown in figure 7, the network has 21 activities and the respective mean expected time for each activity is shown in Table 2.

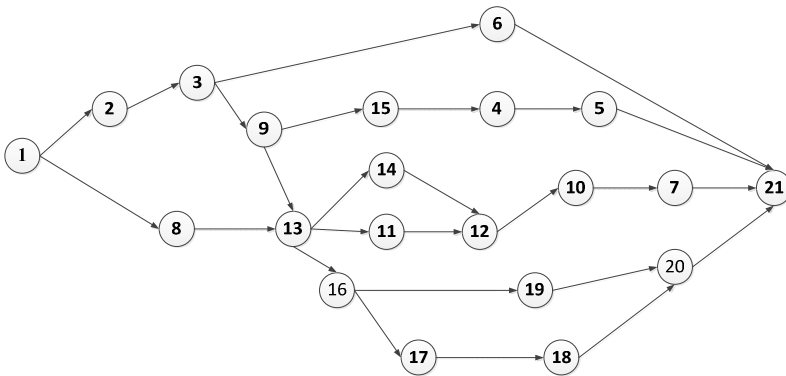


Fig. 7. Network diagram for the villa house project

Table 2. Mean Expected Time for villa house project

Activity	Mean Expect time (Days)	Activity	Mean Expect time (Days)	Activity	Mean Expect time (Days)
1	3.2	8	1.6	15	6.6
2	5.1	9	1.9	16	1.7
3	5.4	10	1.6	17	2.4
4	2.8	11	2.9	18	2.7
5	3.0	12	2.7	19	3.6
6	3.5	13	4.6	20	6.6
7	4.1	14	1.0	21	8.6

Analyzing the above network with Normal, Beta, Weibull and Exponential PDF for critical activity 1, 2, 3, 9, 13, 16, 17, 18, 20 and 21, Table 3 indicates the project completion time for each activity and best fit PDFs is presented Figure 8. The result indicated that each activities uses different PDF and the PCTs varies when there is changes in PDFs. On the other hand, the normal distribution appears to best fitting for activities only for the same PDF throughout the whole activity time.

Table 3. The Project Completion Time of the villa house project using different PDFs

	Probability Distribution Function			
	Normal	Beta	Weibull	Exponential
Project Completion Time [days]	89.71	72.08	77.63	97.09

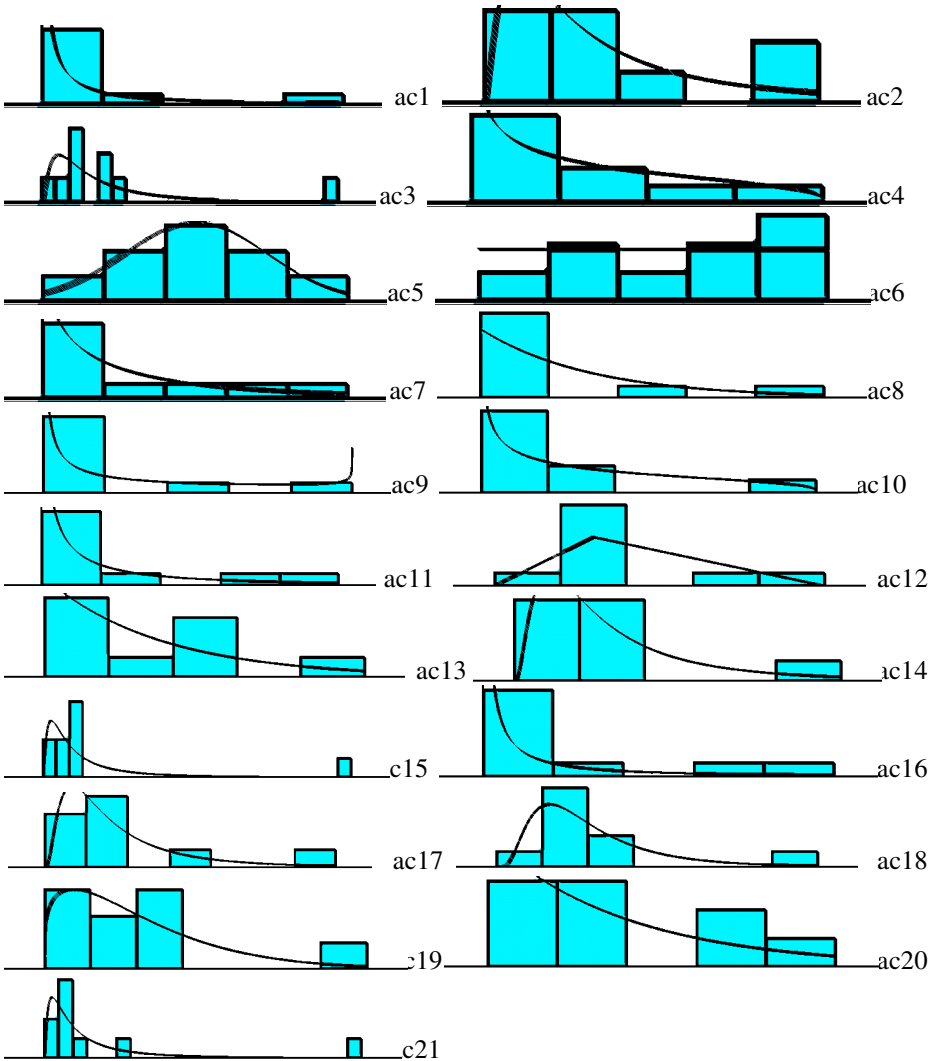


Fig. 8. Best Fit Probability Distribution Function of the 21 activities

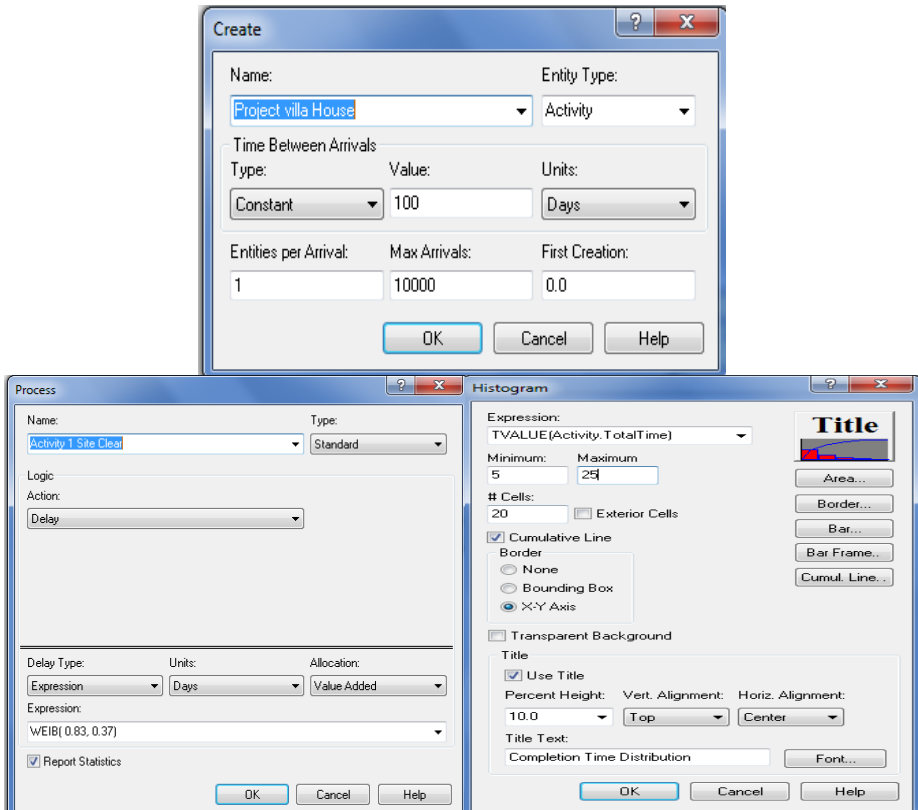
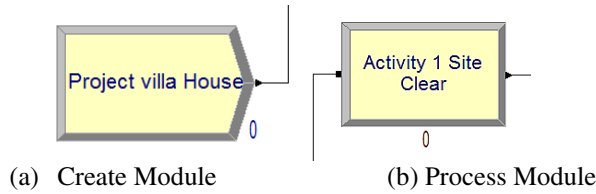


Fig. 9. Snap shot of the ARENA Simulation

The procedure to model the project network using ARENA ver. 14 simulation is demonstrated in figure 9. Create module is modified by setting constant in the Type field and at least 100 in the Value field. The Max Arrivals field sets to be 10,000 to assure the accuracy of the histogram for the completion time distribution, (William J. Cosgrove, 2008)[11]. On the process module, each activity was labeled and the action set to be delay. Delay is an expression based on the best fit PDF with a unit per days.

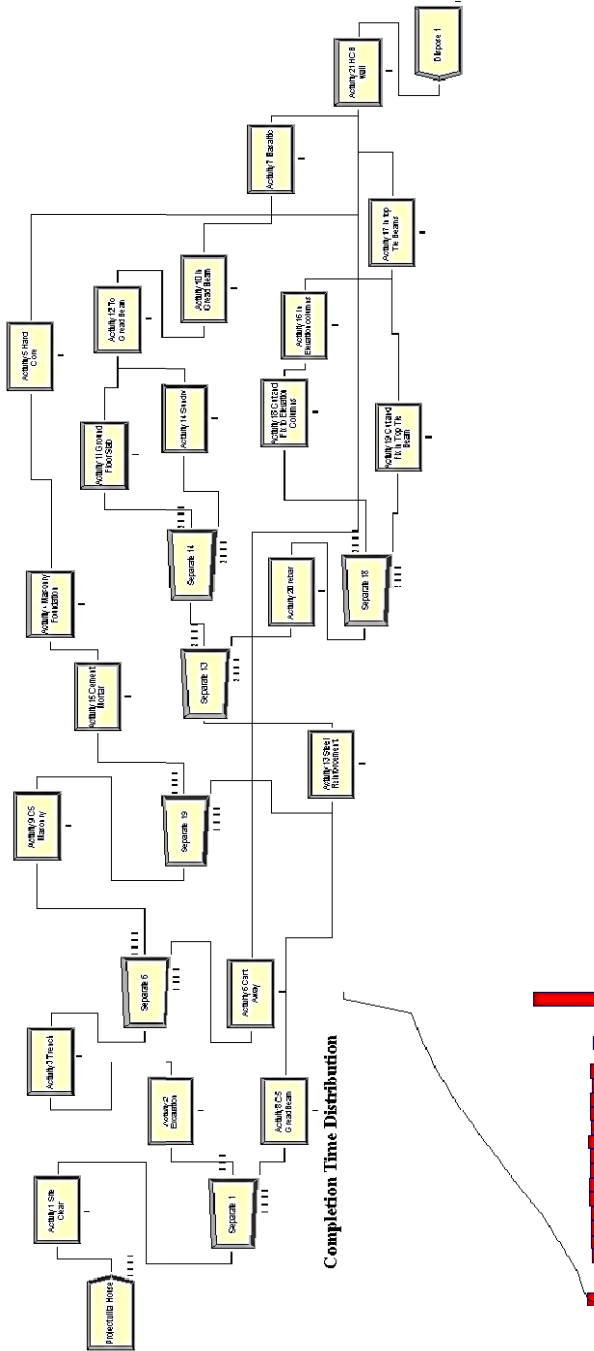


Fig. 10. Completion Time Distribution Simulation Network for The villa house project

As shown in the figure 9, dispose and Separate module groups are used with the default settings because these modules have no impact on Arena network. The last task prior to running the simulation is to construct the completion time histogram. In the Expression field the values used to plot the histogram is based on the research [11]. For the fields minimum, maximum, and # Cells, enter 5, 25, and 20. The first two fields represent estimates of the range of the completion time histogram, and the last field gives the number of time intervals on the histogram. Moreover, all the project activities are represented in Arena by a process module as shown in the figure 10.

Consequently the results obtained from Arena Simulation which is the PCT is plotted using a histogram. The simulation is repeated 10,000 times and it depicts that the normality assumptions are invalid for different activity PDF.

4 Conclusion

In this work two methods of determining completion time distribution were demonstrated. The first based on a hypothetically generated data and the second on an actual construction project. The results of both simulations indicate that different activities display different probability distribution function on a same activity network, thereby leading to different estimation of project completion time. Thus, the current study has demonstrated that, realistic prediction of project completion time shall be carried out based on best fit PDFs rather than simply assuming normality.

References

1. AbouRizk, S.M., Halpin, D.W.: Statistical properties of construction duration data. *J. Constr. Eng. Manage.* 111(4), 525–544 (1992)
2. Adlakha, V., Kulkarni, V.G.: A classified bibliography of research on stochastic PERT networks. *Infor.* 27(3), 272–296 (1989)
3. Dong-Eun Lee, D.A.-B.: The Probability Distribution of Project Completion Times in Simulation-based Scheduling. *KSCE Journal of Civil Engineering* 17(4), 638–645 (2013)
4. Elmaghraby, S.E.: *Activity Networks: Project Planning and Control by Network Models.* Wiley, New York (1977)
5. Fente, J., Knutson, K., Schexnayder, C.: Defining a beta distribution function for construction simulation. In: *Proc. 1999 Winter Simulation Conference*, pp. 1010–1015. IEEE, Piscataway (1999)
6. Ang, A.H.-S., Tang, W.H.: *Probability concepts in engineering planning and design: Volume I - basic principles.* Wiley, New York (1975)
7. Halpin, D.W., Riggs, L.S.: *Planning and analysis of construction operations.* Wiley, New York (1992)
8. Taha, H.A.: *Operations Research: An Introduction*, 7th edn. Prentice Hall (2010)
9. Loostma, F.A.: Network Planning with Stochastic Activity Durations: An Evaluation of PERT. *Statistica Neerlandica* 20, 43–69 (1966)
10. Lu, M., AbouRizk, S.M.: Simplified CPM/PERT simulation model. *Journal of Construction Engineering and Management* 126(4), 219–226 (2000)
11. William, J., Cosgrove, W.: Simplifying PERT Network Simulation with ARENA. *California Journal of Operations Management* 6(1), 61–68 (2008)

Dynamic Simulation of T-Track: Under Moving Loads

Mequanent Mulugeta

Road and Railway Engineering
School of Civil & Environmental Engineering
Addis Ababa Institute of Technology
Addis Ababa University
mequ2me@gmail.com

Abstract. Though, there exists advanced ballastless track forms around the globe several thousand kilometers of 21st century railway lines are under construction with the conventional ballasted tracks. The root reasons mentioned here are huge initial investment and the need of advanced technology for ballastless tracks. Australian and South African experiences is a good example that T-Track economical alternative neither expensive as slab tracks nor problematic as traditional ballasted tracks. The spatial dynamic response of T-track (a.k.a H-Track, Tubular Modular, T-sleeper) has been analyzed by developing continuous double beam model (rail and RC beam) and continuum subgrade for dynamic analysis using a finite element modeling software - ABAQUS. The primary kinematic responses of the track are evaluated based on time and frequency domains under different modeling parameters and train speed. The effect of vehicle speed on the subgrade structure has been investigated under single and an array of moving loads and a comparison is made with ballasted track responses from literature and previous studies. Based on the T-track dynamic analysis, results are characterized by less subgrade deformation under single and multiple moving loads, higher natural frequency ranges and less vibration modes (due to being continuously supported), less vertical rail stress etc. These outputs ultimately indicate that T-track can better be applied and optimized for wide range of axle load and train speed. Specifically, T-Track can suit the climate and topography conditions of Eastern Ethiopian where Great Rift Valley and deserts are dominant.

Keywords: T-Track, Double beam model, Dynamic Analysis, ABAQUS, Spatial responses, Frequency domain, Time domain.

1 Introduction

The loads imposed on the track structure can be form of repeated vertical, lateral and longitudinal forces resulting from traffic and changing temperatures. [16] There are a number of track types and forms (ballasted, modified ballasted track types, ballastless track forms and others developed to satisfy the heavy axle load and high speed demands in the world. T-track (H-Track) is ballastless and sleeperless track structure consists of continuously supported rails, twin parallel reinforced concrete beams and

transverse steel gauge members to connect the beams and railpads. The main advantage of it can be optimized for any axle load and speed and it is found to be sound alternative for track construction in some hostile environments such as deserts and low laying areas. [5], [6]

The innovative ballastless structures remain feasible because of the advantages like resistance to high dynamic force and less maintenance requirements. However, Ballastless track structures lack the resiliency property provided by the ballast bed as in conventional track structures. Resiliency requirement in ballastless tracks is met by inserting resilient elements such as rail pads, placed discreetly or continuously between the rail and the supporting structure. Resilient materials play an important role in ballastless track structures by damping of dynamic impact forces and vibrations caused by the movement of rolling stock.[15] Being continuously supported, T-Track resembles to ladder tracks where the support mechanism is similar. [14]

2 T-Track (Tubular Modular)

T-Track (TMT) is a ballastless track structure developed in South Africa. The track structure consists of longitudinal reinforced concrete beams, which are supported on an engineered foundation and held in position by galvanized steel gauge bars. Fastening systems that are fixed to steel gussets and stirrups, which encircle the concrete beam, hold the rail in place. The track modules are precast off site in lengths of 6m and assembled on site.^{[6], [10]} The gauge bars are not embedded in the beams rather encircle the beams and fixed with steel gussets and stirrups.

2.1 Components of T-Track

a. Gauge Bars

Galvanized Steel gauge bars are used to maintain the gauge/spacing between the rails and connect the concrete beams at a spacing of approximately 3m (depending of specific situations). Stress analysis in gauge bars at tangent, curve and transition section of the track showed that the top of the gauge bar experienced tension regardless of its position in the track. In the transition zone (the largest stress occurred), the outside of the gauge bar was in pure tension, and in the circular curve in pure compression. However, on the straight portion of the track, the outside of the gauge bar experienced, firstly, compression as the train wheel neared the gauge bar, tension as the wheel reached the gauge bar and then compression as the train wheel moved away. [12], [13]

b. Reinforced Concrete Beams

The twin longitudinal RC-beams support the rails continuously and distribute wheel loads down to the formation layers. The beams are precast off site in a production plant at a standard length of 6m module.^[6] The beams can be designed for different load axle and speed. In tubular modular track, the beams have tube holes – so as called Tubular Modular.

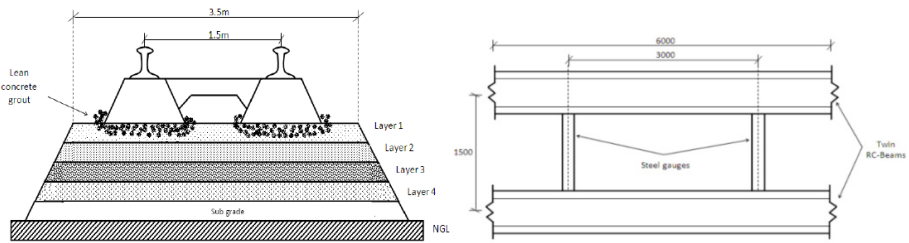


Fig. 1. Typical Cross Section and Plan View of T-Track Structure (mm)

c. Formation Layer

The formation layer is very critical layer responsible to support and transfer the load to natural ground and provide some resiliency to the track. In the similar manner to ballasted track, low frequency vibration are expected be damped by the formation layer. The material specification is required to be high classes of subgrade material such as K30 class. The South African experience suggests that the formation layer is according to the specification of Highway subgrade material specification. Usually multiple composite layers of different property used for formation construction.

2.2 Advantages of T-Track (TMT)

The ballastless, sleeperless and modular nature of the track makes it suitable for both wet and desert conditions where ballast degradation is problematic, also in mines where transportation of the track components is simplified. [6], [10] Extensive testing of Tubular Modular Track by independent experts have shown that rail stresses are reduced by up to 75% compared to ballasted track systems. The continuously supported rail over concrete beams make possible small stress and heavier axle load are applied. Example: 36 tone axle load in Australia and 30 tone in South Africa using 48kg/m rail T-track are constructed. In the technical study, significant savings are also being realized in the components construction costs such as 40% earthworks savings and narrower formation width. T-Track, being ballastless, requires minimal and none plant-intensive maintenance below the rails (60% savings on maintenance as compared to ballasted track). [5, 10, 11, 1] The cost comparison for the earthwork quantity for cut/fill sections between ballasted and T-track shows a remarkable profit for similar conditions; 100% ballast volume; 41% on layer works and 29% on land saved for cut sections; 33% on layer works and 10% on fill saved for fill section of the earthwork. [6]

Table 1. Technical comparison Between T-Track & Ballasted track [6]

<i>T-Track</i>	<i>Ballasted Track</i>
<i>Fully integrated design</i>	<i>Empirical design</i>
<i>Ballastless, sleeperless</i>	<i>Ballasted with sleepers</i>
<i>Continuous support</i>	<i>Discrete support</i>
<i>Reduced rail stress (lighter rail section used)</i>	<i>Higher rail stress (large rail profiles required)</i>
<i>Continuous resilient pad (greater dynamic load disruption)</i>	<i>Resilient supplied by ballast (compromised by contamination of ballast)</i>
<i>Absolute fixed geometry</i>	<i>Unstable geometry</i>
<i>Reduced rolling stock wear, increase comfort</i>	<i>Wear dependent on sleeper and ballast dynamics</i>
<i>Stable and increase passenger comfort</i>	<i>No fixed dynamic geometry, subjected to stability of ballast</i>
<i>Design for desert conditions, drainage friendly</i>	<i>Unsuitable for desert conditions, drainage problematic</i>
<i>Minimal maintenance required</i>	<i>Standard complex maintenance required</i>

3 Track Modelling

For structural analysis of the track, a number of computer packages are available. Especially programs based on finite element method can perform very detailed analysis of displacements, stresses and strains of track components. However, such a modeling of a track structure requires a vast amount of elements, especially under loading condition corresponding to a moving train or under other conditions causing wave propagation and irregularities. ^[2] The significance of dimension in track modeling is up to the analysis accuracy level required in specific track modeling. The complexity and analysis cost also increase as the dimension of the track model increase and the more accurate result will be obtained. Because track structure is a dynamic system, realistic models that comprise inertia and mass have to be considered to fully understand the system performance and response.

3.1 Modelling Parameters

The modeling and analysis parameters used in the paper are Ethiopian railway project data. The model developed in the analysis use the following parameters and values.

<i>Parameters</i>	<i>Symbol</i>	<i>Value</i>
<i>Axle load</i>	<i>P</i>	25 tones
<i>Track Gauge</i>	<i>Standard gauge</i>	1435 mm
<i>Rail</i>	<i>Standard rail T50</i>	50 kg/m
	<i>Elastic Modulus, E</i>	210 GPa
	<i>Poisson Ratio, v</i>	0.25
	<i>Density, ρ</i>	7800 kg/m ³
<i>Railpad</i>	<i>Stiffness, K_p</i>	50 -500 MN/m
	<i>Damping, C_p</i>	75000 N.s/m
<i>RC concrete</i>	<i>Elastic Modulus</i>	300 KPa
<i>subgrade</i>	<i>Elastic Modulus, E_f</i>	50- 120 MPa
	<i>Stiffness, K_f</i>	170 MPa
	<i>Damping, C_f</i>	31 - 100 KN.s/m
	<i>Density, ρ_b</i>	1700-2000 kg/m ³
<i>Formation layer</i>	<i>Formation Stiffness, K_f</i>	170 MPa
	<i>Formation Dampng, C_f</i>	31 KN.s/m

3.2 Simplifications and Assumptions

Here, the model assumes a simplified system, which models the vehicle system’s influence as a mass applied on the track system. The following assumptions made to develop the required models of the track system for analysis.

- The track is symmetry about the center line of the track (middle line)
- The loadings are transmitted/distributed on the two rails symmetrically
- The contact between the rail and wheels assumed as Hertzian contact theory
- Only vertical dynamic movement is studied
- Parameters of the track system are constant and linear in geometry

3.3 Double Layer Beam Models

When the track structure is schematized as a double beam, as shown in *figure 2*, the top beam represents the rail and the bottom beam models the concrete beams. The upper beam has a distributed mass m and a bending stiffness $E_r I_r$. The RC beams are represented as a continuous beam with a mass m_2 and with some non-zero bending stiffness $E_b I_b$ (zero for case of sleeper model in ballasted track).^[4] The rail pads with a stiffness K_p , and a damping C_p and the formation layer can be modeled as elastic foundation and/or as solid continuum media.

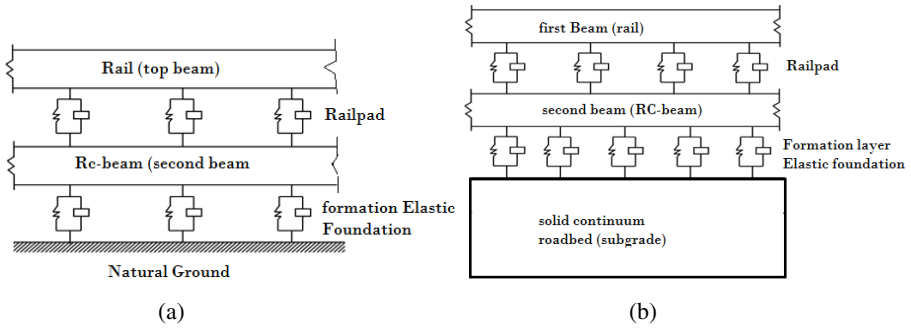


Fig. 2. (a) Simplified Model of T-Track on Elastic foundation Assumption (b) Continuum solid foundation/formation layer

3.4 Modeling of T-Track Structure

The model of the track is an elastically supported continuous double/two-layer beam model consisting of rails, longitudinal beams/sleepers, and formation layer. This model is relatively resembles to the actual cross section and conditions of track structure parts.

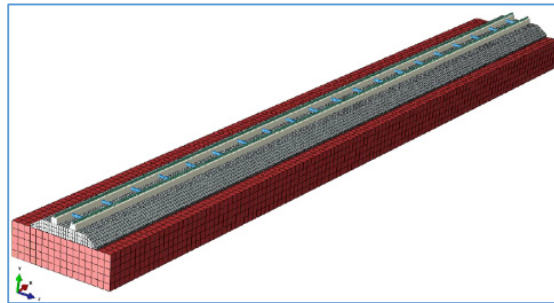


Fig. 3. 3D space Double Beam Model (standard T-Track Module) in ABAQUS

3.5 General Static Response of T-Track

The static load analysis give a general background for simple and direct response of the track. The global beam deflection (rail deflection) and response of the most upper components of the track can be identified. In track static analysis, the beam (rail) stress and deflection are the primary response quantities that can be easily identified.

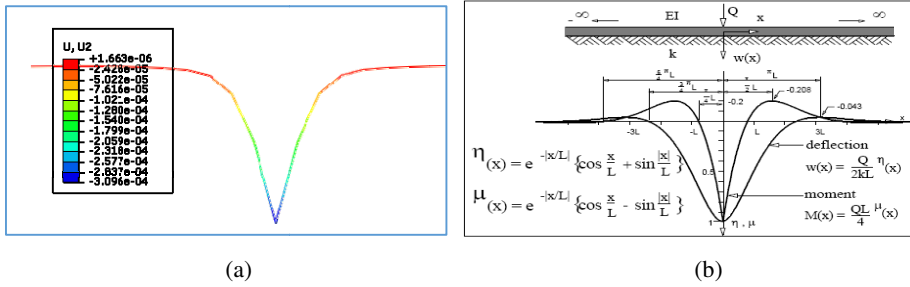


Fig. 4. Static Deflection of Rail (Beam on elastic foundation)a) T-Track b) Ballasted Track

3.6 Steady State Dynamic Analysis

The important response request for beam models were dynamic displacement of the rail (beam) but the velocity and acceleration response of the track should also be the concern. These responses can be quantified or their effect being analyzed with respect to the attenuation capacity of the spring elements connecting the double layer beams. The steady state moving loads considered to model the track and corresponding responses has been evaluated. In addition, railpad stiffness clearly affect the response of the rail and the track deflection. At higher railpad stiffness, the response of the track also improved significantly. The stress and displacement of rails significantly affected by the railpad. During isolated vibration of the beam (rail), the impact of rail on RC-beam is damped out by continuous railpad. That means the rail stress and deflection reduces.

3.7 Natural Frequency Analysis

Natural frequency is a resonance frequency in which the structure tends to vibration in all unrestrained degree of freedom. It is the largest resonance frequency of the structure. A convergence study done to come up with a fully damped resonance frequency. One way of doing this is analyze the system natural frequency with that of the stiffness values. It is very clear that the railpad stiffness increases the response of the track. In the frequency study, the mode of vibration of the system under natural condition, configuration of the connections and global deformations for each component can be identified from the modes. The dynamic response of track is usually associated with the lower modes. However, enough modes should be extracted to provide a good representation of the dynamic response of the track structure. [3], [9] As shown in the figure below the number of modes to be extracted does have a little difference in the vibration of the track. The rail start vibrating in the vertical direction for the 100 modes and to downward direction for the 1000 modes of vibration the rail in the lower frequency range. It seems sensible that the dynamic responses of the track structure usually associated with the lower modes of vibration (below 30Hz).

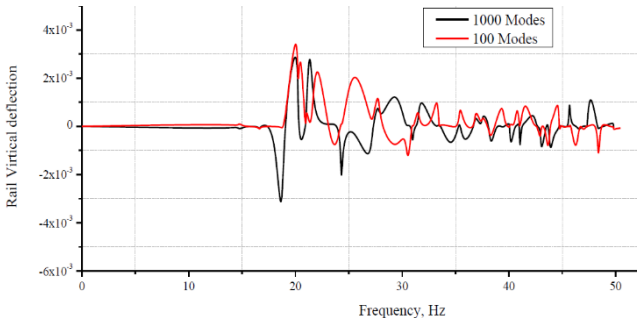


Fig. 5. The effect of Eigen frequency modes on rail vertical deflection for 100 and 1000 modes

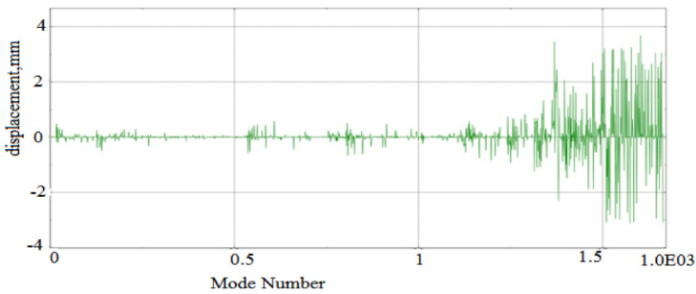


Fig. 6. A single rail Node vertical deflection with in all modes of vibration (left)

3.8 Vibration Modes and Eigen Analysis

Vibration and eigenvalue analysis is useful to analyze the whole structure and global modes of vibration such as the rail vibration can be observed. The simulation view shows which component of the track has dominates the most vibration modes. This is also very helpful to compare field measured data, which usually not clear which part of the track the higher vibration belongs to. At zero frequency, the track structure vibrates in longitudinal and a slight vertical direction.

The railway track structure starts vibrating in in the longitudinal direction and then to lateral direction when the frequency starts to increase from zero. This is true that the primary (most) vibrations of track structure usually found in the first few modes of vibration. The lowest modes of vibration dominated by the lateral vibrations and the modes are mainly the vibration ways of rail on concrete beams. At higher frequencies vertical vibration dominate most excitations. In addition, the gauge bars start vibrating in the middle (about 150Hz) frequencies into vertical and longitudinal directions.

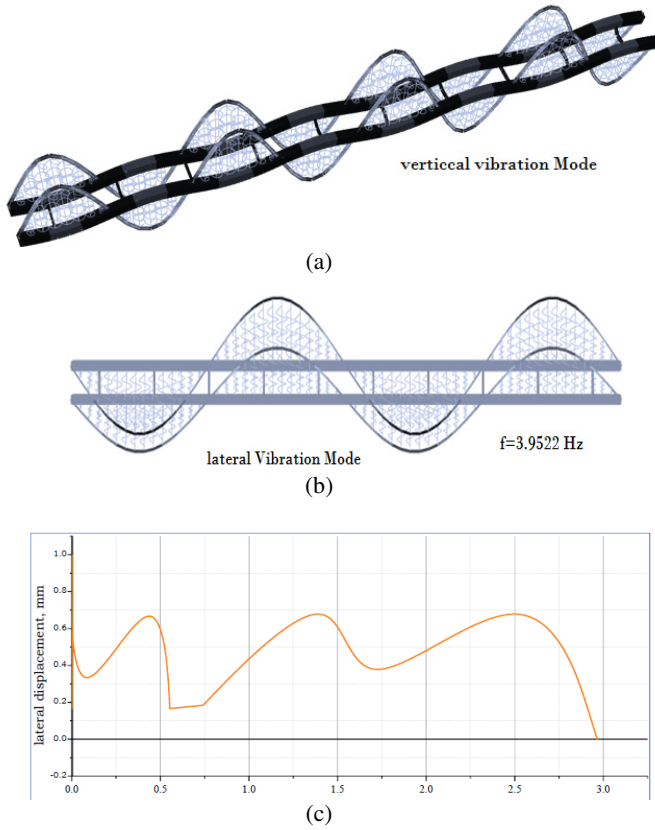


Fig. 7. T-track vibration modes (a) vertical and (b) lateral (c) lateral vibration of T-track with in low frequency and modes

3.9 Dynamic Analysis

The dynamic response of the track system under moving loads at different speeds has been studied. The results populated here will determined the allowable speed for the acceptable frequency and displacement. The contact definition between the wheel and the rail is taken as nonlinear Hertzian contact over closure.

3.10 Attenuation Capacity

Attenuation capacity of the railpad is the damping ability of resonance amplitudes in the track vibration. It is an important indicator of damping capacity of the railpad and other elastic media in the track. The two characteristics of the pad attenuation are based on in-service deflection and the difference in acceleration of the rail to that of the longitudinal reinforced concrete beam. The different in the acceleration between the rail and beam is the attenuation capacity of the railpads. This is an important property of the railpads to isolate the vibration of the rail from being transferred to the beam and down to the formation layer. The relative velocity of rail and beam is an important parameter to examine the dynamic behavior of beam model.

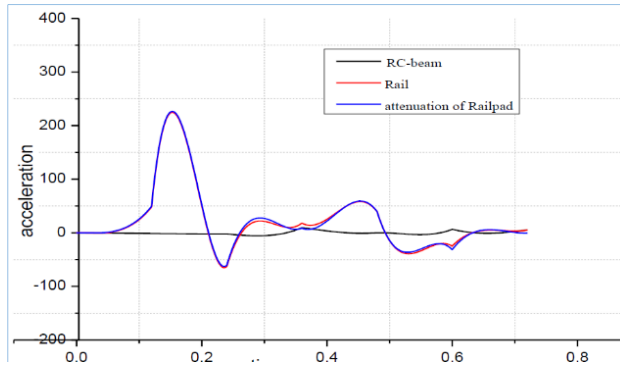


Fig. 8. Acceleration attenuation of Rail pads (time, s and acceleration, m/s²)

At lower modes, the beam and rail vibrate same way but the direction of undulation in most of the cases pattern appears opposite. In such cases, the relative displacement of the two beams increase and stiffer railpads required to attenuate the vibration and decrease the undulation cycle. In other sense, relative resiliencies also required in the system to effectively damp the vibration.

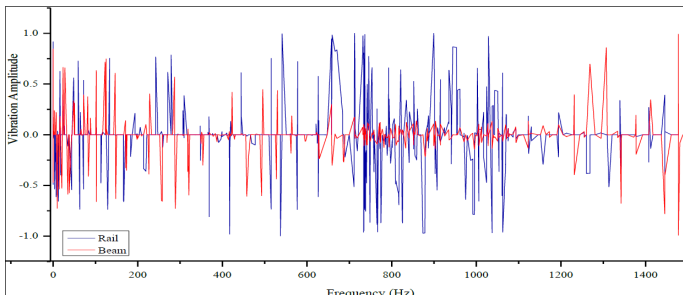


Fig. 9. T-track vertical Vibration (the rail and beam) between the scale of -1 and 1 amplitude Vibration modes and Relative vertical vibration of Rail and Beam

As we can see from the *figure 9* above, the rail vertical vibration pattern dominates the between 700Hz and 1100 Hz frequency range. After 1200 Hz, the vertical deflection of concrete beam dominates the vibration. On the other hand, in the low frequency ranges both beams deflect in a similar fashion. This condition explains the deflection vibration of the rail tends to be damped after about 1050 Hz and in the low range frequencies both beams vibrate in similar fashion up to about 200 Hz. The railpad does its job to damp the vibration after this frequency. It is visible between 550 Hz and 1100 Hz the rail vibrate with larger amplitude than the concrete beam.

Therefore, we can identify about four ranges of frequency that the track likes to vibrate. It is visible that the lower range frequency (*0 to 300 Hz*) vibrations are counteracted by less stiff or soft railpads. In the middle interval (*350 to 700 Hz*) medium stiff springs dominate the vibration. As the frequency intensifies in the high frequency range (*700 to 1100 Hz*), stiffer springs vibrate peak amplitudes and other

springs vibrate with high amplitude. In the last domain higher vibration range, the vibration amplitude tends to damped out where high stiffness springs have small deflection; however track with softer railpad dominate the vibration.

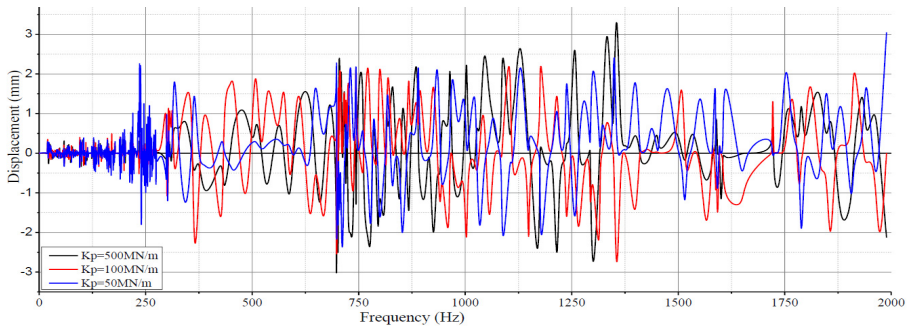


Fig. 10. Track vibration deflection for varying railpads Stiffness

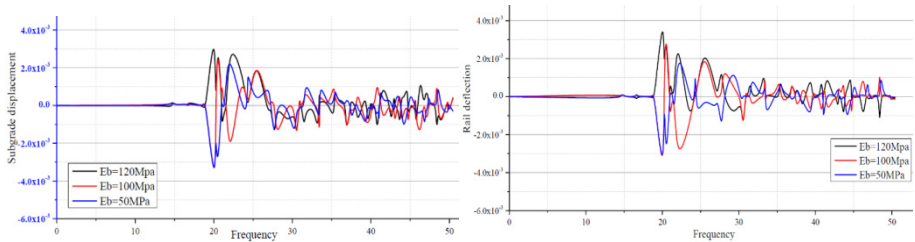


Fig. 11. Effect of subgrade Modulus on rail & subgrade displacement at lower frequency (Hz)

3.11 Study the Effect of Train Speed

Without doubt, railway transport networks have been regaining their importance in recent decades due to the efficiency and environmentally friendly technologies, which has led to increasing train speeds, higher axle loads and more frequent train usage. It was as early as 1938 when De Nie observed dynamic deflection increases due to speed effects in combination with poor subgrade conditions. It has been observed that running high speed trains on railway tracks, constructed on soft ground, induces high levels of vibration in the track and the surrounding area. These vibrations can result in the rapid deterioration of the track structure, causing derailment and ground failure in the worst scenario. The train speed has a direct relation to the subgrade deformation and its stiffness property, K_f . A critical velocity, V_{cr} (a speed that increases the beam displacement dramatically) of the beam-foundation system presented by Kenny [8]:

$$V_{cr} = \sqrt[4]{\frac{4KEI}{\rho^2}}$$

where, K = stiffness, E = Elastic Modulus, I = Area moment, ρ = density.

From the analysis result as shown in the figure 15 (a) to (d) there is a variation in deflection of the roadbed from 33.33m/s to 80m/s train speed.

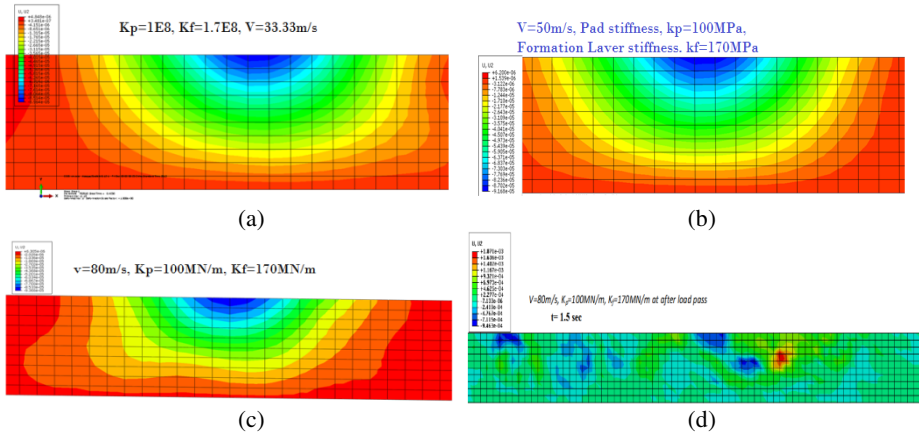


Fig. 12. Deflection of the roadbed (a) 33.33m/s train speed (b) 50m/s train speed (c) and (d) 80m/s train speed

According to the study of Suiker and De Borst (1999), the lowest critical state is mainly determined by the train speed not influenced by other excitation sources such as sleeper spacing in ballasted track. On this basis, the assumption of constant moving loads is reasonable for studying the train speed effect. At a time of 1.5sec and running speed of 80m/s, the deflection of the roadbed appears in scattered regions and local deformation of the roadbed will occur in the track (figure 12 (d)).

The vector diagram shows that the roadbed particles follow a spinning movement just after the load passes that particular location (see figure 13). The amplitude of circular motion of roadbed particles depends on the ratio of wave speed propagated with in the roadbed body to the train speed. When frequency of vibration increases as the train speed increase, the wave speed in roadbed elements also increase. Then failure of roadbed/subgrade occurs when the ratio of the velocities greater than the limit. This needs further study.

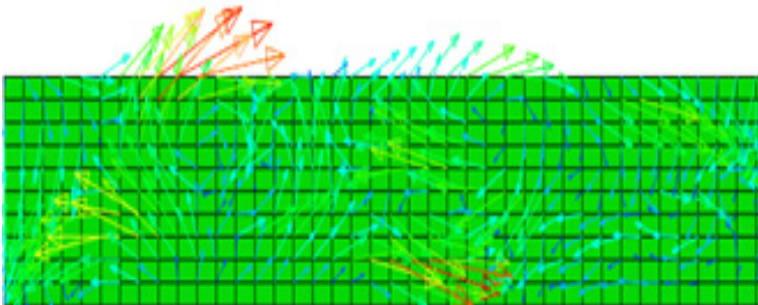


Fig. 13. subgrade particles acceleration under multiple moving loads- a vector diagram

3.12 Effect of Platoon Moving Loads

A moving load is often modeled as a single load with point distribution. In reality, spatial distributions of wheel-rail-track interface is more complex than a point contact. In addition, dynamic response of a beam to an array of moving loads mimics a platoon of vehicles traveling along a railroad. [7] Similar track modeling parameters used to compare spatial response outputs of single wheel load and platoon loads and special concern given to the subgrade layer responses. Taking T50 conventional train parameters a is 2.6m and b is 4.6m, and a 25 tone axle load, the effect of platoon moving wheel load as shown in figures below. The effect of single wheel load model and multiple platoon moving loads is obviously different. Under similar modeling parameters and steady train speed, the kinematic responses of the track are different and multiple load models give higher displacement and stress.

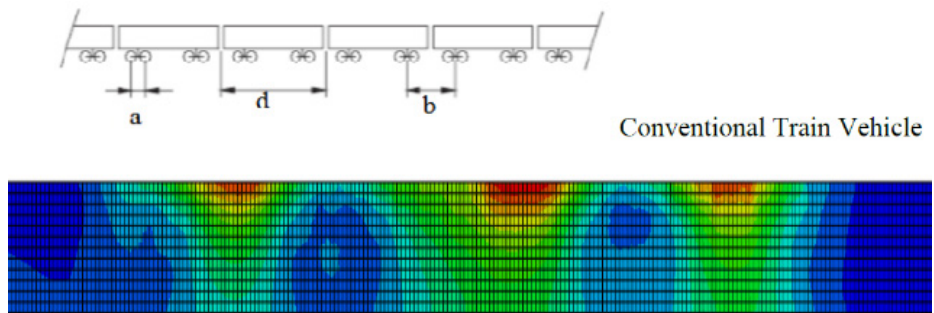


Fig. 14. Roadbed stress under multiple moving wheel loads at 40m/s train speed

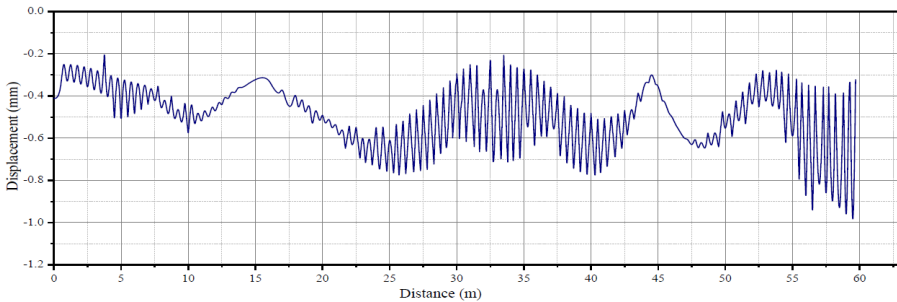


Fig. 15. Roadbed Displacement profile under multiple moving load

4 T-Track for Ethiopian Railway Projects

There are huge technical savings during construction and in-service maintenance as compared to conventional ballasted track. The first technical point is prefabricated and modular nature of the track supper structure. This simplifies the design and construction process and ensure quality of the track with fast construction duration.

The track laying work can be done with in a narrow working space. In addition to economic benefit, dynamic response of the track is found to be stimulating. The rail stress level obtained from the analysis significantly lower than ballasted track.

T-Track is recommended for Ethiopian Railway projects due to the following reasons and advantages.

- The planned and under construction track lines pass through Great Rift Valley and desert regions where vast sand and silt soil dominate. T-track can be the economic solution to alleviate ballast contamination problems
- Due to sand blow in the eastern Ethiopia deserts, clogging of ballast and other granular voids will prevent ease drainage. In case of seasonal rainfall the track will erode and failure of track results. To avoid such causalities, T-track is good alternative.
- Due to less maintenance requirement it is highly recommended
- Due to its suitability, T-track will fit the platform requirement of Addis Ababa light rail transit system
- The technology can be easily adapted without difficulties for the local manpower and technicians.
- In case of track upgrading and replacement, T-track is very simple and cab be effectively upgraded old railway lines can be effectively replaced with T-track.

5 Conclusion

The two parameter model (double beam model) with elastic foundation model of formation layer and solid continuum element result similar results as long as the same constraints and parameters are applied. From the simulation and analysis results the following concluding remarks are forwarded.

- Multiple wheel load model results higher dynamic effect than single moving load so it is realistic to consider the platoon moving loads in modeling so that the ultimate stress and spatial outputs can be accurately calculated. The vibration generated by multiple wheel loads was observed that front wheel starts damped out immediately after the load removed. As the speed of the train increases, another vibration induced on that particular point before the first vibration dissipates. Therefore, a serious of connected vibration amplitudes will be propagated in the track body.
- The direction of wave about a specific point creates an epicenter for the spinning movement (vortex) of subgrade particles. At the local rotational points, very high stress will be generated that will cause roadbed failure.
- The track structure starts to vibrate in longitudinal direction at zero frequency and then to lateral direction, following the vibration way of rails over concrete beams.
- Transverse steel gauge members start vibrating at about middle modes frequency (about 150Hz). This indicates the track's stability (rigidity) in the lateral direction.

- Being continuously supported the rail stress is found to be less that implies smaller rail sections can be used for the same loading and speed to ballasted tracks. In addition, the assembly of rail, beam and steel gauge creates a rigid frame that can bridge irregularities over uneven subgrade.
- The frequency ranges and vibration modes indicate that the track can sustain strong vibrations. This is due to the presence of continuous support railpads.

References

- [1] Zakeri, J.-A., et al.: Effects Of Vibration In Desert Area Caused By Moving Trains. *Journal of Modern Transportation* 20(1) (March 2012)
- [2] Profillidis, V.A.: *Railway Engineering*, Democritus Thrace University, Greece (1995)
- [3] Chopra, A.K.: *Dynamics of Structures*, 3rd edn. University of California at Barkley (2007)
- [4] Esveld, C.: *Track Structures in an Urban Environment*, TU Delft (1997)
- [5] Lubout, N., Gräbe, H.: Performance of Resilient Rail Pads Used in Tubular Modular Track under South African Service Conditions (²Engineer Aurecon and ²University of Pretoria)
- [6] Tubular TrackCompany Profile, <http://www.tubulartrack.co.za>
- [7] Sun, L., Luo, F.: Steady-State Dynamic Response of a Bernoulli–Euler Beam on a Viscoelastic Foundation Subject to a Platoon of Moving Dynamic Loads, Catholic University of America (2008)
- [8] Kenney, J.T.: Steady-State Vibrations of Beam on Elastic Foundation for Moving Load. *Journal of Applied Mechanics* (1954)
- [9] ABAQUS 6.12.1 Documentation
- [10] Tubular Track offers continuous rail support at a competitive price (2005), <http://www.railwaygazette.com>
- [11] Tubular Track – An Array of Advantages, <http://www.railwaysafrica.com>
- [12] Van Schoor, B., Gräbe, H.: An Inside Look At the Stresses Due To Lateral Forces in Tubular Modular Track (1Worley Parsons RSA and 2University of Pretoria)
- [13] Grabe, H.: Track Deflection Measurement, Pretoria University, South Africa
- [14] Wakui, et al.: Technological Innovation in Railway Structure System with Ladder Track System, Railway Technical Research Institute, Tokyo
- [15] Esveld, C.: *Modern Railway Track*, 2nd edn. TU Delft University of Technology (2001)
- [16] Selig, Waters: *Railway Engineering and Substructure Management* (1994)
- [17] Dynatrack, A.: dynamic railway track properties and their quality, TU Delft Technology University, NL (2002)
- [18] Hay, W.W.: *Railroad Engineering*, 2nd edn. University of Illinois, Illinois (1982)
- [19] Dahlberg, T.: *Railway track dynamic– a survey*. Linköping University (2003)
- [20] Esveld, C.: Significance of Track Resilience, TU Delft

Utilizing Text Similarity Measurement for Data Compression to Detect Plagiarism in Czech

Hussein Soori, Michal Prilepok, Jan Platos, and Václav Snášel

Department of Computer Science, FEECS,
IT4 Innovations, Centre of Excellence,
VSB-Technical University of Ostrava,
Ostrava Poruba, Czech Republic
{sen.soori,michal.prilepok,jan.platos,vaclav.snasel}@vsb.cz

Abstract. This paper attempts to apply data compression based similarity method for plagiarism detection. The method has been used earlier for plagiarism detection for Arabic and English languages. In this paper we utilize this method for Czech language text from a local multi-domain Czech corpus with 50 original documents with non-plagiarized parts, and 100 suspicious documents. The documents were generated so that every document could have from 1 to 5 paragraphs. The suspicion rate in the documents was randomly chosen from 0.2 to 0.8. The findings of the study show that the similarity measurement based on Lempel-Ziv comparison algorithms is efficient for the plagiarized part of the Czech text documents with a success rate of 82.60%. Future studies may enhance the efficiency of the algorithms by including combined and more sophisticated methods.

Keywords: similarity measurement, plagiarism detection, Lempel-Ziv compression algorithm, plagiarism in Czech, data compression, plagiarism detection tools.

1 Introduction

The research to find plagiarized texts came as an answer to the massive increase in written data in various fields of science and knowledge and the urge to protect intellectual property of authors. Many text plagiarism methods have been utilized for this purpose.

Similarity detection covers a wide area of research including spam detection, data compression and plagiarism detection. In addition to that, text similarity is a powerful tool for documents processing. Plagiarism detection includes many texts documents such as, students' assignments, postgraduate theses and dissertations, reports, academic papers and articles. This paper investigates the viability of the similarity measurement based on Lempel-Ziv comparison algorithms for detecting plagiarism of Czech texts.

2 Methods of Plagiarism Detection

The similarity measure approach is one of the most used approaches to detect plagiarism. To a certain extent, this method resembles the methods used in information retrieval in that it decides the retrieval rank by measuring the similarity to a certain query.

Generally speaking, we may rank the similarity-based plagiarism detection techniques as: text based techniques that depend on cosine and fingerprints as in [1,2], graph similarity that depend on ontology [3,4] and line matching as in bioinformatics [5].

On the other hand, some machine techniques have demonstrated high accuracy as in the case of k-nearest neighbors (KNN), artificial neural networks (ANN) and support-vector machine learning (SVM).

KNN depends on the idea of X member and its relation to its nearest neighbor. Winnowing algorithm [6] is one of the widely used algorithms in this regard. It depends on the selection of fingerprints of hashes of k-grams. The idea is based on the optimization of results by trying the variations of the n value and how distant or remote is x to its neighbors. In this case, plagiarism is decided on whether or not x falls into the neighboring members, and the number of neighbors. One of the drawbacks of this technique is its slow computational processing because it requires calculation of all the remote and distant members.

The ANN learning classifier functions by modeling the way the brain works with mathematical models. These two are commonly represented (to serve the purpose at hand) by plagiarized and non-plagiarized texts. It resembles the mechanism of human neuron with other neurons that work in another brain layer. In addition to the efficiency of this technique to detect plagiarism, it has also proven very effective in other areas such as speech synthesis.

SVM is a statistical classifier that deals mainly with text. This technique tries to find the boundaries between the plagiarized and non-plagiarized text by finding the threshold for these two classes. The boundary set is called support vectors.

Text similarity measurement can also be used for data compression. Platos et al. [9] utilized similarity measurement in their compression algorithm for small text files. Some other attempts to use the similarity measurement for plagiarism detection include Prilepok et al. [7] who used this measurement to detect plagiarism of English texts and Soori et al. [8] who used the technique to detect plagiarism of Arabic texts. The idea of [7] and [8] was originally inspired by Chuda et al. [10]. In this paper, we utilize the same method to detect plagiarism of Czech texts.

3 Similarity of Text Plagiarism Detection by Compression

The main property in the similarity is a measurement of the distance between two texts. The ideal situation is when this distance is a metric [11]. The distance is

formally defined as a function over Cartesian product over set S with nonnegative real value [12,13]. The metric is a distance, which satisfies four conditions for all:

$$D(x, y) \geq 0 \quad (1)$$

$$D(x, y) = 0; x = y \quad (2)$$

$$D(x, y) = D(y, x) \quad (3)$$

$$D(x, z) \leq D(x, y) + D(y, z) \quad (4)$$

Conditions 1, 2, 3 and 4, above are called: nonnegative, identity, symmetry and the triangle inequality axioms respectively. The definition is valid for any metric, e.g. Euclidean Distance. However, applying this principle into document or data similarity is rather more complicated than that.

The idea behind this paper is inspired by the method used in Chuda et al. [10]. This compression method uses Lempel-Ziv compression algorithm and it is used for text plagiarism detection to detect similarity [16]. The main principle here is the fact that the compression becomes more efficient for the same sequence of data. Lempel-Ziv compression method is one of the currently used methods in data compression for various kinds of data such as, texts, images and audio [17].

4 Creating Dictionaries out of Documents

Creating a dictionary is one task in the encoding process for Lempel-Ziv 78 method [16]. The dictionary is created from the input text, which is split into separate words. If a current word from the input is not found in the dictionary, this word is added. In case the current word is found, then the next word from the input is added to it. This eventually creates a sequence of words. If this sequence of words is found in the dictionary, then the sequence is extended with the next word from the input in a similar way. However, if the sequence is not found in the dictionary, then, it is added to the dictionary with the increased number of sequences property. The process is repeated until we reach the end of the input text. In our experiment, every paragraph has its own dictionary.

5 Comparison of Documents

Comparison of the documents is the main task at hand. One dictionary is created for each and every compared file. Next, the dictionaries are compared to each other. This aims at comparing the number of common sequences in the dictionaries. This number is represented by the parameter in the following formula. The formula below is a similarity metric between two documents.

$$SM = \frac{sc}{\min(c_1, c_2)} \quad (5)$$

Where:

- sc is a count of common word sequences in both dictionaries,
- c_1, c_2 is a count of word sequences in the dictionary of the first or the second document.

The SM value is in the interval between 0 and 1. If $SM = 1$, then the documents are similar. However, if $SM = 0$, then the documents have the highest difference.

6 Experimental Setup

In our experiments we used a local small corpus of Czech texts from the online newspaper iDNES.cz [18]. This corpus contains 4850 words. The topics were collected from various news items including sports, politics, social daily news and science. In our experiment, we needed to have suspicious document collection to test the suggested approach. We created 100 false suspicious and 50 source documents from the local corpus by using a small tool that we designed to create source and suspicious documents as described in the next sub-section.

6.1 Document Creator Tool

The purpose of the document creator tool is to create some documents to be used as testing data. The tool is designed as follows. We split all the documents from the corpus into paragraphs. Next, we labeled each paragraph with a new tagged line for quick reference that indicates its location in the corpus. After that, we created two separate collections of documents from this paragraphs list: a source document collection and a suspicious document collection. Finally, we randomly selected one to five paragraphs from the list of paragraphs. The first group is added to a newly created document and marked as source documents. We repeated this step for all the 50 documents. The collection we have contains 120 different sets of paragraphs.

The same steps were repeated for the group of created suspicious documents. We randomly selected from each suspicious document one to five paragraphs. Then, the tool randomly selected the paragraphs. Each document contains some paragraphs from the source document created earlier, and some unused paragraphs. We repeated these steps for all the 100 documents. To create the list of suspicious documents, we used 227 paragraphs out of the 109 paragraphs from the source documents, and 118 unused paragraphs from our local corpus. For each of the created suspicious document, an XML description file is made. This file includes information about the source of each paragraph in our corpus, its starting and ending point, byte and file name. We repeated this step for all the 100 documents.

6.2 The Experiment

The idea of comparing a whole document where only a small part of it is plagiarized is futile for the reason that the plagiarized part may only be a small chunk hidden among the total volume of text in the new document. Hence, splitting the documents into paragraphs is essential. We chose paragraphs, because we believe they retain some better characteristics than sentences in terms of carrying more sense words, and they are not affected by stop words, i.e. preposition, conjunctives, etc. We separated paragraphs by a line separator and created a dictionary for each paragraph from the source document, according to the steps mentioned earlier. From the fragmentation of the source documents, 120 paragraphs and their corresponding dictionaries were made. These dictionary paragraphs were used as reference dictionaries to be compared against the created suspicious documents.

Similarly, the suspicious document sets were processed in the same manner. The suspicious document were fragmented into 227 paragraphs. Next, a corresponding dictionary was created using the same algorithm, without removing stop words in this step. After that, this dictionary is compared against the dictionaries from the source documents. To accelerate the comparison speed, only a subset of the dictionaries were chosen for comparison for the reason that comparing one suspicious dictionary to all source dictionaries is time consuming. We chose the subset as per the size of a particular dictionary with a tolerance rate of 20%. For instance, if the dictionary of the suspected paragraphs contains 122 phrases, we choose all dictionaries with number of phrases between 98 and 146. The 20% tolerance rate improves the comparison speed significantly. We believe that this tolerance rate does not affect the overall efficiency of the success rate of the algorithm. Accordingly, the paragraph with the highest similarity to each paragraph of the tested paragraph is spotted.

6.3 Stop Words Removal

The removal of stop words increases the accuracy level of text similarity detection. In our experiment, we removed stop words from the text. We used a list of Czech stop words from Semantikoz [19]. The list of the stop words we used contains 138 stop words. The algorithm we used is modified so that after the fragmentation of the text, all stop words are removed from the list of paragraphs and the remaining words are processed by the same algorithm.

7 Results

This paper uses the terms: plagiarized document, to refer to a document which the PlagTool managed to find all its plagiarized parts from the attached annotation XML file; partially plagiarized document, to refer to a document which the PlagTool managed to find parts of it to be plagiarized in the attached annotation XML file (for example, 3 out of 5 parts of the original document); non-plagiarized

Table 1. Results

	Success rate	
Plagiarized documents	76/ 92	82.60%
Partially plagiarized documents	16/ 92	17.40%
Non- Plagiarized documents	8/ 8	100.00%

document, to refer to a document which the PlagTool did not manage to find any plagiarized parts of it in the annotation XML file.

In our experiments we found 86.60% plagiarized documents, 17.40% partially plagiarized documents and 100% non-plagiarized documents.

In case of partially plagiarized documents, we encountered suspicious paragraphs in another document, or paragraphs with higher similarity rate as another paragraph with the same similar content. This case occurs if one of the paragraphs is shorter than the other.

8 The PlagTool and Visualization of Document Similarity in PlagTool

The PlagTool application is divided into four parts. Three parts show document visualisation and the last part shows the result.

In the PlagTool, we have used three visual representations for showing paragraph similarity. These help the user to identify the suspicious documents and identify their locations.

The first representation is a line chart. The line chart shows the similarity for each suspicious paragraph in the document. The user may easily locate the plagiarized parts of the document and the number of the plagiarized parts where higher similarities represent paragraphs with more plagiarized content, as shown in figure 1 where paragraphs two and three are plagiarized.

The second representation is a histogram of document similarity. The histogram depicts how many paragraphs are similar and accordingly how many paragraphs are plagiarized. This is shown in figure 2.

The last visual representation is used to show the similarity in the form of colored highlighted texts. Five colors have been used. The black color is used to show the source document. The red color indicates that the paragraph has a similarity rate greater than 0.2. The orange color shows the paragraphs with lower similarity ranging between less than 0.2 and 0, where the paragraph has only few similar words with the source document. The green text indicates that the paragraph is not found in the source document and accordingly it is not plagiarized. The blue color indicates the location of the plagiarized paragraph in the source document. A snapshot of PlagTool, with the color functions, are depicted in figure 3.

As shown in figures 1, 2, 3, the left table shows all paragraphs from the selected document, starting and ending byte (file location) in the document, the detected highest similarity rates, name of the document with the highest similarity rate and its starting and ending bytes.

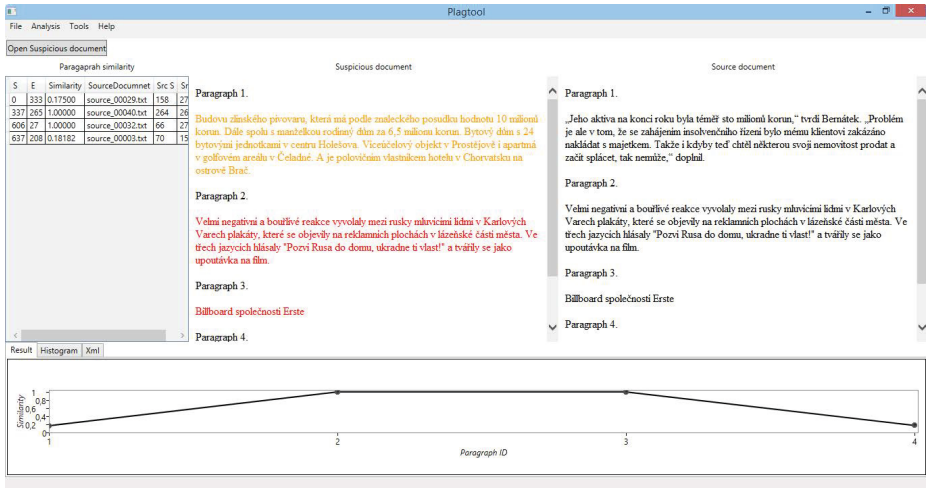


Fig. 1. PlagTool: line chart of document similarity

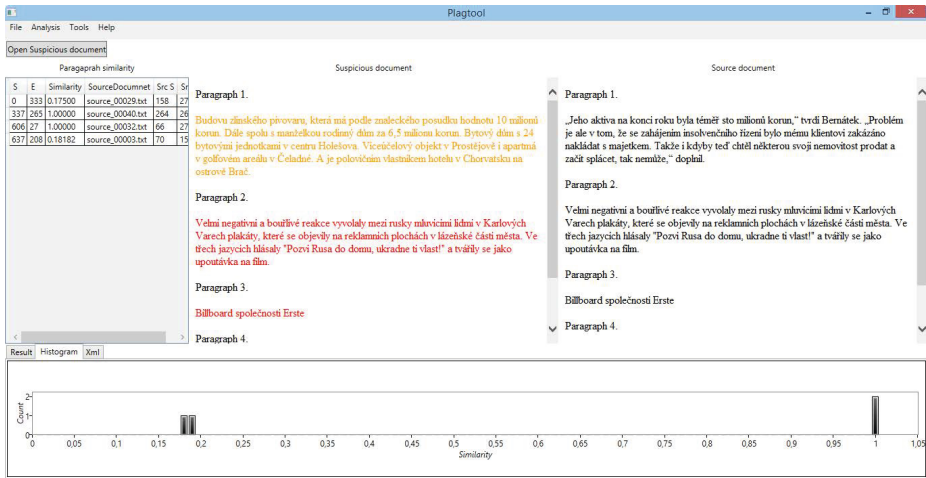


Fig. 2. PlagTool: histogram of document similarities

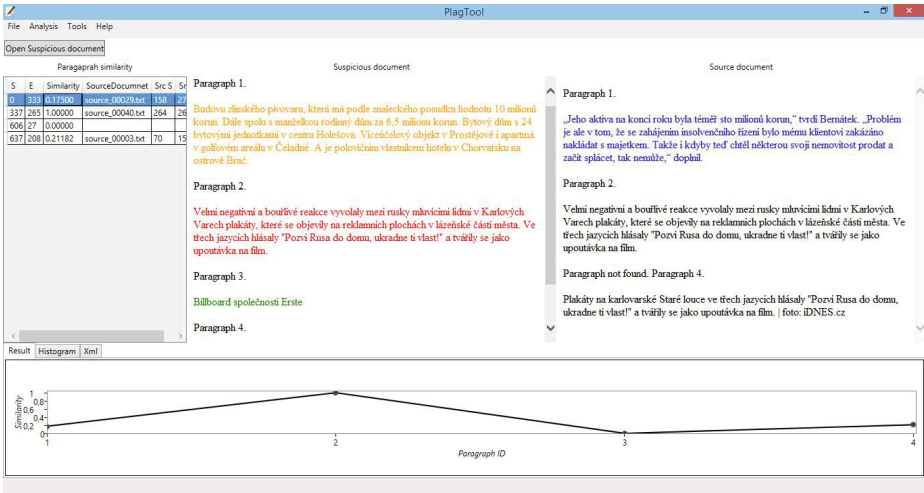


Fig. 3. PlagTool: colored visualization of similarity

The middle table opens suspicious document with color highlighting. The color highlighting helps the user to easily identify the plagiarized paragraphs. The part on the right side shows the content of paragraphs with the highest detected similarity to source documents. This helps the user to compare between the content of the paragraph from the suspicious document and the paragraph from the source document.

The bottom part displays information about the results using three tabs. The first two tabs show various charts of the selected suspicious document. These charts give the user information about the overall rate of plagiarism for the selected document. The last tabs shows information about the suspicious document from the testing data set. It also shows which paragraphs are plagiarized, and from which document they were copied. This information is useful for validating the success rate of plagiarism detection.

9 Conclusion

In this paper, we applied the similarity detection algorithm used by Prilepok et al. [7] for plagiarism detection of English text, and Soori et al. [8] for plagiarism detection of Arabic text on a dataset for Czech text plagiarism detection. We also confirmed the ability to detect plagiarized parts of the documents with the removal of stop words, as well as the viability of this approach for Czech language. The algorithm for similarity measurement based on the Lempel-Ziv compression algorithm and its dictionaries proved to be efficient in detecting the plagiarized parts of the documents. All plagiarized documents in the dataset were marked as plagiarized and in most cases all plagiarized parts were identified, as well as, their original versions, with a success rate of 82.6%.

Acknowledgement. This work was supported by the European Regional Development Fund in the IT4Innovations Centre of Excellence project (CZ.1.05/1.1.00/02.0070) and by Project SP2014/110, Parallel processing of Big data, of the Student Grant System, VSB - Technical University of Ostrava.

References

1. Pera, M.S., Ng, Y.K.: SimPaD: A word-similarity sentence-based plagiarism detection tool on Web documents. *Web Intelligence and Agent Systems* 9(1), 27–41 (2011)
2. Gustafson, N., Pera, M.S., Ng, Y.K.: Nowhere to hide: Finding plagiarized documents based on sentence similarity. In: *Proceedings of the 2008 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology*, vol. 01, pp. 690–696. IEEE Computer Society (2008)
3. Foudeh, P., Salim, N.: A Holistic Approach to Duplicate Publication and Plagiarism Detection Using Probabilistic Ontologies. In: Hassanien, A.E., Salem, A.-B.M., Ramadan, R., Kim, T.-h. (eds.) *AMLTA 2012. CCIS*, vol. 322, pp. 566–574. Springer, Heidelberg (2012)
4. Liu, C., Chen, C., Han, J., Yu, P.S.: GPLAG: detection of software plagiarism by program dependence graph analysis. In: *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 872–881. ACM (August 2006)
5. Chen, X., Francia, B., Li, M., Mckinnon, B., Seker, A.: Shared information and program plagiarism detection. *IEEE Transactions on Information Theory* 50(7), 1545–1551 (2004)
6. Schleimer, S., Wilkerson, D.S., Aiken, A.: Winnowing: local algorithms for document fingerprinting. In: *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data*, pp. 76–85. ACM (June 2003)
7. Prílepok, M., Platos, J., Snasel, V.: Similarity based on data compression. In: Castro, F., Gelbukh, A., González, M. (eds.) *MICAI 2013, Part II. LNCS*, vol. 8266, pp. 267–278. Springer, Heidelberg (2013)
8. Soori, H., Prilepok, M., Platos, J., Berhan, E., Snasel, V.: Text Similarity Based on Data Compression in Arabic. In: Zelinka, I., Duy, V.H., Cha, J. (eds.) *AETA 2013. LNEE*, vol. 282, pp. 211–220. Springer, Heidelberg (2014)
9. Platos, J., Snasel, V., El-Qawasmeh, E.: Compression of small text files. *Advanced Engineering Informatics* 22(3), 410–417 (2008)
10. Chuda, D., Uhlik, M.: The plagiarism detection by compression method. In: *Proceedings of the 12th International Conference on Computer Systems and Technologies*, pp. 429–434. ACM (June 2011)
11. Tversky, A.: Similarity features. *Psychological Review* (84), 327–352 (1977)
12. Cilibrasi, R., Vitányi, P.M.: Clustering by compression. *IEEE Transactions on Information Theory* 51(4), 1523–1545 (2005)
13. Li, M., Chen, X., Li, X., Ma, B., Vitanyi, P.M.: The similarity metric. *IEEE Transactions on Information Theory* 50(12), 3250–3264 (2004)
14. Kirovski, D., Landau, Z.: Randomizing the replacement attack. In: *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2004)*, vol. 5, p. V-381. IEEE (May 2004)

15. Crnojevic, V., Senk, V., Trpovski, Z.: Lossy lempel-ziv algorithm for image compression. In: 6th International Conference on Telecommunications in Modern Satellite, Cable and Broadcasting Service, TELSIKS 2003, vol. 2, pp. 522–525. IEEE (October 2003)
16. Ziv, J., Lempel, A.: Compression of individual sequences via variable-rate coding. IEEE Transactions on Information Theory 24(5), 530–536 (1978)
17. Khoja, S., Garside, R.: Stemming Arabic text. Computing Department, Lancaster University, Lancaster, UK (1999)
18. iDNES.cz (2014), <http://www.idnes.cz/>
19. Semantikoz (2014), <http://www.semantikoz.com/blog/free-stop-word-lists-in-23-languages/>

AJAX Speed Up vs. JSP in Portal – Case Study

David Ježek¹ and Radek Liebzeit²

¹ VŠB – Technical University Of Ostrava, 17. listopadu 15,
Department of Computer Science, Czech Republic
david.jezek@vsb.cz,
<http://www.cs.vsb.cz/jezek/>

² VŠB – Technical University Of Ostrava, 17. listopadu 15,
Center of Information Technology, Czech Republic
radek.liebzeit@vsb.cz

Abstract. Main purpose of the paper is comparison of throughput of AJAX and JSP solution in highly loaded part of JavaEE information system. Prepared performance tests and observed values of response time and utilization of resources push us to investigate cause of low throughput of whole information system infrastructure.

Keywords: performance test, JavaEE, AJAX, JSP, IBM WebSphere Portal, JBoss Cahce.

1 Introduction

A university developed and maintains own information system EdISoN [11] that provides much functionality for students, researchers and employees. Performance of the system is sufficient in most days in years, but there are two big peeks, before each semester. All students of faculty create their individual schedules. To prevent overloading of infrastructure are students separates to group by year of study and study results to 10 groups. Despite this system infrastructure must serve more than 3000 requests in 10-15 minutes.

Students often complain of slowness and freezing of individual schedule creation functionality. There was no acceptable way to increase number of licenses and upgrade of hardware to handle these 15 minutes peeks.

Our work was focused to rewrite the concerned functionality with usage of AJAX technology [8] to save server resources and increase of throughput of system. New solution was put under performance test [1] to compare response time and throughput of system.

After AJAX solution was deployed on server and put to stress by students, solution was evaluated by them via short questionnaire.

The main goal was compare increase effectiveness of new solution based on AJAX technology with previous solution coded via JSP pages on portal server.

2 Examined Application and Overload Peaks

2.1 Application and Known Problems

A department of computer science used to use self-development information system [12] from 1995 to handle common agenda of education and management system. Positive experience caused extension of usage of the system to whole faculty.

As result of the continuous positive experience was decided to use self-development information system for whole university 15 years ago.

Unfortunately design and implementation of existing system was inconvenient to handle so big extend of users and functionality. The development of new information system named EdiSoN started 10 years ago from scratch. Project stakeholders decide to use IBM WebSphere Portal [4] solution with JSP technology [9]. First release Project was deployed in 2005. Maintenance and development of the system is still in progress. The functionality of the system is still extended to handle new futures based on users requirements.

One of main future of system for student is possibility to create individual study schedule. The process contains three main steps. First of all student select their desired subjects. Based on the number of students in particular subjects faculty publish list of lectures and labs with defined time and lecturer and capacity. At the end students can sign in to the preferred lectures and labs until capacity is fulfilled. The last step is most difficult one because faculty has about 1500 students and almost all of theme wants sign in their preferred lectures and labs as soon as possible. That means almost all students create their individual schedules immediately after sign in is launched. Current infrastructure cannot handle so many request therefore student are separated to groups to decrease number of concurrent users. Nevertheless biggest group contain 600 students and many of them complaints about big response time, freezing and other errors typical for overloaded systems.

2.2 Selected Solution

The problem with response time and overall throughput of system was only in functionality of individual schedule creation that is used two times per year. To minimize required effort and resources to resolve the problem stakeholders decide to re-implement only requests and page for individual schedule creation. Old solution with JSP components was replaced with AJAX that still use same business logic implemented in EJB as previous JSP pages.

2.3 Java EE Application and Deployment Description

The information system consists from two WebSphere Portal deployed on two servers . Both of servers are connected to application server [3] that is connected to database server, provided by DB2 server. Web requests are distributed by proxy server based on request subnet and web server load. Hardwer specification is denoted in table 1.

Table 1. Hardwar specification

		Proxy	Portall & 2	Application	Database
HW	CPU	Intel(R) Xeon(R)			
		E5-2660 v2	E5-2665 0	E5-2660 v2	E5-2660 v2
	Count	2	2	12	6
	Freq	2.2GHz	2.4GHz	2.2GHz	2.2GHz
	RAM	4GB	16GB	16GB	24GB
OS	CentOS 6.4	Red Hat Enterprise Linux Server release 6.4			
SW	Nginx 1.6	IBM WS Portal Server 8.0.0.1	IBM WS AS 8.5.5.0	IBM DB2 Wg. Ser. Ed. 10.1	

3 Performance Test Configuration

Performance test infrastructure [7] was assembled from computers of one of our student computer lab. Unfortunately between test agents and the information system was quite wide network infrastructure but it doesn't lead to significant result distortion .

As performance test tool was used IBM Performance Tester 8.2 that is based on open source project Eclipse Test & Performance Tools Platform Project.

Performance test script [2] was designed to simulate students' behavior with is:

1. Login to system
2. Go to page with schedule
3. Select subject.
4. Select favorite courses or labs.

The test script was created quite straightforward to this behavior [7] [6]. After virtual user get page with schedule, AJAX request and response are processed in cycle to simulation of selection of several subjects and courses or labs. Ajax request and response was handled using custom code block that parse response and extract possible subjects and schedule activities and chose randomly one.

The last optimization of script was connected with problem in login phase and information system proxy server setting. Login page with SSO technology has not enough throughputs to handle all virtual users at time so we added startup delay for virtual users to distribute login in time. Proxy server was setup to balance load based on source subnet of request that fails for agents from one subnet. The simplest solution at the time was bypass the proxy server and directly access servers with deployed portals and distribute request using data pool mechanism build in performance test tool.

Scripts for the old pages based on JSP technology was almost same only parsing of possible schedule activities was redesigned to handle standard HTML page content.

Monitoring of system was handled by HypericHQ 5 that monitor standard parameters like processors utilization, memory utilization, amount of disk operation,

and amount of data transfer through network etc. Another important attribute for monitoring was amount of scheduled activities for all students stored in database.

4 Performance Tests Runs

First test was run with 300 virtual users directly to one of the portal with pages using new AJAX solutions. Result was in expectable range of values with comparison to previous experience (Fig. 1, Tab. 2). Now we want to increase load to utilize more of the infrastructure resources. Main resources utilization was on portal and application servers but it was only 20% utilization of processors in average.

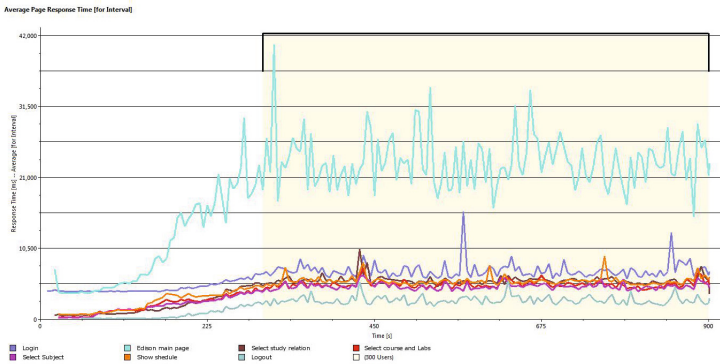


Fig. 1. First test response time graph

Table 2. Performance counter values for first test run

	Response Time [ms]				Rate [per sec]
	Min	Average	Max	σ	
Login	4106	6230.6	27535	2115.6	0.77
Edison main page	3709	19513.5	107634	12224.7	1.54
Select study relation	528	4551.7	14690	2051.5	0
Show schedule	598	4671.4	47331	2101.9	0
Select Subject	133	4097	49782	2074.5	2.31
Select course and Labs	53	4353.6	47946	2656.9	6.15
Logout page	53	2308.3	23413	1761.2	0.77

We run script with 500 users but response time increase to 17s in average for schedule creation and processor utilization remain same. All other monitored resources as memory, disk operation and network were not fully utilized too. We minimize sleep time between requests in performance scripts to gain bigger processor utilization but with no success.

4.1 Search for Bottleneck

Somewhere in system was hidden throttle that cause low throughput of requests and low utilization of all monitored resources. We don't have any track to the problem. We run many experiments with different settings to discover throttle.

Experiments:

1. First pages with login (SSO technology) generate quite long response time. Most of students login into system before startup of schedule creation and wait to proper moment. Startup delay was added for virtual users to distribute login in time and synchronization point was added into script that allow hosted synchronize all virtual clients on all agents. With the synchronization point users starts create schedule at same time as real students does.
Result: After synchronization point system was utilized only by requests related to schedule creation.
2. Split request to both portal servers through redirect script from direct portal access to proxy server.
Result: Proxy server was setup to balance load based on source subnet of request that fails for test agents from one subnet.
Changes: The simplest solution at the time was bypass the proxy server and directly access servers with deployed portals and distribute request using data pool mechanism build in performance test tool to both of them.
3. Test performance of both portal servers together.
Result: Amount of scheduled activities stored in database was same. Response time of requests was almost two times longer.
Changes: Reconfigure Apache server to not wait for synchronization.
4. Log duration of most important methods in processing requests.
Result: Method which save scheduled activity – duration in average 0.2s. All requests call another two methods to list offered activities. Each of the methods was duration 2s in average. Main HQL query last for 10ms. Rest of time was spent to iterate objects of offered courses and labs.
5. Application server utilizes only 50 threads from (80 allowed). Small servlet was implemented to call directly EJB components to utilize more threads and processor's cores.
Result: Utilization of processors and threads was smaller.
6. Source code optimization.
Result: Performance slightly increased.
7. Increase number of dedicated thread in application server and increase of maximum allowed opened files on database server. *Result:* Performance slightly increased. But performance of all other pages in information system was lower. *Changes:* Team experimentally detected number of allowed thread to optimize performance of schedule creation and other pages.

These experiments allows us increase performance of schedule selection, but difference between new AJAX solution and previous JSP solution was small and

do not fulfill our expectation. Utilization of CPU cores was maximally around 50%. Finally we found main bottleneck.

4.2 Main Cause of Low Throughput

We have no glue where to search solution to increase throughput and performance of system. Fortunately one member of our team remembered previous problems with JBoss Cahce [10] so we decide try to switch off the cache and performance of portals and was increased significantly. Maximal utilization of processor cores was over 80%.

JBoss Cache was source of problem in 2010 when after upgrade of application server from version 7.0.0.3 to 7.0.0.5 which cause upgrade of IBM JDK. After that update bigger load of users cause low CPU utilization and long response time. That behavior was observed on operating system SLES 11 but on operating system SLES 10 and Red Hat Enterprise Linux 5 work system without problems.

Another several experiments allow us optimizes setting of all parameters for new configuration of system without JBoss Cache. Response time for pages responsive for creating student’s schedules was decreased from 4s and 21.5s for system with JBoss Cache to 2.2s and 2.8s for system without JBoss Cache (Fig. 2 , Table 3).

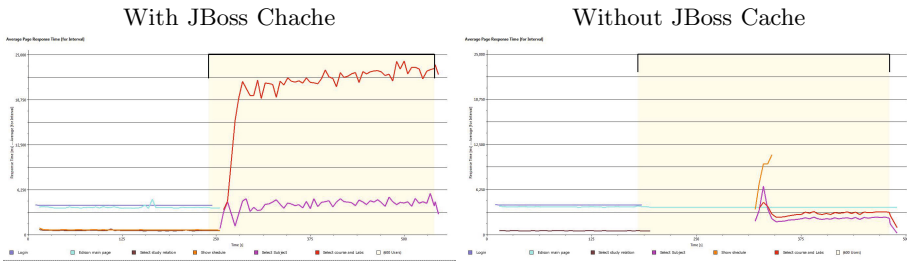


Fig. 2. Performance graph for comparison with and without JBoss Cache

Table 3. Performance counters values with and without JBoss TreeCache

	With JBoss TreeCache				Without JBoss TreeCache			
	Response Time [ms]				Response Time [ms]			
	Min	Avg	Max	σ	Min	Avg	Max	σ
Login	4074	4087	4211	0	4093	4123	4237	0
Edison main page	3587	3813	6916	313	3742	3858	4294	57
Select study relation	506	579	1445	74	497	546	825	29
Show schedule	579	675	1127	50	1788	7928	17558	3484
Select Subject	259	4064	10202	1279	88	2253	13596	945
Select course and Labs	2595	21576	39910	5968	171	2879	13610	658
Throughput average	1388 selected courses per min				7253 selected courses per min			

Another experiment was designed to verify behavior of system under load of 1700 users that corresponds to number of all students in our faculty plus some reserve. Results show that system was able to handle requests in acceptable time (Fig. 3, Table 4).

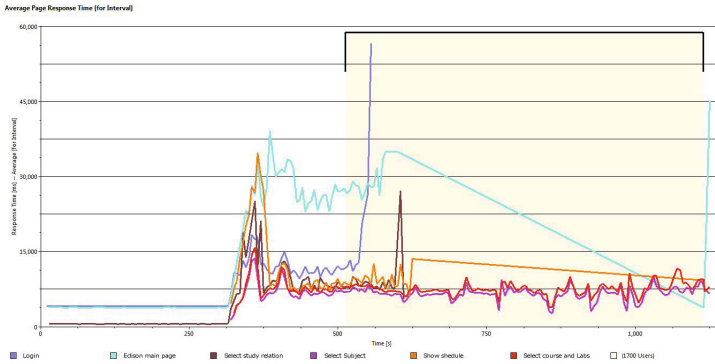


Fig. 3. Performance graph for run with 1700 users

Table 4. Performance counters values for run with 1700 users

	Response Time [ms]			
	Min	Average	Max	σ
Login	4100	7268.2	56563	4557.4
Edison main page	3706	13407.2	51128	12524.8
Select study relation	501	3864.4	33300	4467.6
Show schedule	2230	12287.5	42760	7405.3
Select Subject	83	6673.8	67027	4505.3
Select course and Labs	62	7553	73982	4560.4
Throughput average	10028 selected courses per min			

Now we were able to perform comparison of AJAX and JSF solution that produce reasonable results.

4.3 Comparison of AJAX Speedup to JSP Solution

System was correctly configured and there was observed no strange behavior. Utilization of processor cores, memory, disk and network was proportional to system load. Now we were able to measure speed up of new solution that use AJAX request to obtain data form server. Process of creating HTML code to present data to user was transfer from web portal to client web browser. Old solution use JSP pages that prepare HTML code on server side. Because system run under portal solution not only HTML code for represent new data must

be created but there is also created code for portal wrapper components. Both solutions use same enterprise beans so, gathering data from database generate same load on application server.

Measurement data (Fig. 4, Table 5) shows that speedup for observed pages were 6.4 times and 3.8 times in average. We can deduced that speedup of some pages is around 5 times for AJAX technology that override JSP processing and portal page construction.

Effort required to reimplementation of those pages was 999 man-hour. In comparison with price of better hardware is much more effective especially in case of local peek that use only small part of functionality of information system.

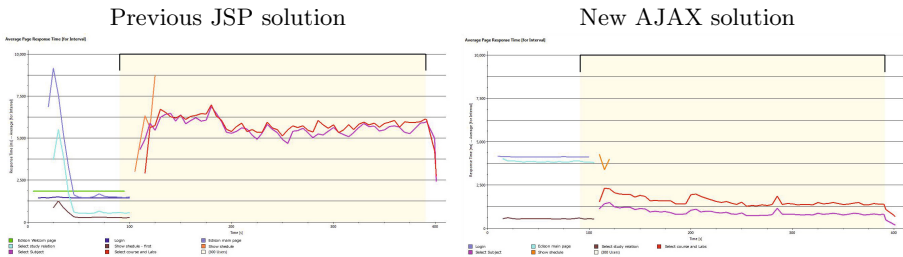


Fig. 4. Performance graph for comparison JSP and AJAX solution

Table 5. Performance counters values for JSP and AJAX solution

	JSP solution				AJAX Solution			
	Response Time [ms]				Response Time [ms]			
	Min	Avg	Max	σ	Min	Avg	Max	σ
Edison Welcom page	1823	1834.3	2028	0				
Login	1432	1450.2	1971	33.3	4107	4126.3	4311	0
Edison main page	1332	3261.6	15622	2937.5	3734	3852.7	4865	88.7
Select study relation	493	1581.1	7527	1768	505	555.3	852	33.3
Show schedule - first	169	464.3	2205	339.3				
Show schedule	2069	4713.3	11434	1770	2391	4096.3	7823	1102.8
Select Subject	2247	5593.9	14478	764.3	45	871.7	4511	344.1
Select course and Labs	2303	5806.6	14493	715.5	5	1504.1	5676	584.6
Throughput average	1388 selected courses per min				7253 selected courses per min			

4.4 Comparison of AJAX Speedup to JSP Solution

System was correctly configured and there was observed no strange behavior. Utilization of processor cores, memory, disk and network was proportional to system load. Now we were able to measure speed up of new solution that use AJAX request to obtain data form server. Process of creating HTML code to

present data to user was transfer from web portal to client web browser. Old solution use JSP pages that prepare HTML code on server side. Because system run under portal solution not only HTML code for represent new data must be created but there is created code for also for portal wrapper components. Both solutions use same enterprise beans so, gathering data from database generate same load in all solutions.

Measurement data shows that speedup for observed pages was 6.4 times and 3.8 times in average. We can deduced that speedup of some pages is around 5 times for AJAX technology that override JSP processing and portal page construction.

Effort required to reimplementation of those requests and page was 34 man-hour. In comparison with price of better hardware is much more effective especially in case of local peek that use only small part of functionality of information system.

4.5 Questioner

Test results show great speedup of system, but experience of students may be different. Experience of users is on of metrics for usability testing [5] [1]. We decide investigate users experience with schedule creation during performance peek by short questionnaire.

Developed questionnaire evaluates experience of users witch tuned system. Questionnaire consists from 10 questions. We send email with link to online questionnaire to 2531 students of our faculty. Students as anonymous questioners fill 811 (32%) questionnaires, 228 (9%) students just open questionnaire and don't respond any questions, rest of students 1492 (58.9%) don't respond to email.

The most important result for us was information about speed of application. Questionnaire contains two question related to speed of application. First was question if speed of application is sufficient (Fig. 5). Second was if speed of application was better (subjectively) with comparison to previous years (Fig. 5).

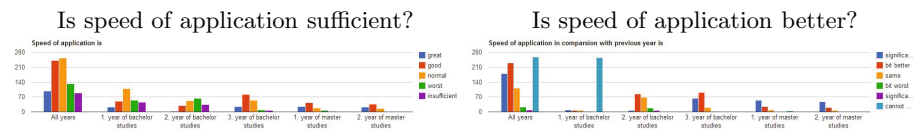


Fig. 5. Questionnaire results

Most unpleasant result from questionnaire was that over 50% of student from first and second year of study get error during schedule creation.

5 Conclusion

Main purpose of our experiments was measure speed-up of newly implemented part of Edison information system for creation of individual student's schedules. Previous solution was implemented with JSP under IBM WebSphere Portal. New solution use AJAX technology to override rendering of full HTML page on server side.

We found badly functional component (JBoss Structural Cache) during the experiments. The component cause poor throughput of system. After switch off of the component response time for database read-only request was two times shorter and response time for database write request was seven times shorter.

Comparison of JSP in portal page and AJAX implementation after switch off JBoss Structural Cache component shows 6 and 4 times shorter response time for both requests, that we was focused on.

References

1. Black, R.: Advanced Software Testing. Guide to the ISTQB Advanced Certification As an Advanced Test Analyst (Rockynook Computing), vol. 1. Rocky Nook (2008)
2. Copeland, L.: A Practitioner's Guide to Software Test Design. Artech House, Inc., Norwood (2003)
3. International Business Machines Corp. Websphere application server, <http://www-03.ibm.com/software/products/en/appserv-was> (accessed: June 6, 2014)
4. International Business Machines Corp. Websphere portal server, <http://www-03.ibm.com/software/products/en/portalserver> (accessed: June 6, 2014)
5. Kan, S.H.: Metrics and Models in Software Quality Engineering, 2nd edn. Addison-Wesley Longman Publishing Co., Inc., Boston (2002)
6. Lang, A.: Ibm websphere portal: Performance testing and analysis (May 2013)
7. Molyneaux, I.: The Art of Application Performance Testing: Help for Programmers and Quality Assurance, 1st edn. O'Reilly Media, Inc. (2009)
8. Olson, S.D.: Ajax on Java. O'Reilly Media (2007)
9. Perry, B.W.: Java Servlet & JSP Cookbook. O'Reilly Media (2004)
10. Surtani, M.: Jboss cache users' guide - a clustered, transactional cache, http://docs.jboss.org/jboss-cache/3.2.1.GA/userguide_en/html_single/index.html (accessed: June 6, 2014)
11. VSB - Technical University of Ostrava. Edison - education information system on net, <http://edison.vsb.cz> (accessed: June 6, 2014)
12. VSB - TUO, Department of Computer Science. Information system katis, <http://katis.cs.vsb.cz/> (accessed: June 6, 2014)

Intelligent Decision Support for Real Time Health Care Monitoring System

Abdelhamid Salih Mohamed Salih¹ and Ajith Abraham^{2,3}

¹ Sudan University of Science and Technology, Faculty of Computer Science, Khartoum, Sudan
hamidsalinh39@yahoo.com

² Machine Intelligence Research Labs (MIR Labs), Washington, USA

³ IT4Innovations -Center of excellence, VSB -Technical University of Ostrava, Czech Republic
ajith.abraham@ieee.org

Abstract. In the health care monitoring, data mining is mainly used for classification and predicting the diseases. Various data mining techniques are available for classification and predicting diseases. This paper analyzes and evaluates various classification techniques for decision support system and for assisting an intelligent health monitoring system. The aim of this paper is to investigate the experimental results of the performance of different classification techniques for classifying the data from different wearable sensors used for monitoring different diseases. The Base Classifiers Proposed used in this work are IBk, Attribute Selected Classifier, Bagging, PART, J48, LMT, Random Forest and the Random Tree algorithm. Experiments are conducted on wearable sensors vital signs data set, which was simulated using a hospital environment. The main focus was to reduce the dimensionality of the attributes and perform different comparative analysis and evaluation using various evaluation methods like Error Metrics, ROC curves, Confusion Matrix, Sensitivity and Specificity. Experimental results reveal that the proposed framework is very efficient and can achieve high accuracy.

Keywords: Classification, Attribute Selected Classifier, Bagging, wearable sensors.

1 Introduction

There is growing need to supply constant Health Care Monitoring (HCM) and support to patients with Chronic Diseases (CD) especially the disabled, and elderly. Wireless sensor networks (WSNs) are used for gathering the information needed. The information may consist of many different sensors such as vital signs (e.g. heart rhythm or blood pressure), etc. Thus, most of the context information can be collected by distributed sensors throughout the environment and even the users themselves [1]. Sensors data is collected from disparate sources and later need to be classified and analyzed to produce information that is more accurate, more complete, or more insightful than the individual pieces. To deal with the large volume of data produced by these special kinds of wireless networks, one approach is the use of Data Mining

techniques. Data mining plays a vital role in various applications such as business organizations, e-commerce, health care industry, scientific and engineering. In the health care industry, the data mining is mainly used for classification and predicting the diseases from the datasets. Various data mining techniques are available for predicting diseases namely Classification, Clustering, Association rules and Regressions. Classification is an important task in data mining. Classification of sensory data is a major research problem in WSNs. Woo et al. [2] proposed ECG signal monitoring system using body sensors. In the research used work One-class support vector machine classifier was used to detect abnormal heart signal values. Patel et al. [3] presented an approach to estimate the severity of symptoms based on accelerometer sensor data. The results of SVM based classification were compared and verified. Korel et al. [4], proposed context awareness Body area sensor network based health-monitoring system to detect abnormal episodes in the signal. Anthony et al. [5] proposed a research work used to recognize various activities of a user using a smart home. The data collected are classified using Multi class SVM and the result is compared for different kernel functions. Some Authors [6] proposed Classification Technique of Human Motion Context based on Wireless Sensor Network. Body sensor nodes are equipped with accelerometer; human motion will cause the waveform of the accelerometer to change accordingly. This change in waveform captured by sensor nodes is then analyzed by PCA (principal component analysis) and SVM (Support Vector Machine) method for clustering and classification. In this paper, Simulation Wearable sensors Data of vital signs are used for developing a decision support system. The rest of this paper is organized as follows. Section 2 describes the methods used for evaluation and the base proposed classifiers. Section 3 presents the Experimental Results and is analyzed in Section 4 followed by discussions in Section 5.

2 Computational Intelligence

2.1 Base Classifiers Used

A) Decision tree algorithm J48

J48 classifier is a simple C4.5 decision tree for classification. It creates a binary tree. The decision tree approach is most useful in classification problem. With this technique, a tree is constructed to model the classification process. Once the tree is built, it is applied to each tuple in the database and results in classification for that tuple [8-9].

B) Logistic Model Trees (LMT)

A logistic model tree (LMT) [10] is an algorithm for supervised learning tasks, which is combined with linear logistic regression and tree induction. LMT creates a model tree with a standard decision tree structure with logistic regression functions at leaf nodes. In LMT, leaves have a associated logic regression functions instead of just class labels.

C) Random Forest

Random forest [11] is an ensemble classifier that consists of many decision tree and outputs the class that is the mode of the class's output by individual trees. Random Forests grows many classification trees without pruning. Then each decision tree classifies a test sample and random forest assigns a class, which have maximum occurrence among these classifications.

D) Random Tree

A random tree is a tree formed by stochastic process. Types of random trees include Uniform spanning tree, Random minimal spanning tree, Random binary tree, Random recursive tree, Treap, Rapidly exploring random tree, Brownian tree, Random forest and branching process [12].

E) Meta-learning

Meta-learning is about learning from learned knowledge [13]. The idea is to execute a number of concept learning processes on a number of data subsets, and combine their collective results through an extra level of learning. Meta-learning aims to compute a number of independent classifiers by applying learning programs to a collection of independent and inherently distributed databases in parallel. The “base classifiers” computed are then collected and combined by another learning process. The most popular meta-learning algorithms are bagging and boosting. Bagging [14] is a method for generating multiple classifiers (learners) from the same training set. The final class is chosen by, e.g., voting.

F) PART

Rule-based learning, especially decision trees (also called classification trees or hierarchical classifiers) is a rule generator that uses J48 to generate pruned decision trees from which rules are extracted [15].

G) IBK

The lazy IBk (commonly known as K- nearest neighbor) is one of classification algorithms that uses distance weighting measures with capability of various attributes like Date attributes, Numeric attributes, Unary attributes, Nominal attributes, Missing values, Binary attributes and Empty nominal attributes. K-nearest neighbours classifier can select appropriate value of K based on cross-validation and also do distance weighting [16-17].

2.2 Attribute Selection

It is often an essential data processing step prior to applying a learning algorithm. Reduction of the attribute space leads to a better understandable model and simplifies the usage of different visualization technique. Attribute selection reduces dataset size by removing irrelevant and redundant attributes. It finds a minimum set of attributes such that the resulting probability distribution of data classes is as close as possible of original distribution. Attribute evaluator method and the search method Best first evaluates the

worth subset of attributes by considering the individual predictive ability of each attribute [18]. In the preprocessing step, we have changed the class attribute to Abnormal or Normal where an 'Abnormal' specifies class 1 and a 'Normal' Specifies class 0.

2.3 Cross-Validation Method

In this paper, we applied a 10-fold cross validation test option. Cross-Validation (CV) is a statistical method of evaluating and comparing learning algorithms by dividing data into two segments: one used to learn or train a model and the other used to validate the model. The basic form of CV is k-fold CV. In k-fold CV the data is first partitioned into k equally (or nearly equally) sized segments or folds. Subsequently k iterations of training and validation are performed such that, within each iteration a different fold of the data is held-out for validation while the remaining k -1 folds are used for learning. The advantage of K-Fold Cross validation is that all the examples in the dataset are eventually used for both training and testing.

2.4 Methods Used for Evaluation of Algorithms

We evaluate our classifiers by measuring their performance by various methods and performance matrices. The following methods are used in our experiments.

- Evaluation of time to build a model for each classifier.
- Mean Absolute Error (MAE):
- Root Mean Squared Error (RMSE)
- Kappa Statistics (KS)
- ROC curves. Additionally the AUC (Area Under ROC Curve) is taken under consideration.
- Confusion Matrix

3 Experimental Results

3.1 Data Set and Simulation of Hospital Environment

We simulated the environment of Baraha Medical City in Shambat, Khartoum North, Sudan using the framework reported in [23-24]. It is situated in a 600 Sq. meter lot with a garden within the compound. The hospital has five floors with a 75-bed capacity and provides complete medical services for patients. The Hospital receives patients who suffer from chronic diseases such as heart diseases, asthma, diabetes and abnormal blood pressure etc. Also people in post-surgery state needs continuous monitoring of their health condition, especially the vital signs, until their health status becomes stable. In our simulation, we allocated 6 chronic ill patients in each floor (total 30 patients) as we focused only on the monitoring and providing medical service for patients with chronic or terminally ill diseases. Depending on the critical condition of the patient, each patient was attached with several sensors. For thirty patients, there were a total of 300 readings at any measuring instant. Depending on the criticality of the patient's condition, when a sensor finds values that fall in the

danger zone an automated alarm is triggered notifying the nurses and doctors through mobile network or Wifi systems [23]. In this project, our main task is to develop a decision support system that could assist the hospital management to assess the situation of the hospital as Normal or Abnormal (too many medical emergencies) so that more medical help could be sorted.

We apply attribute selection method to reduce the number of the attributes. All 300 attributes were labeled as *A, B, C, Z, ...* and *KN*. We investigated several classifiers using WEKA [7] and finally managed to reduce to 6 attributes: *AK, CM, CP, CW, FJ* and *KN*. We found that cross-validation give the best classification with 10 Fold. Then the overall accuracy for all classifiers was done. We selected classifiers with classification accuracy between 90% to 100% as the proposed Base Classifiers. The Base Classifiers Proposed in our investigation in this paper are IBk, Attribute Selected Classifier, Bagging, Random Committee, PART, J48, LMT, Random Forest, Random Tree. The aim of this paper is to investigate the experimental results of the performance of different classification techniques for the simulation wearable sensors dataset. The performance factors used for analysis are accuracy and error measures. The accuracy measures are TP rate, F Measure, ROC area, Sensitivity and Specificity. The error measures are Mean Absolute Error, Root Mean Squared Error and Kappa Statistics. In the preprocessing step we have changed the class attribute to Abnormal or Normal where a 'Abnormal' specifies 1 class and a 'Normal' Specifies 0 class. Table 1 depicts the various error metrics analyzed in the data set. It is inferred from Table 1 that Random Tree has the least MAE and highest Kappa Statistic value. Random Tree is an appropriate model for classifying the hospital situation in a minimal span of time with higher accuracy.

Table 1. Performance Measures comparison

Algorithm	MAE	RMSE	KS	Correctly Classified
IBk	0.0978	0.3104	0.8062	673 90.3356 %
Attribute Selected Classifier	0.1008	0.2631	0.8384	685 91.9463 %
Bagging	0.1527	0.2609	0.8089	674 90.4698 %
Random Committee	0.0643	0.1931	0.9004	708 95.0336 %
PART	0.101	0.264	0.8355	684 91.8121 %
J48	0.0865	0.2518	0.8574	692 92.8859 %
LMT	0.0854	0.2454	0.844	687 92.2148 %
Random Forest	0.0961	0.219	0.8843	702 94.2282 %
Random Tree	0.051	0.2258	0.8977	707 94.8993 %

Table 2. Classifier performance in term of recall precision, *f* measure and false alarm rate

Algorithm	Recall	Precision	F-measure	False alarm rate
IBk	0.916	0.905	0.911	0.085
Attribute Selected classifier	0.914	0.914	0.913	0.076
Bagging	0.893	0.905	0.898	0.084
Random Committee	0.938	0.957	0.947	0.038
PART	0.924	0.908	0.914	0.080
J48	0.9185	0.9316	0.924	0.061
LMT	0.905	0.9316	0.917	0.062
Random Forest	0.927	0.951	0.938	0.044
Random Tree	0.9383	0.9544	0.945	0.041

As an example of classifier error illustration, Figure 1 depicts the Classifier error of Random Committee. The blue crosses indicate the Normal class and red crosses indicate the Normal class and squares indicate not classified. Table 2 depicts the classifier performance of each classifier in term of recall precision, *f* measure and false alarm rate. It is inferred from that Random Committee model has the highest precision and lowest false alarm rate, and the same recall as Radom Tree. Table 3 depicts the algorithm performance of each classifier in term of recall precision and *f* measure for Normal class is summarized. It is inferred from Table 3 that Random Committee model has the highest precision and also high recall. Figure 2 depicts the Area under ROC of Random Committee classifier with highest area under Roc. Tables 4 depict the classifier performance of each classifier in term of recall precision, and *f* measure for abnormal class is summarized. It is inferred from Table 4 that Random Committee model has the highest precision. Table 5 depicts the classification performance of each classifier in term of Sensitivity and Specificity with the Random Committee model having the highest Specificity and also high Sensitivity. Random Committee model also has the highest accuracy and the IBK model has the lowest accuracy.

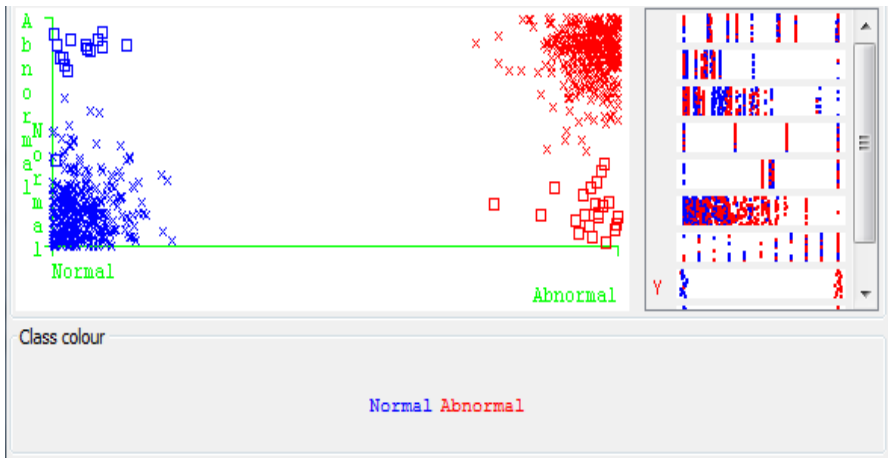
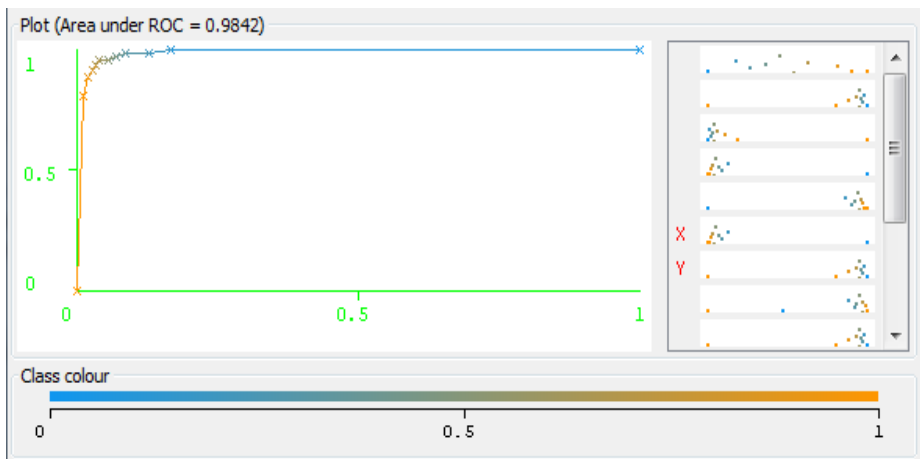


Fig. 1. Classifier error of Random Committee

Table 3. Classification performance for Normal class

Classifiers	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area
IBk	0.883	0.109	0.878	0.883	0.881	0.891
Attribute selected classifier	0.906	0.122	0.869	0.906	0.887	0.926
Bagging	0.880	0.112	0.875	0.880	0.878	0.953
Random Committee	0.943	0.071	0.922	0.943	0.932	0.984
PART	0.926	0.124	0.869	0.926	0.897	0.955
J48	0.903	0.109	0.881	0.903	0.892	0.934
LMT	0.886	0.094	0.894	0.886	0.890	0.948
Random Forest	0.937	0.084	0.909	0.937	0.923	0.973
Random Tree	0.932	0.074	0.919	0.932	0.925	0.929

**Fig. 2.** Area under ROC of Random Committee classifier**Table 4.** Classification performance of each classifier for abnormal class

Classifiers	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area
IBk	0.901	0.094	0.915	0.901	0.908	0.910
Attribute selected classifier	0.878	0.094	0.913	0.878	0.895	0.926
Bagging	0.888	0.120	0.893	0.888	0.891	0.953
Random Committee	0.929	0.057	0.948	0.929	0.938	0.984
PART	0.876	0.074	0.930	0.876	0.902	0.955
J48	0.891	0.097	0.912	0.891	0.901	0.934
LMT	0.906	0.114	0.899	0.906	0.903	0.948
Random Forest	0.916	0.063	0.943	0.916	0.929	0.973
Random Tree	0.926	0.068	0.938	0.926	0.932	0.929

Table 5. Classification performance of each classifier in term of Sensitivity and Specificity

Classifiers	Sensitivity	Specificity	Accuracy
IBk	0.8907	0.914	0.9033
Attribute Selected Classifier	0.941	0.923	0.9194
Bagging	0.8932	0.915	0.9046
Random Committee	0.938	0.961	0.9503
PART	0.924	0.92	0.9221
J48	0.918	0.938	0.9288
LMT	0.905	0.937	0.9221
Random Forest	0.927	0.955	0.9422
Random Tree	0.938	0.958	0.9489

4 Discussions

Empirical results indicate that the execution time of Random Committee algorithm is lowest for classification in comparison with the rest of classification algorithms, and the LMT algorithm has the higher execution time. The MSE error of the classification values for Random Committee is lower in comparison with the rest of the based proposed classifiers, and the Meta bagging classifier has higher MSE error in comparison with the rest of the base proposed classifiers. In terms of recall precision, f measure and false alarm rate the Random Committee model has the highest precision and lowest false alarm rate, and the same recall as Random Tree. In term of recall precision and f- measure for Normal class it is inferred that Random Committee model has the highest precision and also high recall. With higher true positive rate and minimum false rate also with higher ROC Area when the classification is Normal class in comparison of the rest of the classifiers. Attribute Selected Classifier has the lower precision in comparison with the rest. Also from the performance of each classifier in term of recall precision and f measure for abnormal class, Random Committee model has the highest precision and also high recall (with higher true positive rate and minimum false rate), also has highest ROC Area in comparison with other classifiers. While PART classifier has the lowest precision the same as Attribute Selected Classifier but with highest in ROC Area compare with Attribute Selected Classifier. From Sensitivity, Specificity and Accuracy perspective, the Random Committee model has the highest Specificity and also high Sensitivity the same as Random Tree but with highest accuracy of all the classifiers. While IBK classifier has the lowest Sensitivity, Specificity and Accuracy compare with the rest of the classifiers. To sum up, from the execution and accuracy point of view, Random committee model can be identified as the best choice for analysis and detection model among all the other classifier algorithms. Random committee provides an advantage that with a reduced feature set a better classification performance and is able to offer a better decision support system.

5 Conclusions

The main goal of this paper is to evaluate nine based classifier algorithms to develop a decision support system to classify the situation of an emergency hospital based on the Vital Signs from Wearable Sensors. We reduced the number of attributes from 300 attributes to 6 attributes. We explored and evaluated the models with various methods of evaluation based on Error Metrics, ROC curves, Confusion Matrix, Sensitivity and Specificity. We compared the performance of the entire classifiers and empirical results illustrate that Random committee classifier with selection attribute method gives better accuracy, error rate and reduced false alarm rate and with the highest Sensitivity and Specificity.

References

1. Salih, A., Abraham, A.: A Review of Ambient Intelligence Assisted Healthcare Monitoring. *International Journal of Computer Information Systems and Industrial Management (IJCISIM)* 5, 741–750 (2013) ISSN 2150-7988
2. Woo, S.-M., Lee, H.-J., Kang, B.-J., Ban, S.-W.: ECG signal monitoring using one-class support vector machine. In: *Proceedings of the 9th WSEAS International Conference on Applications of Electrical Engineering*, Penang, Malaysia, March 23-25, pp. 132–137 (2010) ISSN 2150-7988
3. Patel, S., Lorincz, K., Hughes, R., et al.: Monitoring motor fluctuations in patients with Parkinson's disease using wearable sensors. *IEEE Transactions on Information Technology in Biomedicine* 13(6), 864–873 (2009)
4. Korel, B.T., Koo, S.G.M.: Addressing Context Awareness Techniques in Body Sensor Networks This paper appears. In: *Advanced Information Networking and Applications Workshops (2007)*, doi:10.1109/AINAW.2007.69, ISBN: 978-0-7695-2847-2
5. Fleury, A., Vacher, M., Noury, N.: SVM-Based Multi-Modal Classification of Activities of Daily Living in Health Smart Homes: Sensors, Algorithms and First Experimental Results. *IEEE Transactions on Information Technology in Biomedicine* 14(2), 274–283 (2010)
6. Kim, N.J., Hong, J.H., Cha, E.J., Lee, T.S.: Classification Technique of Human Motion Context based on Wireless Sensor Network (2005)
7. Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., Witten, I.H.: *The WEKA Data Mining Software: An Update*. *SIGKDD Explorations* 11(1) (2009)
8. Quinlan, R.: *C4.5: Programs for Machine Learning*. Morgan Kaufman, San Mateo (1993)
9. Danham, M.H., Sridhar, S.: *Data mining, Introductory and Advanced Topics*, 1st edn. Person Education (2006)
10. Sharma, A.K., Sahni, S.: A Comparative Study of Classification Algorithms for Spam Email Data Analysis. *IJCSE* 3(5), 1890–1895 (2011)
11. Landwehr, N., Hall, M., Frank, E.: Logistic Model Trees. *Machine Learning* 95(1-2), 161–205 (2005)
12. Leo, B.: Random Forests. *Machine Learning* 45(1), 5–32 (2001)
13. Ush Rani, M., Prasanna Kumari, G.T.: A Study of Meta-Learning in Ensemble Base Classifier. *IRACST, Engineering Science and Technology: An International Journal (ESTIJ)* 2(1) (2012) ISSN: 2250–3498

14. Chan, P.K., Stolfo, S.J.: Experiments on multi-strategy learning by meta-learning. In: Bhargava, B.K., Finin, T.W., Yesha, Y. (eds.) *CIKM*, pp. 314–323. ACM (1993)
15. Breiman, L.: Bagging predictors. *Machine Learning* 24(2), 123–140 (1996)
16. Witten, I.H., Frank, E.: *Weka machine learning algorithms in java*. In: *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*, pp. 265–320. Morgan Kaufmann Publishers (2000)
17. Aha, D., Kibler, D.: Instance-based learning algorithms. *Machine Learning* 6, 37–66 (1991)
18. Altman, N.S.: An introduction to kernel and nearest-neighbor nonparametric regression. *The American Statistician* 46(3), 175–185 (1992)
19. Hall, M.A.: *Correlation based Feature Selection for Machine Learning* (1999)
20. Autio, L., Juhola, M., Laurikkala, J.: On the neural network classification of medical data and an endeavor to balance non-uniform data sets with artificial data extension. *Computers in Biology and Medicine* 37(3), 388–397 (2007)
21. Berry, J., Linoff, G.: *Data Mining Techniques: For Marketing, Sales, and Customer Relationship Management*, 2nd edn. Wiley, Indianapolis (2004)
22. Fraile, J., Bajo, J., Corchado, J., Abraham, A.: Applying wearable solutions in dependent environments. *IEEE Transactions on Information Technology in Biomedicine* 14(6), 1459–1467 (2010)
23. Corchado, J., Bajo, J., Abraham, A.: GerAmi: Improving the delivery of health care in geriatric residences. *IEEE Intelligent Systems* 23(2), 19–25 (2008)

An Efficient Segmentation Algorithm for Arabic Handwritten Characters Recognition System

Mohamed A. Ali

Computer Science Dept., Faculty of Science, Sebha University, Sabha, Libya
fadeel1@Sebhau.edu.ly

Abstract. If the pre-processing phase, in optical character recognition systems, is the heart of the recognition process, the segmentation stage is the “aorta” of this heart. This paper introduce a reliable segmentation technique for Arabic handwritten script. Number of techniques like; script height, character width, pen thickness and word/subword gaps are used to design an efficient segmentation algorithm. the algorithm performs diacritics removal, word/subword segmentation, ascender characters segmentation, descender characters segmentation and finally embedded characters segmentation.

Keywords: Off-line optical character recognition, Arabic character recognition, segmentation algorithm, coarse and fine segmentation, handwritten character recognition.

1 Introduction

In most optical character recognition system, failing to segment a handwritten text correctly will definitely result in poor recognition, no matter how well the following and previous stages are designed. It is very obvious that a considerable share of recognition errors is attributed to the segmentation stage. For this reason, many researchers have concentrated their researches in improving segmentation methodology [1]. In fact, development of new segmentation algorithm is one of our commitments for contribution in the field of Arabic handwritten characters recognition.

Character recognition of off-line Arabic handwritten script remains a challenging problem in pattern recognition compared with that of machine printed [8].

2 Methodology

In this paper number of techniques have been developed to assist the process of diacritics, word, subword and characters segmentation;

2.1 Diacritics Removal

One of the main features of Arabic script is existence of diacritics [6]. There are six types of diacritics, divided in terms of prerequisites into two groups; the first group

includes dots and short-Alif (Fig. 1-a), diacritics of this group are crucial for read and comprehend Arabic script properly. The second group includes diacritics like Hamza and Madda shown in Fig. 1-b, without which an expert reader can easily read and comprehend an Arabic script using context hints help.

Presence of diacritics may affect the estimation of text-line height as shown in Fig. 2-a where h_r and h_d are the real and distorted heights of text-line respectively. Fig. 2-b and Fig. 2-c show the presumable points of segmentation with and without presence of dots respectively.

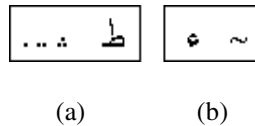


Fig. 1. Diacritics strokes in Arabic writing

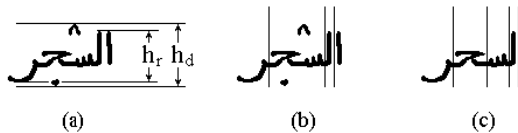


Fig. 2. Effect of presence of dots on segmentation process

Semi-Alif above a loop (e.g. Tta ط), in a main loop, which may or may not touch the loop. If the semi-Alif does not touch the loop, as shown in Fig. 1a, it will be considered as a diacritic and the algorithm will remove it otherwise it will be considered as a part of the main stroke and will not be removed.

In fact, Presence of diacritics in Arabic script, in both machine printed as well as handwritten form, make achievement of an efficient segmentation and classification processes a hard job. Hence, we decided to temporary remove these diacritics till the recognition stage where we shall retrieve them back for efficient recognition.

In fact, we use the Algorithm illustrated in section 2.5 for segmentation of diacritics as well as words and subwords, the only difference in case of diacritic segmentation is that the segmented diacritic along with its attributes (includes their numbers, types and positions) is directly applied to classification and recognition stages and no further segmentation or thinning processes are applied.

2.2 Character Width Estimation

Estimation of suitable standard width is an essential choice and it is not a simple task. It is necessary to know the character width in order to validate the location of the candidate segmentation/end points. Thinking about the problem of finding good character width estimation has lead us to the following approaches; the first approach involves performing statistical studies of an Arabic document in order to get an approximate measure of the average character width. However, this approach depends

heavily on the training set and can produce undesirable results if a different font size is being used. Second approach uses the width of whole word and the average number of characters in a word to estimate the character width. This approach can produce unpredictable results if the number of characters is far away from this average. The third approach, which we provide in this paper, takes the average height of the character Alif as a reference to estimate the character width. The main merit of this approach is that, it locally estimates the character width, in other word, it estimates the character width in a word reside in that particular text-line.

We have made statistic on character width of isolated printed Arabic characters (at 12 pt.) as shown in Fig. 3, and found that the least width for Alif is (1 pixel) whereas the highest is for Ssad (16 pixels). Although this statistic is for printed form of characters yet we could use it to estimate the average character width in handwritten form.

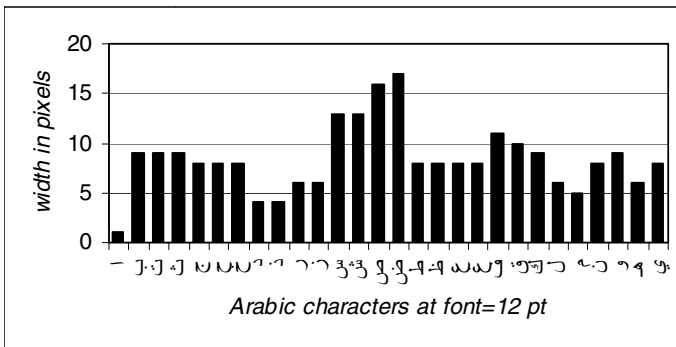


Fig. 3. Arabic characters vs. their widths

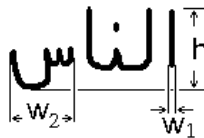


Fig. 4. Character width estimation

Since a writer may scale up or down the characters width depend on his handwriting. The calculated width average is around 9 pixels. We notice also that the height of Alif is 10 pixels (at font 12 pt.). Now, since the average width is almost the same as that of Alif height, we have chosen the height of Alif as a reference for estimating the average width of characters in handwriting. Experimentally, we found the character width varies from few pixels w_1 in case of “Alif” to as many as the height of Alif which roughly equals to $2/3$ of text-line height ($w_2 \approx 2h/3$, h : text-line height) w_2 as in case of “Seen” in the same word shown in Fig. 4. Word-to-character segmentation depends heavily on good estimation of character width.

2.3 Pen Thickness Estimation

Estimation of stroke thickness is one of the pillars that our segmentation algorithm is based on. We have developed a simple procedure that can estimate the thickness of stroke. The procedure starts by scanning the image of text-line from top right to bottom left looking for the widest rectangle which implicitly has the widest word or subword, once that rectangle is detected we vertically scan it from right to left in a (vertical histogram), the black pixels in each column are counted with two conditions; it stops counting if the next pixel is white or if the count is greater than 20% of the text-line height which has been calculated earlier. Putting these conditions serves two purposes. First, it saves scanning time since the scanning will be, in most cases, from the top of the rectangle up to the base-line area where most of the characters reside, and in some cases even before that in case of scanning through ascending characters (e.g., Alif, Lam Kaf $\text{ا} \text{ ل} \text{ ك}$) where the second condition is fulfilled. Second, to make certain that we don't add up black pixels in two or more segments in the same column like in the case of Jeem (ج) and hence false estimation is obtained. Now after we scan the whole rectangle we end up with number of black pixels/column values ranging from one pixel (in case of noise or edge point of stroke) to number of pixels less than or equal to 20% of the text-line height. Final step is to find the most frequent value and store it as a pen thickness. Pen thickness can be used in determination of whether a gap between two successive characters (*inter-characters*) is deliberately there or it is due to noise or discontinued writing flow. If thickness of a stroke is P_{th} pixels then after thinning process we can say that the output medial skeleton (one pixel width) has lost $(P_{th}-1)/2$ pixels from both sides of its contour of its original shape before thinning. Now if we have two adjacent strokes which are separated by a gap approximately equal to the pen thickness, then this gap will be widened (as a result of thinning) by as much as twice of its original space. Hence our algorithm task is to join these two or any strokes separated by a gap $\leq 2 \times P_{th}$. Pen thickness was used here as a reference dimension to distinguish among word/subword, a character Alif or a diacritic. Stroke thickness is also used as a key feature (knowledge source) in classification and recognition stages.

2.4 Gaps in Text-Line

Gaps in a text-line can be between subwords (*inter-subword*) and between words (*inter-word*), Fig. 5. In Arabic script (printed/handwritten) *inter-word* gaps and *inter-subword* gaps are remarkably unequal. Generally, gaps between words are more spacious than that between subwords. This is always true in case of machine-printed text. In handwritten script, however, the *inter-word* gaps and *inter-subword* gaps are interchangeably irregular. There are two salient aspects to our approach: (i) it is not dependent on a large set of training data collected across a variety of writers, and (ii) the metric measuring the magnitude of gaps between components is statistically computed using the features of the phrase image under examination.

We have developed a simple procedure by which we could distinguish between the *inter-word and inter-subword* gaps. The procedure starts by applying vertical histogram to the text-line currently under process. The width of the gaps (0-pixels zone) between words and subwords are recorded as GW_1, GW_2, \dots, GW_n . Now we first find the minimum and maximum values of GW and save them as GW_{\min} and GW_{\max} respectively.

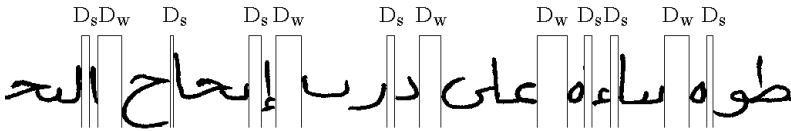


Fig. 5. Inter-words and inter-subwords gaps in a Text-Line

We found experimentally that the distance between words (*inter-word*) D_w and the distance between subwords (*inter-subword*) D_s (as shown in Fig. 5) are governed by the following statements:

if $d = GW_{\max} - GW_{\min}$

and $\Delta d = (GW_{\max} - GW_{\min})2/3$

$$\text{then } GW_{\max} \geq D_w \geq (GW_{\max} - \Delta d) \tag{1}$$

$$(GW_{\max} - \Delta d) > D_s \geq GW_{\min} \tag{2}$$

Where d is the difference between GW_{\max} and GW_{\min} , Δd is the save margin (that we should add to GW_{\min} or subtract from GW_{\max}) in order to fix a width range for D_w and D_s . The save margin is found experimentally. The only condition here is that GW_{\min} should be equal to or greater than “Pen thickness”, if this condition is not true then we need to make $GW_{\min} = \text{Pen thickness}$ first, when we calculate d and Δd , then we proceed to equations 1 and 2. Putting this condition here prevent the system from inter-subwords underestimation. Underestimation occurs when GW_{\min} tends to become one pixel width which, in turn, makes d almost equal to GW_{\max} , then D_s become even smaller. If this happened then the probability of considering an inter-subword as an inter-word is higher. In other words, numbers of subwords are mistakenly considered as words which makes the task of recognition is even harder

The locations of gaps in the text-line are recorded. We call the information regarding gaps locations and whether these gaps are inter-subword or inter-word as *gap knowledge* source and it is denoted by GP-KS. GP-KS later participates in matching and recognition as we shall see in recognition part of this research.

2.5 Fine Segmentation

Fine segmentation is a further step in segmentation process [10], where the output of the coarse segmentation stage (simply word or subword) is analysed and searching for character segmenting points is performed.

2.5.1 An Overview of the Proposed Algorithm

Our approach is based on two main steps; segmenting ascending and descending characters first followed by segmenting the rest of characters (embedded) using vertical histogram. Ascending, descending and embedded characters are explained in the next section. We have introduced certain measurements that would assist in locating segmentation points like; stroke width estimation, word height estimation and average character width estimation. The algorithm, initially, start segmenting the ascending and descending characters due to three significant reasons; first, because ascending characters may prevent vertical histogram from detecting an underneath segmentation point, like in case of Kaf and Nabira (ك ن) shown in Fig. 6 or assign a segmentation points at a level much higher than base-area level. Second, it eases the detecting and segmenting of characters in ligature. Third, it might leads in some cases to early character recognition during segmentation and hence the overall time cost is reduced.

Since we rely mainly on structure features (topology) of Arabic characters for our segmentation algorithm, therefore there will be less mathematical implementation which means less program complexity.

2.5.2 Structural Properties of Arabic Script

Arabic script has properties which make it unique compared with the Latin, Chinese and Japanese script. Some of these properties facilitate the task of segmentation as we shall see later in this paper. However, some other properties clearly show why the Arabic characters segmentation is considered as one of the toughest job in optical character recognition systems in general. The structural properties upon which our fine segmentation algorithm is based can be illustrated as follow:

- 1) Text-line of Arabic script (both printed and handwritten) can be divided horizontally into three zones; upper, lower and middle zones as shown in Fig. 6. Middle zone is the base-area [4], [8]. Basic characters shapes of Arabic script can fall in one, two or sometime even three of these zones. The Basic characters shapes (BCS) of Arabic script can be accordingly clustered into four classes as to the zones they cover.

Class-I:BCSs cover only the middle zone of a word such as (ك ن ف م ه ل) (ك ن ف م ه ل). We call these BCS as embedded characters since they are confined by the upper and lower limits of the base-area.

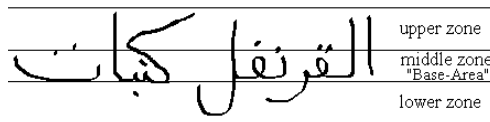


Fig. 6. Three zones specified for Arabic handwriting

Class-II:BCSs cover both the upper zone and the middle zone of base-area; we call these characters as ascending characters (e.g. ط ك ل).

Class-III: Basic characters shapes that cover both the lower zone and the middle zone of base-area; we call these characters as descending characters (e.g. ح ن م و ي).

Class-IV: BCS that covers all three zones (e.g. ل). In Arabic cursive scripts, several characters of a word are usually merged to form a connected component, with each character touching or connecting with its right (precedent) and left (subsequent) adjacent characters, as a result all Arabic characters may take different shapes according to its position in the word. This is particularly true with characters that may have a right, left or both sides connection with other characters. This structure is a rather common feature of Arabic characters.

- 2) Second feature in Arabic script states that two characters may touch each other (e.g. كك), this feature however put more constrains on characters segmentation process.
- 3) The important and useful property of Arabic script (printed/ handwritten) is that the existence of one of the following characters (ك ه م ن ل ك ق ع ط ص س ح ت) gives absolute indication that this position is end of current word (not subword) and beginning of new word. Hence this property, in particular, is very useful in our segmentation and classification processes.
- 4) Some character can be indistinguishable from a three connected characters(*trigraph*) at the image level, for instance, the letter 'س' and the trigraph 'تينا' have the same basic shape 'س' after diacritics removal. In order to deal with this ambiguity, in most handwriting recognition systems the segmentation algorithm tries to find all the real letter boundaries with as few extraneous ones as possible; then, the recognition system finds the best matching word in the lexicon. In our fine segmentation algorithm the problem of ambiguity has been, alternatively, resolved by using the basic characters shapes.
- 5) Embedded Characters (printed/handwritten) like ح ت ف و م ه ه ل ن ح ن س ص, present almost the same height, although this varies (in case of handwriting) from person to another.

3 Segmentation Model

The input to this phase of segmentation, as we mentioned earlier, is a single subword, image coming from the previous phase, with all its diacritics been removed. The significance of diacritic removal has been introduced in section 2.1, and hence the input of this stage is a word consisting of basic shape of Arabic alphabets.

The next step is *sandwich scanning* which is a novel technique we have introduced in this paper. To the best of our knowledge this technique has not been used before at least in Arabic OCR system. Sandwich scanning is a method in which a scanning is performed in two areas above (upper zone) and below (lower zone) the

base-area. Now the scanning in the upper zone as well as in lower zone is started from top-right to bottom left corner of each zone. In other word, the scanning in the upper zone starts from the top-right coordinate of the subword image to the upper-left coordinate of the base-area and, likewise, the scanning in the lower zone starts from the bottom-right of the base-area to the bottom-left of the subword image. If a pixel is detected in the upper zone then we proceed for *ascending characters* extraction algorithm otherwise we proceed to the lower zone searching for a pixel (i.e. descending character). Likewise, if a pixel is detected in the lower zone then we proceed for *descending characters* extraction algorithm otherwise we proceed to the base-area zone where we can apply the *embedded characters* segmentation. While either of extraction algorithms is in progress, ascending and descending characters are truncated and saved in separate images leaving behind only characters inside the base-area (embedded characters) which, in further step, will be segmented using vertical histogram and saved in separate images as we shall see in the next section.

Performing this type of scanning technique guarantee the proper detection of ascending and descending characters according to their sequence as the writing goes from right to left. We have used this technique to detect and segregate ascending and descending characters first, then we detect and segregate the embedded characters. Moreover, using this type of scanning guarantee that a word is sequentially segmented into characters from right to left, so that in recognition stage the same sequence is used to retrieve and recognize the characters being segmented in this stage.

3.1 Ascenders and Descenders

Existence of both ascending and descending characters (from now on we call them *ascenders* and *descenders* respectively) in a word complicates segmenting that word into its constitutional characters using conventional vertical histogram technique [13]. For instance, the word (كنبات) in Fig. 6 has its ascender 'ك' overshadow the embedded characters 'ن' and 'ب' so if we apply the vertical histogram, the first minima will occurs after the character 'ب' which is not desired segmentation point. Therefore we had to find a mean to segment that ascender 'ك' first then we go for segmenting the other characters. Empirically it has been noticed that the main body of a word is not uniform in height over the entire width of the word. As an example, one can consider the word shown in Fig. 6, the letter 'ر' in that word is almost entirely below the base-area so we call this letter as descender, while the letter 'ا' is emerging above the top reference line of base-area so we call this letter as ascender. Obviously we do not want some letters like 'ق' or 'ف' to be falsely considered as either ascenders or descenders. Consequently we need to set some threshold values for detection of ascenders and descenders. These thresholds have been fixed experimentally and are expressed as a percentage of the main text-line height [4], [9]. The two thresholds are; ascending area threshold A_{th} and descending area threshold D_{th} .

3.2 Segmentation of Ascenders

A segmentation model was designed for segmenting the ascenders of the Arabic characters shown in Fig. 7.

According to our assumption, the detection of an ascender should occur in the upper zone of the word/subword currently under process. The model moves from the top-right to bottom-left in the upper zone until it finds a black pixel, from first detected black pixel the model will follow all connected pixels in upper zone as well as those in base-area. Now to avoid character over segmentation we put two conditions; tracing of pixels will stop (assign segmentation point SP) if the number of successive pixels in direction-4 (according to Freeman code) reach $0.2 * \text{height}$ of the word or if the pixels start changing its moving direction from direction-4 to either of direction-2, 3, 5 or 6. Once one of these conditions is fulfilled and the current column has no more pixels in direction-6, the traced connected pixels are cut and saved in separate array and named as ascender in the same manner we followed in coarse segmentation. The same steps were followed again to detect and segment other ascenders if any. Once all ascenders are segregated from the main word/subword the next step is descenders segmentation.



Fig. 7. The basic shapes of ascenders

3.3 Segmentation of Descenders

One of the most remarkable features in Arabic script states that all descenders are always come either in isolated form or end-letter form and are not connected to their succeeding (left) character. This feature in particular had significantly helped us in designing an efficient descenders segmentation model. This model is similar to that for the ascenders segmentation model. The descenders (in their basic shapes) are Ya, Waw, Haa, Noon, Meem, Lam, Ain, Ssad, Seen, Ra and Hha as shown in Fig. 8.

A model was designed to segment descenders from the word currently under segmentation process. The module will use similar scanning algorithm to that used in case of ascenders segmentation, scanning in this case is from bottom-right corner of lower zone to top-left corner of the same zone. The module will try to detect first black pixel it came across in the zone. If the first detected pixel is at P1, then from P1 the module will trace the descender stroke/s connected to P1 in all directions looking for three successive pixels in direction-0 and located inside the base-area, and each of them on a column of one pixel height, if this condition is fulfilled then the far east limit of the character is found, the module, on the other hand, will also search for three successive pixels in direction-4 and located inside the base-area, and each of them on a column of one pixel height, if this condition is fulfilled then the far west limit of the character is found and it is most probably a character of 'هـ', however if the later conditions are not fulfilled then we will make use of the significant feature

we mentioned in the beginning of this section which state that all descenders are always come either in isolated form or end-letter form and are not connected to their succeeding left-character, therefore the west limit of the whole word/subword will be taken as the west limit of the current descender. The question here is what about characters like seen 'س' and saad 'ص' where there is a possibility that only their cavities are segmented i.e. 'س', actually it does not really matter because this mis-segmentation can be compensated in recognition module where diacritics are used to assure whether a cavity like this belongs to character like س, ش, ص, ض or ن. Finally we use the same naming strategy used in ascenders segmentation modules to name the output subimages of descenders segmentation module. Fig. 9 illustrates the segmentation of descender 'و' from the subword 'سندکُو' .

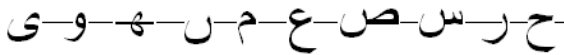


Fig. 8. The basic shapes of ascenders

3.4 Segmentation of Embedded Characters

Finally it is time to segment connected embedded characters into individual characters. Embedded characters are given in section 2.5.2 (class-I). The technique used for embedded characters segmentation is vertical projection. The vertical projection is obtained by calculating total runs of vertical pixels for each column of the image where black pixels exist. Firstly, a vertical Projection of all pixels in base-area is created. Now since we are dealing with a thinned image then the minimum value of vertical projection is one pixel high. Next we examine all columns that have more than one pixel (and comes immediately after column of one pixel height) and mark them as possible segmentation point PSP. Locations of PSP indicate the beginning of embedded characters, however as we already explained the embedded character may be connected to either or both of its adjacent characters by 'Maddah'. Maddah is expected to be between any two successive characters. A midpoint between any two successive PSP is assigned as segmentation point SP. One can ask what about character basic shape like 'س' which has four SPs, so it might be one character 'س' or three different characters like 'سنت' ? The answer is yes this character is segmented into three characters basic shapes at this stage and then it is the recognition module responsibility to decide whether it is one or three characters depending on the presence of diacritics above or below this/these characters.

The segmentation process is repeated until the full length of the line is exhausted. Fig. 9 gives a good understanding of the entire process. All segmented characters are saved in separate files, the name of each file includes x-y coordinates indicating a position on text-line from where this character was segmented and the character 'E' at the end of file name to indicate that this image file contains a character of type 'Embedded'.

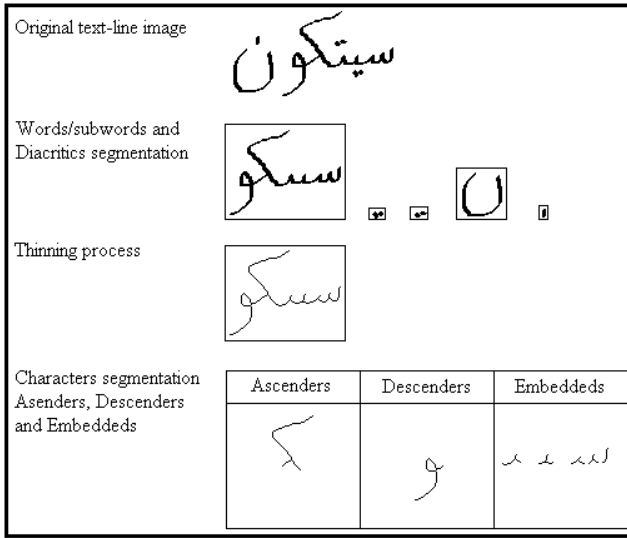


Fig. 9. Fine Segmentation Steps

4 Conclusion

Numbers of handwritten samples are used to test the different types of segmentation modules given in this section, i.e. text-lines segmentation, coarse segmentation and fine segmentation. Samples vary from having very neat to very bad Arabic handwriting script. Text-lines and coarse segmentations scored higher accuracies compared with fine segmentation. Segmentation-recognition of Alif and diacritics segmentation are very easy and effective processes. Proper base-area estimation is one of the main factors which tremendously affect fine segmentation. Structural Properties of Arabic script as well as other knowledge sources are used effectively in segmentation process. Other knowledge sources are extracted during segmentation process e.g. character coordinates with respect to base-area inter-words and inter-subwords gaps, and pen thickness are effectively utilized, hence, a very satisfied results were obtained. Very neat segmentation process is promising for a good characters recognition stage.

References

1. Motawa, D., Amin, A., Sabourin, R.: Segmentation of Arabic Cursive Script. In: Proc. of the 4th Inter. Conf. on Doc Analysis and Recognition, vol. (2), pp. 625–628 (1997)
2. Amin, A., Mari, J.: Machine Recognition and correction of printed Arabic text. IEEE Trans. Syst. Man Cybern. SMC 19(5), 1300–1306 (1989)
3. El-Gowely, K., El-Dessouki, O., Nazif, A.: Multi-phase Recognition of Multi-font Photo-script Arabic Text. In: Proc. 10th Conf: on Pattern Recognition (1990)

4. Al-Yousefi, H., Udpa, S.: "Recognition of Arabic characters", IEEE Trans. Patt. Analysis Machine Intell. PANH (1992)
5. Elkhaly, F., Sid-Ahmed, M.A.: Machine Recognition of Optically Captured Machine Printed Arabic Text. Pattern Recognition (1990)
6. Azeem, S.A., Ahmed, H.: Effective Technique for the Recognition of Writer Independent Off-Line Handwritten Arabic Words. In: Int. Conf. on Frontiers in Handwriting Recognition (ICFHR), pp. 594–599 (2012)
7. Ali, M.A.: Base-Area Detection and Slant Correction Techniques Applied for Arabic Handwritten Characters Recognition Systems. In: International Conference on Artificial Intelligence and Pattern Recognition (AIPR 2009), Orlando, USA, pp. 133–138 (2009)
8. Bushofa, B.F.M., Spann, M.: Segmentation of Arabic Characters Using their Contour Information. In: Proc. 13th DSP 1997, Santorini, Greece, vol. 2, pp. 683–686 (1997)
9. Cheung, A., Bennamoun, M., Bergmann, N.W.: An Arabic optical character recognition system using recognition-based segmentation. Pattern Recog. 34, 215–233 (2001)
10. Shukla, M.K., Banka, H.: An Efficient Segmentation Scheme for the Recognition of Printed Devanagari Script. IJCST 2(4), 529–533 (2011)
11. Lawgali, A., Bouridane, A., Angelova, M., Ghassemlooy, Z.: Handwritten Arabic Character Recognition: Which Feature Extraction Method. International Journal of Advanced Science and Technolog 34, 1–8 (2011)
12. Broumandnia, A., Shanbehzadeh, J., Rezakhah Varnoosfaderani, M.: Persian/arabic handwritten word recognition using M-band packet wavelet transform. Image and Vision Computing Archive 26(6), 829–842 (2008)
13. Elzobi, M., Al-Hamadi, A., Al Aghbari, Z., Dings, L., Saeed, A.: Gabor Wavelet Recognition Approach for Off-Line Handwritten Arabic Using Explicit Segmentation. In: Choraś, R.S. (ed.) Image Processing and Communications Challenges 5. AISC, vol. 233, pp. 241–250. Springer, Heidelberg (2014)

Graph Drawing Using Dimension Reduction Methods

Tomáš Buriánek¹, Lukáš Zaorálek¹, Václav Snášel^{1,2}, and Tomáš Peterek²

¹ VSB - Technical University of Ostrava, FEECS, Department of Computer Science,
17. listopadu 15, 708 33, Ostrava-Poruba, Czech Republic
{tomas.burianek.st1, lukas.zaoralek, vaclav.snasel}@vsb.cz

² IT4innovations, VSB - Technical University of Ostrava, Ostrava 708 33,
Czech republic
tomas.peterek@vsb.cz

Abstract. Graphs are common data structures in computer science used to capture relations between set of objects. Graph drawing is a visual representation of the graph in more readable form, usually with vertices projected into \mathbb{R}^2 space. Dimension reduction methods are designed to transform original data in high-dimensional space into a new data in lower-dimensional space. This work is focused on selected dimension reduction methods and their use to obtain sensible graph drawing from the original graph. Results of the dimension reduction methods are discussed and compared to each other and also to the classical approach presented by Kamada and Kawai.

Keywords: Graph Drawing, Kamada and Kawai, Dimension Reduction, Principal Component Analysis, Stochastic Neighbor Embedding, Factor Analysis, Diffusion Maps.

1 Introduction

Today, there is a wide area of graph drawing such as cartography, social networks or bioinformatics. Graph drawing has origins in geometric graph theory, mathematical theory and visualization of informations. These applications can bring graphs containing thousands of vertices and edges for various analysis (e.g. data mining, clustering). Ordinary force-directed graph drawing methods have its limitations. The basic Kamada and Kawai algorithm is good for small graphs. But for bigger graphs with more than hundred of vertices it results to poor drawings. The main reason of worse results for bigger graphs is that physical models has many local minimums where method can stuck and stop improving its result [19]. This paper proposes method for graph drawing based on dimension reduction of the original graph representation. Task of dimension reduction may be performed by many existing methods. Each dimension reduction method has its advantages and drawbacks and each has different results of graphs drawings. This has the benefits in the variability of the method. Appropriate dimension reduction method may be selected by suitability to the task of graph drawing.

This paper has the following structure: First, theory of graphs and graph drawing is introduced. Then problem of dimension reduction is described followed by description of four selected methods for dimension reduction in more detail. In the third section, proposed method for graph drawing using dimension reduction methods is presented. Then experiments and results performed by selected dimension reduction methods are discussed with final conclusion.

2 Graph Drawing

Graph is an abstract data structure used to model the relations and processes in a connected network. Example of use is in physical, biological, social and information systems. In this paper, the main focus is on undirected unweighted graphs. Undirected graph is represented as an ordered pair $G = (V, E)$, where elements of the set $V = \{1, \dots, n\}$ are the *vertices* and the set $E = \{e_1, \dots, e_m\}$ consists of *edges* and sets V and E are finite. Each edge e_j is represented as a pair (u, v) where $u, v \in V$. This describes edge related with two vertices as a connection between them. Graph theory is discussed in more details in: [10][13].

Graph drawing is part of combinatorial optimization problems [9]. As described in [5], graph drawing is a mapping of each vertex $i \in V$ from given graph G to point $p_i \in P$, where $P = \{p_1, \dots, p_n\} \subset \mathbb{R}^2$. Each edge (u, v) of graph G described by two vertices is mapped to the straight line connecting points p_u and p_v . To obtain positions of vertices in a plane, many different graph drawing methods may be used such as: Force-based algorithms [17][12][11], Spectral algorithms [20], Tree layout algorithms [15] and Orthogonal layout methods [18].

In the following subsection, force-based spring model introduced by Kamada and Kawai will be explained.

2.1 Kamada and Kawai Spring model

In [17] Kamada and Kawai presented their algorithm for graph drawing, where the graph is represented as a dynamic system, where points $p_1, \dots, p_n \in \mathbb{R}^2$ (here are presented as particles) are connected by springs. Good distribution of vertices is associated to the dynamically balanced spring system. Finally, the amount of imbalance is obtained as the total energy of springs presented in equation:

$$E = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{1}{2} k_{ij} (\|p_i - p_j\| - l_{ij})^2, \quad (1)$$

where l_{ij} is the original length of the spring between particles p_i and p_j corresponds to the desirable length between them in the drawing and is obtained as follows:

$$l_{ij} = \frac{L}{d_{ij}}. \quad (2)$$

The distance d_{ij} between two vertices v_i and v_j in a graph is defined as the length of the shortest path between them in the graph. L is the desired length

of a single edge in the projected plane and it is determined by distance of the farthest pair in the graph:

$$L = \frac{L_0}{\max_{i < j} d_{ij}}, \quad (3)$$

where L_0 is side length of the projected square area. The variable k_{ij} is the spring strength between particles p_i and p_j and is obtained from equation:

$$k_{ij} = \frac{K}{d_{ij}}, \quad (4)$$

where K is a constant strength. The desired graph drawing is obtained by minimizing the energy in a spring model described in equation (1). Further detail is described in [17].

3 Dimension Reduction

The problem of dimension reduction for a set of variables $X = \{x_1, \dots, x_n\}$ where $x_i \in \mathbb{R}^D$ is to find a lower dimensional representation of this set $Y = \{y_1, \dots, y_n\}$ where $y_i \in \mathbb{R}^d$ and where $d < D$ (often $d \ll D$) in such a way to preserve content of the original data as much as possible [28][7]. Dimension reduction methods can be used as a visualization tool to show multi-dimensional data in \mathbb{R}^2 or \mathbb{R}^3 space that is more convenient for humans. The high-dimensional data is processed to extract only important features by reducing redundant components to obtain precise representation. It is also an important task in data preprocessing for analysis in the classification [26].

In the following subsections, several dimension reduction methods are discussed.

3.1 Principal Component Analysis

The Principal Component Analysis (PCA) [1] extracts the most important information from dataset. In other words, the PCA computes variables principal component, which are linear combination of the original variables. Also the PCA reduces size of the original feature space [29]. Let matrix X represents graph dataset (connections of vertices in graph with each other), then:

$$Z = XA. \quad (5)$$

In this context, the matrix Z consists of principal components which are linear combinations of columns of X and the matrix A is an orthogonal matrix. In more details, let $x_1, x_1, \dots, x_n \in \mathbb{R}^p$. Reconstruction of data in \mathbb{R}^q to \mathbb{R}^p is defined as follows:

$$f(\lambda) = \mu + v_q \lambda \quad (6)$$

where $\mu \in \mathbb{R}^p$ is the mean, v_q is a $p \times q$ matrix with q orthogonal unit vectors and $\lambda \in \mathbb{R}^q$ is the low-dimensional data points.

3.2 Stochastic Neighbor Embedding

Stochastic Neighbor Embedding (SNE) is a probabilistic approach to place objects with high-dimensional representation to lower-dimensional space regarding to preserve neighbor identities. This method do approximation of probability distribution of neighboring points in the high-dimensional space with their probability distribution in a lower-dimensional space [16].

The asymmetric probability in high-dimension q_{ij} that object i picks as its neighbor j from set S is done by equation with Gaussian neighborhoods:

$$q_{ij} = \frac{\exp(-d_{ij}^2)}{\sum_{k \in S, k \neq i} \exp(-d_{ik}^2)}, \tag{7}$$

where d_{ij}^2 is dissimilarity between two points x_i and x_j in high-dimensional space. It may be described in problem definition or may be computed using the scaled Euclidean distance:

$$d_{ij}^2 = \frac{\|x_i - x_j\|^2}{2\sigma_i^2}, \tag{8}$$

where σ_i is a variance and could be set by hand or found by binary search. Lower-dimension probability r_{ij} that point i picks as its neighbor point j is obtained by equation:

$$r_{ij} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \in S, k \neq i} \exp(-\|y_i - y_k\|^2)}, \tag{9}$$

where the dissimilarity between points y_i and y_j is with fixed variance (set to $\frac{1}{2}$). Finally the aim of dimension reduction in SNE is to match probability distributions q_{ij} and r_{ij} as much as possible. It is done by minimizing cost function which is sum of Kullback-Leiber divergences [21] between these two distributions over neighbors for each object from set S :

$$C = \sum_{i,j \in S, i \neq j} q_{ij} \log \frac{q_{ij}}{r_{ij}} = \sum_{i \in S} KL(Q_i || R_i). \tag{10}$$

Minimization could be done by gradient descent approach using this function:

$$\frac{\partial C}{\partial y_i} = 2 \sum_{j \in S, j \neq i} (y_i - y_j)(q_{ij} - r_{ij} + q_{ji} - r_{ji}). \tag{11}$$

3.3 Factor Analysis

Factor Analysis is a linear method for analyzing the relations through a set of observed variables. Main objective is to reveal some common factors in a set of variables. For observed sample $x = (x_1, \dots, x_D)$ with mean $\mu = (\mu_1, \dots, \mu_D)$ there is unobserved common factor $f = (f_1, \dots, f_d)$ where $d < D$ satisfying condition:

$$x_i = \sum_j^d l_{ij} f_j + \mu_i + \epsilon_i \tag{12}$$

where l_{ij} is definition of constants called *factorloadings* and ϵ_i is added independent error. The constraints for factors are $E(f) = 0$ and $Cov(f) = I$ to preserve factors to be uncorrelated. The result is that factors f_j are independently stretched, rotated and translated to obtain input sample x . Dimension reduction is obtained using factor loadings that are gained from dependencies in observed variables. For more information, Factor Analysis in details is described in [14][8].

3.4 Diffusion Maps

Diffusion Maps is non-linear method which originates from dynamical systems. It enables to discover the underlying manifold from sampled data with some probability distribution. In [6] it was presented as a method based on defining Markov random walk on the graph of the data, where eigenfunctions of Markov matrices can be used to construct diffusion maps representing hidden complex geometric structure. More details about Markov chain and random walk on a graph is explained in [25] and in [2]. Basis of diffusion maps is that the Euclidean distances in diffusion space are approximately equal to diffusion distances in the original space [23]. In process of dimension reduction, first matrix of weights W is constructed by equation using Gaussian kernel function between each data point x_i and x_j :

$$W(i, j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right), \tag{13}$$

where σ is variance of the Gaussian distribution. Then the Markov matrix M is obtained by normalization of the matrix W :

$$M(i, j) = \frac{W(i, j)}{\sum_k W(i, k)}. \tag{14}$$

This defines transition probability from one data point to another data point in a dynamical process. Using spectral decomposition of the matrix M the eigenvalues and eigen-vectors are computed. The sequence of positive eigen-values is sorted in a descending order:

$$1 = \lambda_0 > \lambda_1 \geq \lambda_2 \dots \tag{15}$$

Largest eigen-value λ_0 is trivial because of full connectivity of graph. Therefore, this eigen-value λ_0 and its eigen-vector ψ_0 are excluded from dimensionality reduction. Dimensionality reduction is processed to obtain new features vector y_i in lower dimension d by selecting only d largest eigen-values λ_1 to λ_d and corresponding eigen-vectors ψ_{i1} to ψ_{id} :

$$y_i = \begin{bmatrix} \lambda_1 \psi_{i1} \\ \lambda_2 \psi_{i2} \\ \lambda_3 \psi_{i3} \\ \vdots \\ \vdots \\ \lambda_d \psi_{id} \end{bmatrix} \tag{16}$$

4 Proposed Method

Proposed method of graph visualization in \mathbb{R}^2 space is based on dimension reduction of the original graph. To given graph G with vertices $V = \{1, \dots, n\}$ the distance matrix $T \in \mathbb{R}^{D \times D}$ is computed, where $D = |V|$. Distance between two vertices is presented as the shortest path on a graph between them. Finding shortest paths for all pairs in a undirected unweighted graph is a complex problem. It could be done by algorithmic improvement of breath-first search strategy proposed by Chan [4] or by ADP algorithm proposed in [27]. Obtained distance matrix T is then high-dimensional representation of a graph. Then the dimension reduction method on matrix T is performed and a new matrix is obtained $P \in \mathbb{R}^{D \times d}$, where a new dimension d is in this case equals to 2 to be able draw graph in a plane. Vertices $1, \dots, n$ are now drawn as a points $p_1, \dots, p_n \in \mathbb{R}^2$. Because of connectivity preservation, each edge of a graph is drawn as a straight line between two points representing two connected vertices. Qualitative measurement of a new graph representation is acquired by using the modified version of the original energy function proposed by Kamada and Kawai desribed in section 2.1. For comparable measurement between different graphs with $\frac{n^2-n}{2}$ considered springs between n vertices, equation is then modified as follows:

$$E^{norm} = \frac{1}{n^2 - n} \sum_{i=1}^{n-1} \sum_{j=i+1}^n k_{ij} (\|p_i - p_j\| - l_{ij})^2, \quad (17)$$

5 Experiments and Results

Experiments were performed on a graph datasets obtained from the 10th DIMACS Implementation Challenge website [22]. Particularly graphs *jazz*, *celegans* and *email* were taken from Alex Arenas Website [3]. Graphs *adjnoun*, *lesmis*, *dolphins*, *karate*, *football* and *polbooks* were taken from website of Network Data [24]. Graph *graph₃₀* is manually generated small graph with 30 vertices. Properties of graphs are showed in table 1.

For each dataset, graph drawings using dimension reduction were acquired by all presented methods. In figures 1 and 2 graph drawings are showed for graphs *graph₃₀* and *lesmis* to compare results of different dimension reduction methods.

The qualitative measures were gained for each result of the method using Kamada and Kawai energy function from equation 17 proposed in section 4. Lower energy for given graph drawing is assumed as an indicator of a better representation of the graph in the plane.

For *graph₃₀* and *karate* graphs the best results were achieved by classical method - Kamada and Kawai. For a bigger graphs Kamada and Kawai method has worse results than the dimension reduction methods. Minimal energy for *email*, *celegans*, *jazz*, *football*, *adjnoun*, *lesmis*, *dolphins* graphs were obtained by PCA dimension reduction method. For *polbooks* graph the best representation was done by SNE method.

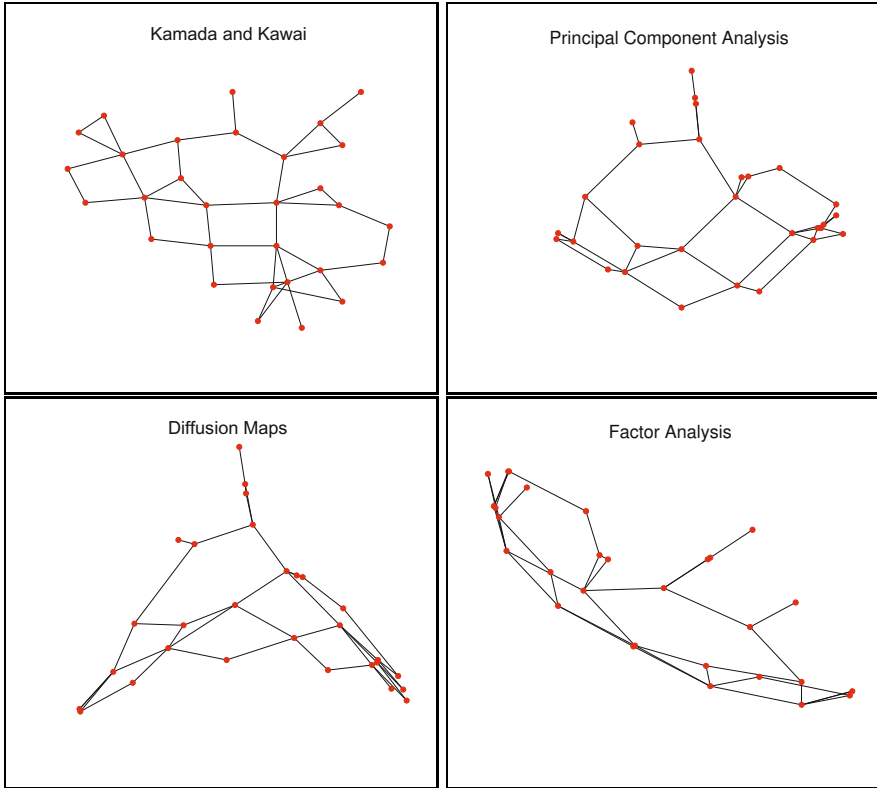


Fig. 1. Example of *graf*₃₀ graph drawing for selected methods

Table 1. Table of graphs properties

<i>graphs</i>	<i>vertices</i>	<i>edges</i>
email	1133	5451
celegans	453	2025
jazz	198	2742
football	115	613
adjnoun	112	425
polbooks	105	441
lesmis	77	254
dolphins	62	159
karate	34	78
graph ₃₀	30	43

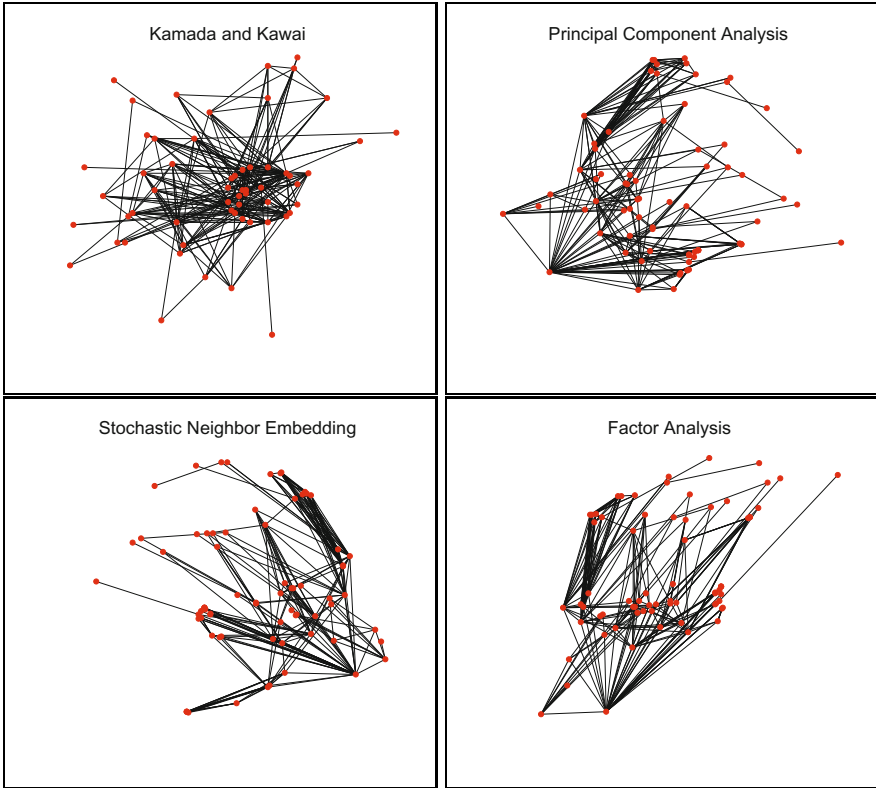


Fig. 2. Example of *lesmis* graph drawing for selected methods

Table 2. Table of Energies for each method and graph computed by modified Kamada and Kawai energy function. (Each value in the table is multiplied by 10^{-3} .)

	SNE	DiffMaps	FA	PCA	Kamada & Kawai
email	5.22	7.73	2.61	1.66	9.48
celegans	5.77	8.86	4.27	2.64	20.46
jazz	5.18	11.21	3.82	2.97	26.79
football	5.06	5.74	4.65	4.57	16.41
adjnoun	5.62	14.07	7.32	4.91	9.68
polbooks	1.91	3.14	3.33	2.51	4.16
lesmis	9.98	16.64	13.01	9.61	11.04
dolphins	2.43	2.70	1.97	1.93	1.97
karate	5.73	4.53	5.22	5.96	2.16
graf ₃₀	21.34	1.17	1.98	1.13	0.38

6 Conclusion

In this work a new graph drawing approach based on dimension reduction of graph representation was presented. All selected methods have been applied to different graphs and example of graph drawings were illustrated in figures. For comparative quality measurement of graph drawings among methods the modified version of Kamada and Kawai energy function was newly introduced. All results were caught in a table and then discussed. This paper presents that a new approach based on dimension reduction methods may be used to obtain good graph drawings. For larger graphs with number of vertices $|V|$ more than about 100 a new approach outperformed standard *Kamada and Kawai* method with best results for *PCA* dimension reduction method.

Acknowledgement. This article has been elaborated in the framework of the IT4Innovations Centre of Excellence project, reg. no. CZ.1.05/1.1.00/02.0070 funded by Structural Funds of the European Union and state budget of the Czech Republic. The work is partially supported by Grant of SGS No. SP2014/110, VŠB - Technical University of Ostrava, Czech Republic. This work was also supported by the Bio-Inspired Methods: research, development and knowledge transfer project, reg. no. CZ.1.07/2.3.00/20.0073 funded by Operational Programme Education for Competitiveness, co-financed by ESF and state budget of the Czech Republic.

References

1. Abdi, H., Williams, L.J.: Principal component analysis. *Wiley Interdisciplinary Reviews: Computational Statistics* 2(4), 433–459 (2010)
2. Aldous, D., Fill, J.A.: Reversible markov chains and random walks on graphs (2002), unfinished monograph, recompiled 2014
3. Arenas, A.: Alex arenas website (2009), <http://deim.urv.cat/~aarenas/data/welcome>
4. Chan, T.M.: All-pairs shortest paths for unweighted undirected graphs in $o(mn)$ time. In: *Proc. ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 514–523 (2006)
5. Cohen, R.F., Battista, G.D., Tamassia, R., Tollis, I.G.: A framework for dynamic graph drawing. *Congressus Numerantium* 42, 149–160 (1992)
6. Coifman, R.R., Lafon, S.: Diffusion maps. *Applied and Computational Harmonic Analysis* 21(1), 5–30 (2006); Special Issue: Diffusion Maps and Wavelets
7. Cunningham, P.: Dimension reduction. In: *Machine Learning Techniques for Multimedia*, pp. 91–112. Springer (2008)
8. Cureton, E., D’Agostino, R.: *Factor Analysis: An Applied Approach*. Taylor & Francis (2013)
9. Díaz, J., Petit, J., Serna, M.: A survey of graph layout problems. *ACM Comput. Surv.* 34(3), 313–356 (2002)
10. Diestel, R.: *Graph Theory*. Electronic library of mathematics. Springer (2006)
11. Forrester, D., Kobourov, S.G., Navabi, A., Wampler, K., Yee, G.V.: Graphael: A system for generalized force-directed layouts. In: Pach, J. (ed.) *GD 2004. LNCS*, vol. 3383, pp. 454–464. Springer, Heidelberg (2005)

12. Fruchterman, T.M.J., Reingold, E.M.: Graph drawing by force-directed placement. *Softw. Pract. Exper.* 21(11), 1129–1164 (1991)
13. Gross, J.L., Yellen, J.: *Graph Theory and Its Applications, Second Edition (Discrete Mathematics and Its Applications)*. Chapman & Hall/CRC (2005)
14. Harman, H.: *Modern Factor Analysis*. University of Chicago Press (1976)
15. Herman, I., Society, I.C., Melançon, G., Marshall, M.S.: Graph visualization and navigation in information visualization: a survey. *IEEE Transactions on Visualization and Computer Graphics* 6, 24–43 (2000)
16. Hinton, G.E., Roweis, S.T.: Stochastic neighbor embedding. In: *NIPS*, pp. 833–840 (2002)
17. Kamada, T., Kawai, S.: An algorithm for drawing general undirected graphs. *Information Processing Letters* 31(1), 7–15 (1989)
18. Kaufmann, M., Wagner, D. (eds.): *Drawing Graphs*. LNCS, vol. 2025. Springer, Heidelberg (2001)
19. Kobourov, S.G.: Spring embedders and force directed graph drawing algorithms. *CoRR abs/1201.3011* (2012)
20. Koren, Y.: Drawing graphs by eigenvectors: theory and practice. *Computers and Mathematics with Applications* 49(11-12), 1867–1888 (2005)
21. Kullback, S., Leibler, R.A.: On information and sufficiency. *The Annals of Mathematical Statistics* 22(1), 79–86 (1951)
22. Meyerhenke, H.: 10th dimacs implementation challenge - graph partitioning and graph clustering (2012)
23. Nadler, B., Lafon, S., Coifman, R.R., Kevrekidis, I.G.: Diffusion maps, spectral clustering and reaction coordinates of dynamical systems. *Applied and Computational Harmonic Analysis* 21(1), 113–127 (2006); Special Issue: Diffusion Maps and Wavelets
24. Newmann, M.: Network data (2013), <http://www-personal.umich.edu/~mejn/netdata/>
25. Norris, J.: *Markov Chains*, No. č 2008. Statistical and Probabilistic Mathematics. Cambridge University Press, Cambridge Series in (1998)
26. Plastria, F., De Bruyne, S., Carrizosa, E.: Dimensionality reduction for classification. In: Tang, C., Ling, C.X., Zhou, X., Cercone, N.J., Li, X. (eds.) *ADMA 2008*. LNCS (LNAI), vol. 5139, pp. 411–418. Springer, Heidelberg (2008)
27. Seidel, R.: On the all-pairs-shortest-path problem in unweighted undirected graphs. *J. Comput. Syst. Sci.* 51(3), 400–403 (1995)
28. Sorzano, C.O.S., Vargas, J., Pascual-Montano, A.D.: A survey of dimensionality reduction techniques. *CoRR abs/1403.2877* (2014)
29. Zaoralek, L., Peterek, T., Dohnálek, P., Gajdos, P.: Comparison of feature reduction methods in the task of arrhythmia classification. In: *IBICA*, pp. 375–382 (2014)

Design of a Single Stage Thermoelectric Power Generator Module with Specific Application on the Automotive Industry

Yidnekachew Messele, Eyerusalem Yilma, and Rahma Nasser

School of Mechanical and Industrial Engineering in Addis Ababa Institute of Technology, Addis Ababa University, Addis Ababa, Ethiopia

Abstract. In a world where fossil fuels dominate as energy sources, the need for an environmentally friendly and economically and commercially viable renewable energy source is dire. As a result, different green technologies have been developed for the generation of energy. One such promising technology is the generation of electrical power from waste heat. As waste heat recovering techniques, thermoelectric generator (TEG) technologies utilization in automotive industry is attempted from many aspects. This paper focuses on analytical investigations on the dynamics of thermoelectric generators at low temperature (waste heat). This research proposes to design and simulate a single stage TEG module to contribute to the further development of TEGs as reasonable energy sources for the consumer market. The design process involved calculating and optimizing the energy balance across the heat absorber, minimizing heat losses, analyzing heat transfer through the thermoelectric elements, and analyzing the electrical power system.

Keywords: Single stage, Thermoelectric- generator, Waste heat recovery.

1 Introduction

In recent years, an increasing concern of environmental issues of emissions, in particular global warming and the limitations of energy resources has resulted in extensive research into novel green technologies, the major one being the generation of electrical power using cleaner technologies. Thermoelectric power generators have emerged as a promising alternative green technology due to their distinct advantages. Thermoelectric power generation offer a potential application in the direct conversion of waste heat energy into electrical power where it is unnecessary to consider the cost of the thermal energy input. The application of this alternative green technology in converting waste-heat energy directly into electrical power can also improve the overall efficiencies of energy conversion systems [1].

Automotive industry is one of the main application fields of TE technologies. One of the main reasons is that a large portion of all generated energy from combustion engine is emitted as waste heat. For a typical gasoline fueled internal combustion engine vehicle, only about 25% of the fuel energy is utilized for vehicle mobility and

accessories the remainder is lost in the form of waste heat and coolant, as well as friction and parasitic losses [2]. The high quality of the waste heat is demonstrated by the high temperature characteristic of vehicle internal environment, which provides possibility of desired large temperature gradients for TEGs. Exhaust gas system is an example of waste heat harvesting location.

A thermoelectric power generator is a solid state device that provides direct energy conversion from thermal energy (heat) due to a temperature gradient into electrical energy based on “Seebeck effect”. The thermoelectric power cycle, with charge carriers (electrons) serving as the working fluid, follows the fundamental laws of thermodynamics and intimately resembles the power cycle of a conventional heat engine. Thermoelectric power generators offer several distinct advantages over other technologies [3].

- They are extremely reliable (typically exceed 100,000 hours of steady-state operation) and silent in operation since they have no mechanical moving parts and require considerably less maintenance;
- They are simple, compact and safe;
- They have very small size and virtually weightless;
- They are capable of operating at elevated temperatures;
- They are suited for small-scale and remote applications typical of rural power supply, where there is limited or no electricity;
- They are environmentally friendly;

2 Background

2.1 Thermoelectricity and Thermoelectric Modules

A heat engine is any device that operates continuously or cyclically and that converts heat to work. Examples include internal combustion engines, power plants, and thermoelectric devices.

We propose to consider the analogies between a classical steam engine and a thermoelectric material. The analogy is that in both systems, the entropy is transported by a fluid, which is here a gas of electrons, also called the “Fermi gas”. At first this Fermi gas can be considered to be a perfect gas, with no interactions between particles. Then the equivalent “partial pressure” p of the fluid in the system is the electrochemical potential μ_e .

$$\mu_e = \mu_c + eV \tag{1}$$

Where; μ_c = the chemical potential,

e = the particle’s charge

v = the electrical potential.

Then the “gas” equivalences for the steam and thermoelectric engines are:

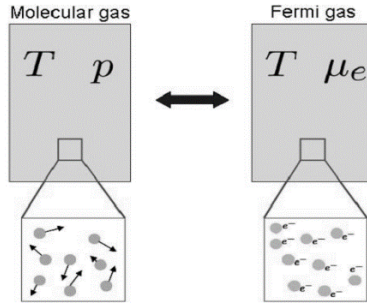


Fig. 1. Comparison between traditional and thermoelectric heat engines

Consider the generation of electrical power by the application of the Seebeck effect that states when a temperature gradient is established within a material, a corresponding voltage gradient is induced. The Seebeck coefficient is a material property representing the proportionality between voltage and temperature gradients and, accordingly, has units of volts/K. For a constant property material experiencing one-dimensional conduction, as illustrated in Figure 2,

$$E_1 - E_2 = S(T_1 - T_2) \tag{2}$$

Where; E_1 = Voltage at heat source
 E_2 = Voltage at heat sink
 S = Seebeck coefficient

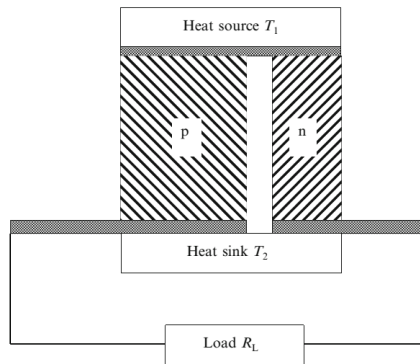


Fig. 2. Simple thermocouple used as a generator [4]

The thermal EMF is equal to $(\alpha_p - \alpha_n)(T_1 - T_2)$ where α_p and α_n are the Seebeck coefficients of the semiconducting pellets, and this gives rise to a current I that may be expressed as;

$$I = \frac{(\alpha_p - \alpha_n)(T_1 - T_2)}{R_p + R_n + R_L} \tag{3}$$

Where; R_p = resistance of p-type thermoelement
 R_n = resistance of n-type thermoelement
 R_L = load resistance

Thence, the power delivered W to the load is;

$$W = I^2 R_L = \left[\frac{(\alpha_p - \alpha_n)(T_1 - T_2)}{R_p + R_n + R_L} \right]^2 R_L \tag{4}$$

The useful power reaches its maximum value when the load resistance is equal to the generator resistance. However, even if there were no loss of heat through thermal conduction, the efficiency could then never exceed 50%. An increase in the load resistance reduces the power output but increases the efficiency. It may be shown that the efficiency becomes a maximum when the ratio, M , of the resistance of the load to that of the generator is given by;

$$M = \frac{R_L}{R_p + R_n} = (1 + ZT_m)^{1/2} \tag{5}$$

Where; ZT_m = the dimensionless figure of merit
 The efficiency of the module can be given as;

$$\eta = \frac{(T_1 - T_2)(M - 1)}{T_1(M + T_2 / T_1)} \tag{6}$$

If ZT_m were much greater than unity, M would also be very large and the efficiency would approach $(T_1 - T_2)/T_1$, which is the value for the Carnot cycle. In Figure 3, the variation of the efficiency with the dimensionless figure of merit for a thermoelectric generator in which the source and sink are at 400 and 360 K, respectively is shown.

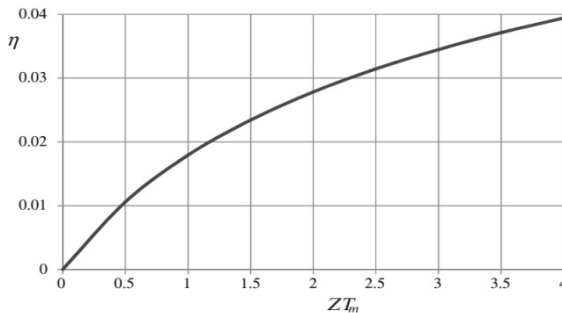


Fig. 3. Plot of efficiency against dimensionless figure of merit for the heat source at 400 K and the heat sink at 360 K [4]

2.2 Thermoelectric Materials and the Figure of Merit

The efficiency of a thermoelectric generator is governed by the thermoelectric properties of the generator materials and the temperature drop across the generator. The temperature difference, ΔT between the hot side (T_h) and the cold side (T_c) sets the upper limit of efficiency through the Carnot efficiency $\eta = \Delta T / T_h$. The thermoelectric material governs how close the efficiency can be to Carnot primarily through the thermoelectric figure of merit, z , defined by [7];

$$z = \frac{\alpha^2}{k\rho} \tag{7}$$

The relevant materials properties are the Seebeck coefficient α , the thermal conductivity k , and electrical resistivity ρ , which all vary with temperature [7].

Established thermoelectric materials conveniently fall into three categories depending upon their temperature range of operation. Bismuth telluride and its alloys have the highest figures-of-merit, are extensively employed in refrigeration, and have a maximum operating temperature of around 450 K. Alloys based on lead telluride have the next highest figures-of-merit with silicon-germanium alloys having the lowest. Lead telluride and silicon germanium are used in generator applications with upper operating temperatures of around 1000 and 1300 K, respectively [5].

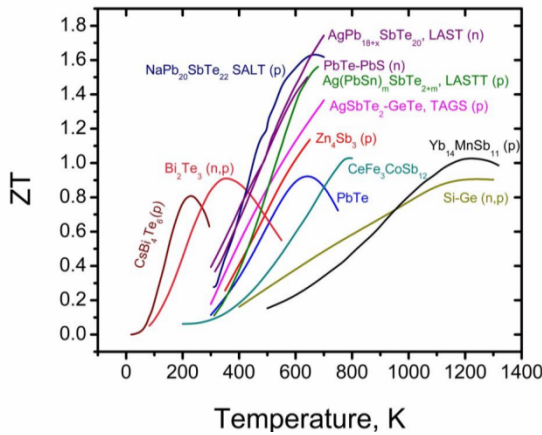


Fig. 4. Figure-of-merit values for different thermoelectric materials [12]

2.3 Composition and Specification of a Thermoelectric Power Generator Module

As shown in Fig. 5, it is composed of two ceramic plates (substrates) that serve as a foundation, providing mechanical integrity, and electrical insulation for n-type (heavily doped to create excess electrons) and p-type (heavily doped to create excess holes) semiconductor thermo-elements [6].

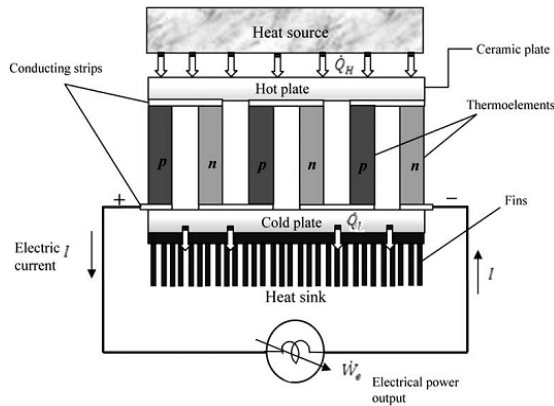


Fig. 5. Schematic diagram of a typical Single-stage thermoelectric power generator [9]

In thermoelectric materials, electrons and holes operate as both charge carriers and energy carriers. There are very few modules without ceramic plates, which could eliminate the thermal resistance associated with the ceramic plates, but might lead to mechanical fragility of the module. The ceramic plates are commonly made from alumina (Al_2O_3), but when large lateral heat transfer is required, materials with higher thermal conductivity (e.g. beryllia and aluminum nitride) are desired. The semiconductor thermo-elements (e.g. silicon-germanium SiGe , lead-telluride PbTe based alloys) that are sandwiched between the ceramic plates are connected thermally in parallel and electrically in series to form a thermoelectric device (module). More than one pair of semiconductors are normally assembled together to form a thermoelectric module and within the module a pair of thermo-elements is called a thermocouple. The junctions connecting the thermo-elements between the hot and cold plates are interconnected using highly conducting metal (e.g. copper) strips as shown in Fig 5 [9].

3 Design of the Proposed TEG Module

3.1 Design Concepts

Heat Source and TEG Module Geometry

Figure 6 shows the variation in temperature across different areas of an automotive waste heat exhaust system.

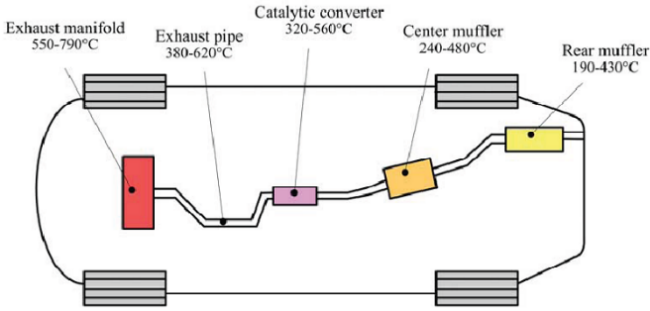


Fig. 6. Typical temperature distribution of an automobile exhaust gas system

3.2 Thermoelectric Materials Design and Selection

Thermoelements Selection

It has been found convenient to introduce a quantity known as the power factor that contains both the Seebeck coefficient and the electrical conductivity. The power factor is defined as $(\alpha^2 \times \sigma)$ and is useful because alpha and sigma are the parameters that are most strongly dependent on the carrier concentration. The other quantity that is involved in the definition of the figure of merit is the thermal conductivity; lambda λ . Lambda is less dependent on the concentration of the charge carriers since it is often dominated by the lattice contribution. Thus, the carrier concentration that yields the maximum power factor for a given material is usually close to that which gives the highest figure of merit.

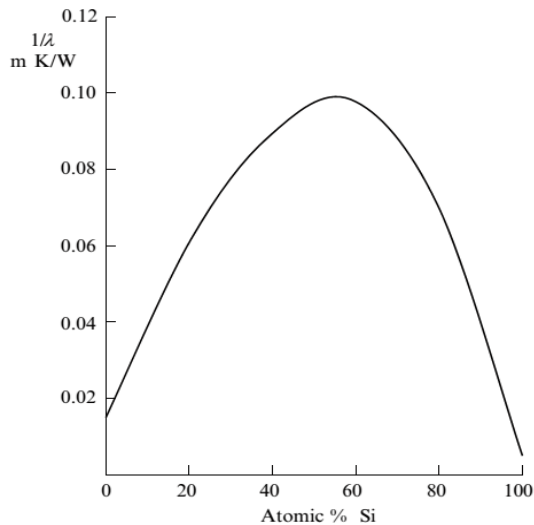


Fig. 7. Thermal resistivity of silicon-germanium alloys at 300 K. Schematic plot based on the data of Steele and Rosi

Thus due to cost and manufacturing considerations the Si-Ge alloy semiconductors are chosen for the n and p type semiconductor thermo-elements, without a considerable loss in efficiency and power factor.

The lattice conductivities of silicon and germanium at 300 K are 145 and 64Wm-1K-1, respectively. The value of λL falls rapidly on adding germanium to silicon, as is apparent from Fig.7 in which the thermal resistivity is plotted against the concentration of silicon in germanium for the whole range of Si-Ge alloys.

Substrate Material

The substrate material must be a good electrical insulator and at the same time be a good thermal conductor. It also serves as a foundation for the TEG module, providing mechanical integrity for the p and n type semiconductor thermo-elements. There are different materials that can be utilized as a substrate material for a thermoelectric module, but ceramics are the most widely used.

Table 1. Typical electrical conductivity values for different ceramic materials [19]

<i>Material</i>	<i>Electrical Conductivity</i> [[Ω-m) ⁻¹]
Graphite	$3 \times 10^4 - 2 \times 10^5$
Ceramics	
Concrete (dry)	10^{-9}
Soda-lime glass	$10^{-10} - 10^{-11}$
Porcelain	$10^{-10} - 10^{-12}$
Borosilicate glass	$\sim 10^{-13}$
Aluminum oxide	$< 10^{-13}$
Fused silica	$< 10^{-18}$

As can be seen from the table above the ceramics are very good electrical insulators with very low electrical conductivity values. Comparing electrical conductivity values and considering easy availability and cost the top and bottom substrate materials are selected to be made from aluminum oxide.

3.3 Heat Transfer Analysis

If the material of Figure 8 is installed in an electric circuit, the voltage difference induced by the Seebeck effect can drive an electric current I, and electric power can be generated from waste heat that induces a temperature difference across the material. A simplified thermoelectric circuit, consisting of two pellets of semiconducting material, is shown in Figure 8.

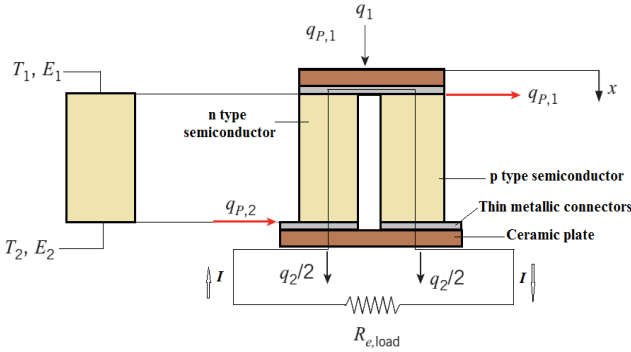


Fig. 8. A simplified thermoelectric circuit consisting of one pair (N=1) of semiconducting pellets

In addition to inducing an electric current I , thermoelectric effects also induce the generation or absorption of heat at the interface between two dissimilar In addition to inducing an electric current I , thermoelectric effects also induce the generation or absorption of heat at the interface between two dissimilar

$$q_p = I(S_p - S_n)T = I(S_p - n)T \tag{8}$$

Points of Considerations for the TEG Module Heat Transfer Analysis

- One dimensional, steady-state conduction.
- Internal heat generation in the semiconductor plates.
- Constant thermo-physical properties.
- Constant temperature of the hot region.
- Thermal resistances of the thin metallic conductors (coppers) are assumed to be negligible due to their high thermal conductivity values and very small thickness.
- Contact resistance at the junctions of the different layers of materials is assumed to be negligible if sufficient pressure is applied between the top and bottom plates of the TEG module.

The expression for the power developed inside the module considering only a single pair of thermoelements is given by;

$$P = 2K_m A_m \left[\left(\frac{\dot{q}}{K_m} ((L_c + L_{cu}) - C_1) \right) - \left(\frac{\dot{q}}{K_m} ((L_c + L_{cu} + L) - C_1) \right) \right] + I(S_p - n)(T_1 + T_2) - 2I^2 R_e \tag{9}$$

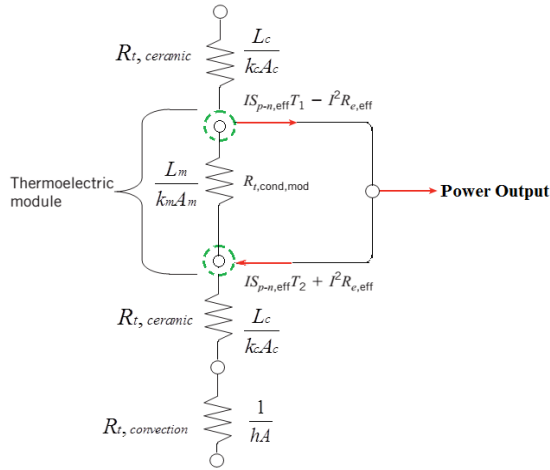


Fig. 9. Equivalent Thermal Circuit

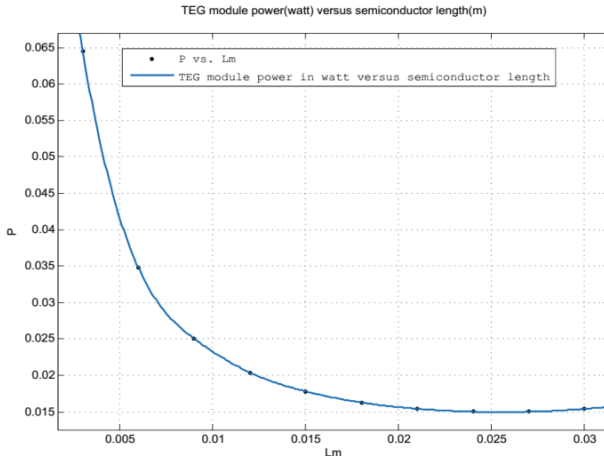


Fig. 10. TEG module power versus semiconductor length

An iteration is done by holding constant the hot and cold side temperatures at their respected values and varying the values of L_c , (thickness of ceramic plates), (0-9mm), L_{cu} , (thickness of the copper plates) and L , (the length of the thermoelements). It is reasonably assumed that L_{cu} , is half of L_c . The corresponding length of the thermoelements (L) is determined based on the imposed condition that the total height of the module is 30mm. The other parameters are obtained by incorporating in the iteration scheme the corresponding relationships and equations. Thus, for the selected row of geometric and thermoelectric parameters of the module and at the maximum temperature gradient specified (4000c), the module produces in total 0.017801watts of power or 58Volts and its efficiency is 0.33%. Note that the total number of thermoelements is given in section 3.4.

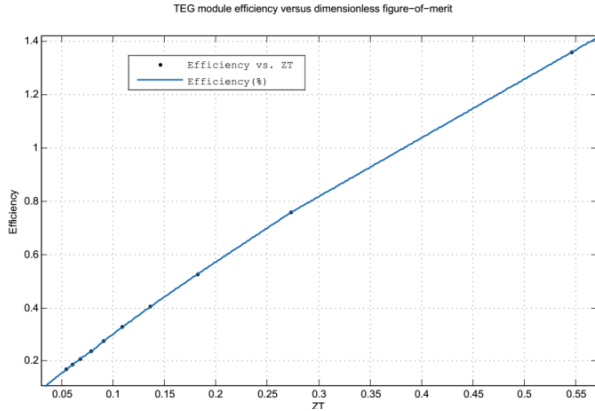


Fig. 11. TEG module efficiency versus dimensionless figure of merit

As can be seen from the efficiency versus dimensionless figure-of-merit plot, it closely resembles the behavior of the curve obtained experimentally for different thermoelectric materials. (Refer figure 3).

3.4 Configuration of Thermoelements

The number of pairs of thermo-elements depends on the geometric configuration of the pairs.

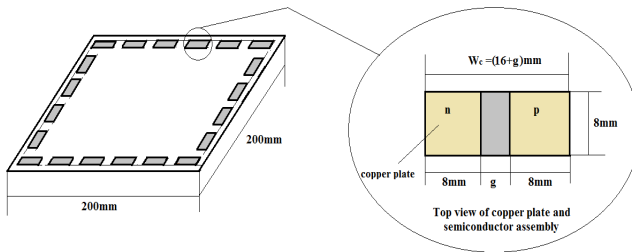


Fig. 12. 3-dimensional configuration of the copper plate assembly

The total number of copper plates that can be fitted in the square ceramic plates can be approximated by the formula;

$$N = \left[\frac{n(16 + g) + (n - 1)g}{8 + g} \right] * (n - 2) + (4n) \tag{10}$$

Where n denotes the number of copper plates that can be fitted onto the module without exceeding the width. (i.e. 200mm) and g is the gap between the n and p type semiconductor pellets. Thus, the total number of copper plates is selected to be 116.

3.5 Solidworks Heat Transfer Simulation

The 3D Solidworks model of the module with the top ceramic removed is given below.

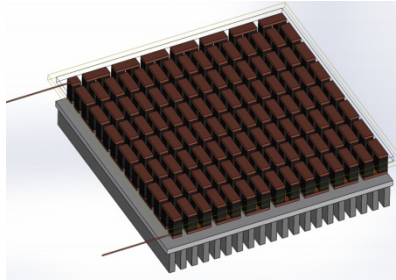


Fig. 13. 3D model of the assembled TEG module

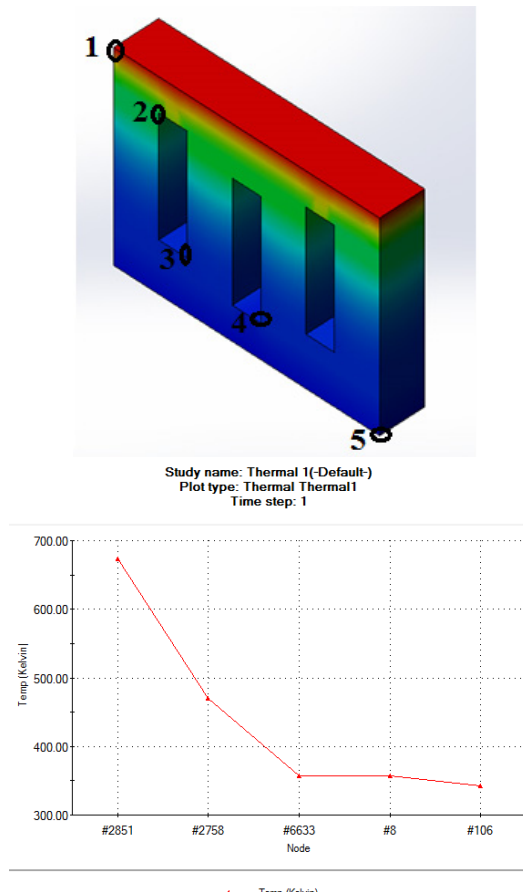


Fig. 14. Temperature values at 5 different node points on the module

The heat transfer computer simulation attached herewith describes the performance of the TEG module under normal operating conditions. A simple model of the TEG module comprising of two pairs of semiconducting pellets (thermocouples) is used to simulate the steady state heat transfer phenomena inside the module. The model is loaded with the maximum temperature that can be provided by the heat source (automotive muffler).

From the steady state heat transfer simulation it can be seen that the maximum effective temperature gradient between the top and bottom surfaces of the semiconductor plates is 112°C . Note that the maximum temperature gradient between the hot and cold sides of the module in the theoretical analysis was considered to be 80°C . This temperature difference subsists due to assumptions that has not been considered such as; contact resistance between the different plates, different environmental effects and the resistance of the copper plate's etc. The semiconductor plates are generating heat internally at a temperature value of around 430K. Thus, the values theoretically predicted in the simulation for the temperature gradient may be reduced to the anticipated gradient value (80°C) in the specification. In the future, improvement on power generation can be achieved by further enhancing the thermoelectric properties of the semiconducting pellets.

4 TEG Module Electric Circuit Performance Analysis

Theoretical electrical circuit performance analysis for TEG modules are very helpful in demonstrating the performance of the module in regards to power generation. Different analysis software's including Simulink can be used to study the circuit performance of the TEG module. In this chapter such a software model is used to estimate the various performance parameters of the module. The steps utilized for the study include;

- The representation of the module as an electrical circuit and the derivation of the mathematical equations representing the system using Kirchhoff's' voltage law.
- The representation of the already formulated mathematical equations as a block diagram in the Matlab, Simulink environment.
- The evaluation and comparison of the simulation results with the theoretical predictions

4.1 Mathematical Modelling

The generator can be modeled using the electrical schematic shown in Figure 15, where the open circuit voltage, the source resistance and the capacitance represent the

single stage TEG module and the load is a resistor with the output as the voltage drop across this resistor.

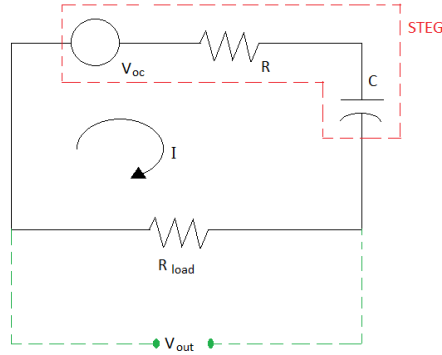


Fig. 15. Electrical circuit representing the single stage TEG module

4.2 Simulink Modelling

In general, the mathematical equations representing a given system that serves as the basis for a Simulink model can be derived from physical laws. The TEG module starts to produce power, theoretically, as soon as a temperature gradient starts to develop between the two extremes of the thermoelement faces. In the Simulink model such a very small temperature gradient is modeled with a step input from the commonly used blocks, to outset current and thus power generation inside the module. This step input is shown to be connected with a unity gain to verify that the input is not amplified (this can be eliminated) and then to an adder where the system parameters converge. The differential equation is represented by the series of blocks to the right of the adder and the input to the system is represented by the series of blocks to the left of the adder as described above.

As can be seen from the simulation result above, the module does not produce voltage for one second. When the temperature gradient develops with time, the module will eventually start to produce voltage. Furthermore, the TEG module produces a voltage of magnitude around 4.25V for the first 10seconds of operation. Given enough time until the system reaches the expected temperature gradient between the hot and cold ends of the module, the voltage is expected to reach its peak value of 58V. (Note that this value of voltage can be computed from values of power and current of the module obtained in the previous thermo-electric analysis.)

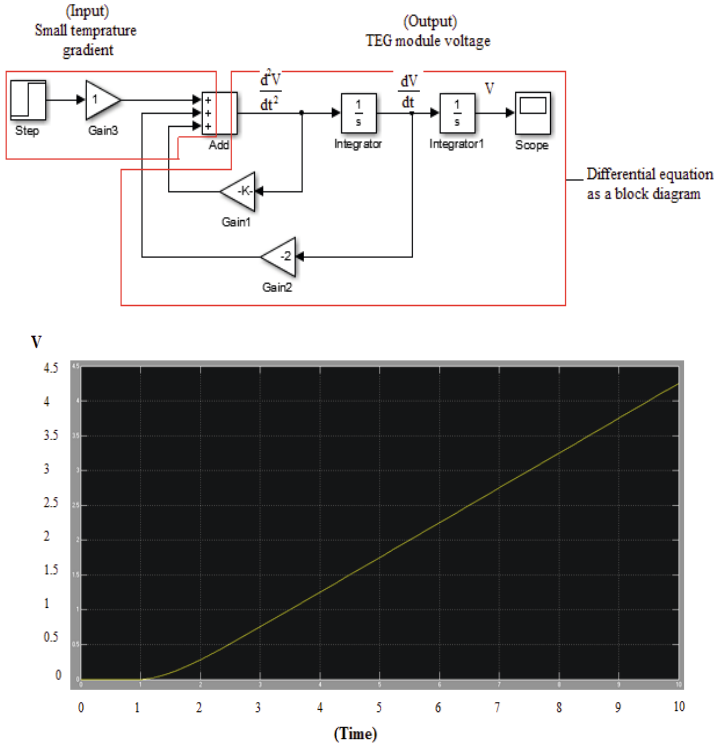


Fig. 16. Simulation result showing the variation of the TEG module’s voltage with time

5 Conclusion

The design of a single stage TEG module is presented as a solution, owing to the current situation of the world with regard to uncovering green technologies that are cost effective and simple. Also the different ways of integrating these devices in to an automotive for alternative power generation and their feasibility for future applications has been established.

Different thermoelectric materials has been studied to be employed on the TEG module. Si-Ge alloys, due to their satisfactory thermoelectric properties at the temperature range considered (i.e. 593-673K) and their easy availability in the market has been chosen as the p and n type thermoelements.

Under the assumption of a non-fluctuating temperature gradient, a rigorous theoretical analysis showed that the TEG module is able to produce a voltage of about 58V. With an increase in temperature gradient from zero to its maximum value, a Simulink study shows that the TEG module has a voltage output of about 4.5V for the first 10seconds of operation. Given sufficient time the module can reach its peak value of 58V in power generation.

Under normal operating conditions, the efficiency of the TEG module is 0.33%. The total power output of the module is 0.017801Watt. A Solidworks heat transfer simulation showed that the effective temperature gradient between the top and bottom faces of the thermoelements to be 112⁰C. However, different thermophysical and environmental conditions that are not considered in the simulation study can reduce the anticipated value to the design value (i.e. 80⁰C).

A systematic design procedure considering both the thermal and electrical phenomenon inside the module is presented which can be used on a similar study in the future.

References

1. Rowe, D.M.: Thermoelectric waste heat recovery as a renewable energy source. *Int. J. Innov. Energy Syst. Power* 1, 13–23 (2006)
2. Li, M.: Thermoelectric Generator Based DC-DC Conversion Network for Automotive Applications. Master of Science Thesis Stockholm, Sweden (2011)
3. Ismail, B.I., Ahmed, W.H.: Thermoelectric Power Generation Using Waste-Heat Energy as an Alternative Green Technology. Department of Mechanical Engineering, Lakehead University, Canada, Component Life Technology, Atomic Energy of Canada Ltd., Canada
4. Julian Goldsmid, H.: Introduction to thermoelectricity. Series in Material Science, vol. 121
5. Tritt, T.M.: Thermoelectric Materials: Principles, Structure, Properties, and Applications. Copyright, Elsevier Science Ltd. (2002)
6. Jeffrey Snyder, G., Caillat, T.: Using the Compatibility Factor to Design High Efficiency Segmented Thermoelectric Generators. Jet Propulsion Laboratory/California Institute of Technology
7. Jeffrey Snyder, G.: Small Thermoelectric Generators
8. Boukai, A.I.: Thermoelectric Properties of Bismuth and Silicon Nanowires. Degree of Doctor of Philosophy Thesis, California Institute of Technology, Pasadena, California (2008)
9. Shawwaf, A.: Optimization of the electric properties of thermoelectric generators. Department of Automatic Control, Lund University (December 2010)
10. Watzman, S.: Design of a Solar Thermoelectric Generator. Undergraduate Honors Thesis, Ohio State University
11. Flores-Livas, J.A.: Thermoelectricity: An Introduction. Laboratoire de Physique de la Matière Condensée et Nanostructures, Université Claude Bernard Lyon 1
12. Kandylas, I.P., Stamatelos, A.M.: Engine exhaust system design based on heat transfer computation. Laboratory of Applied Thermodynamics, Mechanical Engineering Department, Aristotle University of Thessaloniki, Thessaloniki, Greece
13. Bergman, T.L., Lavine, A.S., Incropera, F.P., Dewitt, D.P.: Fundamentals of heat and mass transfer, 7th edn.
14. Wolverine tube heat transfer data book
15. http://www.KELK.co.jp/english/thermo/pdf/popup_01.pdf

16. Goupil, C., Seifert, W., Zbrocki, K., Muller, E., Snyder, G.J.: Thermodynamics of Thermoelectric Phenomena and Applications
17. Karris, S.T.: Introduction to Simulink® with Engineering Applications. Orchard Publications
18. Kreith, F., Boehm, R.F., et al.: Heat and Mass Transfer. Mechanical Engineering Handbook. CRC Press, Boca Raton (1999), Ed. Frank Kreith
19. Callister Jr., W.D., Rethwisch, D.G.: Fundamentals of Materials Science and Engineering, An Integrated Approach. John Wiley & Sons, Inc. (2007)

Ethiopian Livestock Husbandry Cluster Identification Using FUZZY-AHP Approach

Netsanet Jote, Birhanu Beshah, and Daniel Kitaw

Addis Ababa Institute of Technology, School of Mechanical and Industrial Engineering,
Addis Ababa, Ethiopia
netsijote@gmail.com, birhanu.beshah@aait.edu.et,
danielkitaw@yahoo.com

Abstract. The problems of leather sector in Ethiopia starts from animal husbandry stage. This calls for intervention options as early as possible in the supply chain. In this paper livestock husbandry cluster is proposed to mitigate the problems of Ethiopian leather sector at animal husbandry stage. The first and the most important stage of industrial clustering procedure is identifying best area for cluster development. Livestock husbandry cluster identification is a strategic decision with uncertainties. To handle the uncertainties, Fuzzy-AHP based livestock husbandry cluster identification is proposed. Up to now, there is no research conducted on Fuzzy-AHP for livestock husbandry cluster identification. Therefore, the aim of this paper is to identify livestock husbandry cluster in Ethiopia using Fuzzy-AHP. As a result, three alternatives (i.e. West Gojjam, East Gojjam and North Shewa) and six quantitative and qualitative criteria (i.e. geographical proximity, sectorial concentration, market potential, support services, resource potential and potential entrepreneurs) are found. Finally, North Shewa is selected as best area for livestock husbandry clusters. A sensitivity analysis is also performed to justify the results.

Keywords: Fuzzy-AHP, Livestock husbandry, Cluster identification.

1 Introduction

Ethiopia is relatively well endowed in its livestock base in Africa, and stands seventh in the world in cattle, ninth in sheep, and eighth in goats [12]. This amount covers 2.88% global share of livestock population. This enormous population of livestock promises an ample opportunity for the development of the leather industry in the country.

Even if, Ethiopia has a large livestock population, the quality of hide and skin supplied to tanneries is poor quality. This is mainly due to poor animal husbandry practices, poor handling of the hides and skins at the slaughter facilities/abattoirs, poor storage and preservation methods and inadequate enforcement of the existing laws related to hides and skins. According to research done by Leather Industry Development Institute (LIDI) on eight tanneries of Ethiopia, 90% of defects of sheep skin, 85% of defects of goat skin and 59% of defects of cattle hide occur at livestock

husbandry stage [1]. This problem of quality caused by defects extends to each of the subsequent stages of processing of the leather, thus ultimately determining the price paid to the primary producer and of the semi processed or an end product. To solve this problem, livestock husbandry cluster is proposed as a solution. In Ethiopia there are successful clusters which give firms the ability to grow together. Some of these are: Merkato footwear cluster, Guleli Handloom Cluster, Addis Ababa Ready-made Garments Cluster and Mekelle Metal and Wood Enterprises Cluster.

A cluster is a concentration of interconnected, geographically close businesses operating together within the same commercial sector and whose activities rely on certain local specificities such as availability of natural resources, centres for technological development (through universities, research centres, technology parks, or a technology-based industry), and a consolidated productive structure for all tiers of the productive chain of the region [15]. Most researchers in the area believe that a cluster approach is the most feasible approach for developing small enterprises and large industries [19]. Similarly, developing livestock husbandry cluster will generate export standard meat product and quality by products. The first and the most important stage of clustering procedure is identifying best area for cluster development.

Cluster identification process is a complex process that involves both qualitative and quantitative, often conflicting criteria. It is also a strategic decision with uncertainties. In our previous papers [13] and [14], we compare the advantage and the disadvantage of existing cluster identification tools. The pitfall of the existing methods is that, they cannot handle uncertainties. To handle the uncertainties Fuzzy-AHP approach is selected. In this paper we use Fuzzy-AHP approach to identify livestock husbandry cluster for the first time.

The rest of the paper is organized as follows: Section two explores the literature review; Section three presents research methodology. Section four presents results and discussions. Finally, Section five presents the conclusions.

2 Literature Review

Analytic Hierarchy Process (AHP), introduced by Saaty, is a useful and practical tool that provides the ability to incorporate both qualitative and quantitative factors in the decision-making process [5]. Any complex problem can be decomposed into several sub-problems using AHP in terms of hierarchical levels [4]. One of the main advantages of the AHP method is the simple structure and design which represent human mind and nature [21]. But, it is generally criticized by the use of a discrete scale of 1-9 which cannot handle the uncertainty and ambiguity present in deciding the priorities of different attributes [5]. To overcome these problems, several researchers integrate fuzzy theory with AHP to improve the uncertainty. The use of Fuzzy-AHP for multiple criteria decision-making requires scientific approaches for deriving the weights from fuzzy pair-wise comparison matrices [20]. Recently, Fuzzy-AHP has been widely used to solve multi-criteria decision problems, such as: [9], [6], [7], [18]. However, up to now, no research has been conducted on Fuzzy-AHP for identification of livestock husbandry cluster.

Fuzzy theory is composed of three key factors: fuzzy set, membership function, and fuzzy number to change vague data into useful data efficiently [11]. The merit and strength of using fuzzy approach is to express the relative importance of the alternatives and the criteria with fuzzy numbers instead of using simple crisp numbers as most of the decision-making problems in the real world takes place in a situation where the pertinent data and the sequences of possible actions are not precisely known [11]. Triangular and trapezoidal fuzzy numbers are usually used to capture the vagueness of the parameters which are related to select the alternatives [8]. Triangular Fuzzy Numbers (TFN) are expressed with boundaries instead of crisp numbers for reflecting the fuzziness as decision-makers select the alternatives or pair-wise comparisons matrix. This paper used TFN to prioritize livestock husbandry cluster areas with fuzziness.

The surveyed literatures show that the vast majority of the Fuzzy-AHP applications use a simple extent analysis method proposed by Chang (1996) [3]. The extent analysis method is used to consider the extent of an object to be satisfied for the goal, that is, satisfied extent [9]. In the method, the “extent” is quantified by using a fuzzy number. The basics of the Extent Analysis Method on Fuzzy-AHP are introduced in different researches [9], [7]. As illustrated above, extent analysis method on Fuzzy-AHP is used to solve multi-criteria decision problems. The paper used this approach to identify livestock husbandry cluster.

3 Methodology

The methodology is based on the approach described in our previous papers [13] and [14]. The main aim of this paper is to apply Fuzzy-AHP for livestock husbandry cluster identification. The methodology consists of five main steps. In the first step, best origins or alternatives for livestock husbandry cluster are identified. In the second step, cluster identification criteria are selected. In the third step, weights of cluster selection criteria are calculated using the Fuzzy-AHP process. During the fourth step, the alternative ranking results are calculated and the best origin (alternative) for the livestock husbandry cluster is determined. Finally, sensitivity analysis is provided.

4 Result and Discussion

4.1 Best Origins (Alternatives) for Livestock Husbandry

Experts from the Ministry of Agriculture (MoA), the Ministry of Trade (MoT), the Leather Industry Development Institute (LIDI), the Addis Ababa Abattoirs Enterprise, and the Ethiopian Leather Industries Association (ELIA) have been interviewed. Raw hide and skin collectors, big raw hide and skin traders and raw hide and skin selectors in some tanneries have also been interviewed to strengthen the assessment. Information collected through interviews covers about the Ethiopian leather sector, problems and challenges of leather sector at each stage of the supply chain, suitable

solutions, best origin areas for sheep skin and goat skins and cattle hides, characteristics of Ethiopian hide and skin, quality measurement parameters of hide and skin, etc. the interviews were carried out face-to-face and using the telephone. From the interview conducted, 23 raw hide and skin source areas were identified. From these areas, to select best origins/areas for raw hide and skin, questionnaires are distributed to 30 tanneries out of which 21 were filled and returned. The questionnaires were distributed through personal contact, and e-mail. The questionnaire has two parts. In the first part, the experts asked to grade the natural quality of hide and skin in the given areas. In the second part, based on the given grade, they are asked to measure the quality of hides and skins using five parameters (i.e. natural heritage, feeding practice, husbandry practice, free from diseases and absence of branding practice). The questionnaire and the interviews were concerned with best origins/areas for sheep and goat skins and cattle hides production in Ethiopia. But the detailed analysis of the study was restricted to sheep skin. This is because; Ethiopian sheepskin has a reputation for its fibre strength and other qualities attractive to the international market. It has also a higher relative price, and leather manufacturers use it to produce most valuable finished products that are in high demand. It can be used to produce best quality dresses, gloves and shoe uppers that sell in high end retail outlets in Europe, the United States and other developed countries.

Aggregating questionnaire results by doing frequency analysis, three best origins for sheep skin production; namely West Gojjam, East Gojjam and North Shewa, were selected. Farmers found in three areas are used for analysis. These areas are identified as best origins for sheep husbandry in Ethiopia. From these three origins/areas the best area for cluster using Fuzzy-AHP methodology has been identified.

4.2 Identification of Cluster Selection Criteria

The detailed analysis for cluster selection criteria is given in our previous papers [13] and [14]. The same criteria are used to in this study also. These are Geographical proximity (GP), Sectorial Concentration (SC), Market Potential (MP), Support Services (SS), Resource Potential (RP) and Potential Entrepreneurs (PE).

To select best sheep husbandry cluster area, Fuzzy-AHP approach is introduced. The method allows a complex decision to be structured into a hierarchy descending from an overall objective to various criteria, sub-criteria and so on until the lowest level. First, the overall goal of the cluster identification problem has been identified which was “best area for sheep husbandry cluster”. To identify the best cluster, as explained above, six criteria are selected by experienced experts. Finally, the three sheep husbandry areas are laid down at the last level of the hierarchy. Fig 2 shows the hierarchical structure of the objective, criteria and alternatives.

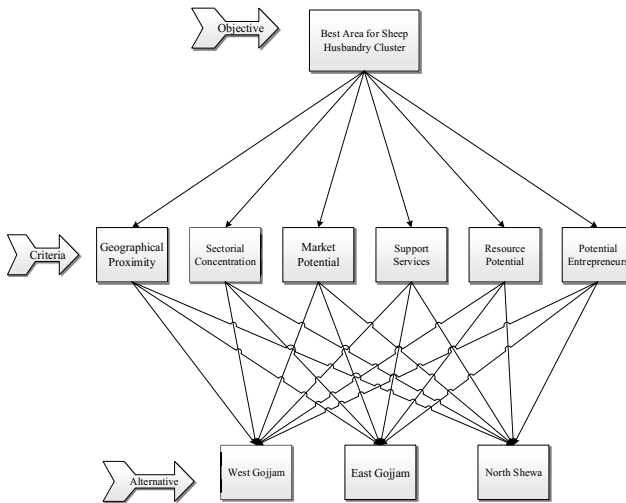


Fig. 1. The hierarchical structure, the alternatives and the criteria

4.3 Pair-wise Comparison

After identifying the criteria, the different priority weights of each criterion were calculated using the Fuzzy-AHP approach. The comparison of the importance of one criterion over another was achieved by the help of the questionnaire. The questionnaire facilitates the answering of pair-wise comparison questions. The preference of one measure over another was decided by the experience of the experts.

Expert used the linguistic variables to compare the criteria with respect to the main goal. Then the linguistic variables were converted to triangular fuzzy numbers. Table 1 shows the linguistic variables and their corresponding triangular fuzzy numbers.

After the pair-wise comparison matrices were formed, the consistency of the pair-wise judgment of each comparison matrix was checked, using the calculation method of consistency index and consistency ratios in crisp AHP.

Each triangular fuzzy number, $M = (l, m, u)$ in the pair-wise comparison matrix was converted to a crisp number using $M\text{-crisp} = (4 * m + l + u) / 6$. After the fuzzy comparison matrices were converted into crisp matrices; the consistency of each matrix was checked by the method in crisp AHP [9].

After calculating the consistency ratios of the entire matrix and making it below 0.1, the next step is to calculate the weight vector for each factor lying at different levels of the hierarchy using Chang's extent analysis approach.

Table 1. Triangular fuzzy Scale[3],[2],[10]

Linguistic Scale	Triangu- lar fuzzy scale	Triangular fuzzy recip- rocal scale	Explanation
Equally important	(1,1,1)	(1,1,1)	Two elements contribute equally
Moderately important	(2/3,1,3/2)	(2/3,1,3/2)	One element is slightly favored over another
Strongly important	(3/2,2,5/2)	(2/5,1/2,2/3)	One element is strongly favored over another
Very strongly important	(5/2,3,7/2)	(2/7,1/3,2/5)	An element is very strongly favored over another
Extremely important	(7/2,4,9/2)	(2/9,1/4,2/7)	One element is the highest favored over another

The fuzzy evaluation matrix with respect to the goal with triangular fuzzy numbers can be seen in Table 2. Before calculating the weights, the consistency of the comparison matrixes is checked. All of the comparison matrixes are consistent. Because of space limitation, the processes of the consistency check are not shown in this research.

In order to find the priority weights of each criterion, Chang’s extent analysis approach is used. First the fuzzy synthetic extent values of the attributes were calculated by using Eq. (1)

The value of fuzzy synthetic extent with respect to the i^{th} object is defined as:

$$S_i = \sum_{j=1}^m M_{gi}^j * [\sum_{i=1}^m \sum_{j=1}^m M_{gi}^j]^{-1} \tag{1}$$

The different values of fuzzy synthetic extent of the six different main criteria were denoted by $S_{GP}, S_{SC}, S_{MP}, S_{SS}, S_{RP}$ and S_{PE} .

Table 2. The fuzzy evaluation matrix with respect to the goal

	GP	SC	MP	SS	RP	PE
GP	(1,1,1)	(2/3,1,3/2)	(2/3,2,5/2)	(2/3,2,5/2)	(3/2,1,3/2)	(5/2,3,7/2)
SC	(2/3,1,3/2)	(1,1,1)	(2/3,1,3/2)	(2/3,1,3/2)	(2/3,1,3/2)	(3/2,2,5/2)
MP	(2/5,1/2,2/3)	(2/3,1,3/2)	(1,1,1)	(1,1,1)	(2/3,1,3/2)	(2/3,1,3/2)
SS	(2/3,1/2,2/3)	(2/3,1,3/2)	(1,1,1)	(1,1,1)	(2/3,1,3/2)	(2/3,1,3/2)
RP	(3/2,1,3/2)	(2/3,1,3/2)	(2/3,1,3/2)	(2/3,1,3/2)	(1,1,1)	(2/3,1,3/2)
PE	(2/7,1/3,2/5)	(2/5,1/2,2/3)	(2/3,1,3/2)	(2/3,1,3/2)	(2/3,1,3/2)	(1,1,1)

$$S_{GP}=(7.84,10,12.5)*(1/51.41,1/38.8,1/29.89)= (0.15,0.26,0.42)$$

$$S_{SC}=(5.2,7,9.5)*(1/51.41,1/38.8,1/29.89)= (0.1,0.18,0.32)$$

$$S_{MP}=(4.41,5.5,7.17)*(1/51.41,1/38.8,1/29.89)= (0.09,0.14,0.24)$$

$$S_{SS}=(4.41,5.5,7.17)*(1/51.41,1/38.8,1/29.89)= (0.09,0.14,0.24)$$

$$S_{RP}=(4.37,6,8.5)*(1/51.41,1/38.8,1/29.89)= (0.09,0.15,0.28)$$

$$S_{PE}=(3.7,4.8,6.6)*(1/51.41,1/38.8,1/29.89)= (0.09,0.12,0.22)$$

The value of S_i has been compared individually and the degree of possibility of $M_2 = (l_2, m_2, u_2) \geq M_1 = (l_1, m_1, u_1)$ are identified using Eq. (2).

$$= \begin{cases} 1, & \text{if } m_2 \geq m_1, \\ 0, & \text{if } l_1 \geq u_2, \\ \frac{l_1 - u_2}{(m_2 - u_2) - (m_1 - l_1)} & \text{otherwise,} \end{cases} \tag{2}$$

Thereafter, the minimum degree of possibility $V(M \geq M_i), i = 1,2,3, \dots, k$) has been determined using Eq. (3).

$$V(M \geq M_1, M_2, \dots, M_k) = V[(M \geq M_1) \text{and } (M \geq M_2) \text{and } \dots \text{and } (M \geq M_k)] = \min V(M \geq M_i), i = 1,2,3, \dots, k \tag{3}$$

$$\begin{aligned} \min V(M_{GP} \geq M_{SC}, M_{MP}, M_{SS}, M_{RP}, M_{PE}) &= \min(1,1,1,1,1) = 1 \\ \min V(M_{SC} \geq M_{GP}, M_{MP}, M_{SS}, M_{RP}, M_{PE}) &= \min(0.68,1,1,1,1) = 0.68 \\ \min V(M_{MP} \geq M_{GP}, M_{SC}, M_{SS}, M_{RP}, M_{PE}) &= \min(0.43,0.78,1,0.94,1) = 0.43 \\ \min V(M_{SS} \geq M_{GP}, M_{SC}, M_{MP}, M_{RP}, M_{PE}) &= \min(0.43,0.78,1,0.94,1) = 0.43 \\ \min V(M_{RP} \geq M_{GP}, M_{SC}, M_{MP}, M_{SS}, M_{PE}) &= \min(0.54,0.86,1,1,1) = 0.54 \\ \min V(M_{PE} \geq M_{GP}, M_{SC}, M_{MP}, M_{SS}, M_{RP}) &= \min(0.33,0.67,0.87,0.87,0.81) = 0.45 \end{aligned}$$

Therefore, the weight vector shown below was found using Eq. (4):

$$\begin{aligned} W' &= ((d'(A_1), d'(A_2), \dots, d'(A_n))^T \\ W' &= (1, 0.68, 0.43, 0.43, 0.54, 0.33)^T \end{aligned} \tag{4}$$

Finally, the weight vectors have been normalized using Eq. (5) and the relative weights of the 6 criteria are obtained.

$$\begin{aligned} W &= (d(A_1), d(A_2), \dots, d(A_n))^T \\ W &= (0.28, 0.2, 0.13, 0.13, 0.16, 0.1)^T \end{aligned} \tag{5}$$

The final weights for Geographical Proximity (GP), Sectorial Concentration (SC), Market Potential (MP), Support Services (SS), Resource Potential (RP) and Potential Entrepreneurs (PE) were found to be 0.28, 0.2, 0.13, 0.13, 0.16 and 0.1, respectively. It has been conclude that the most important criteria for livestock husbandry cluster identification process is geographical proximity criteria as it has the highest priority weight. Sectorial concentration is the next preferred criteria. This result is supported by Porter’s (1998, 1990) cluster definitions [16],[17].

4.4 Prioritize and Rank the Alternatives (Origins)

The same calculations were applied to the other pair-wise comparison matrices and the priority weights of the three alternatives with respect to Geographical Proximity (GP), Sectorial Concentration (SC), Market Potential (MP) , Support Services (SS), Resource Potential (RP) and Potential Entrepreneurs (PE) criteria illustrated in table 3, 4, 5,6,7 and 8 respectively.

Table 3. Fuzzy Comparison matrix of the alternatives with respect to GP criteria

	West Gojjam	East Gojjam	North Shewa	Priority Weight
West Gojjam	(1,1,1)	(2/3,1,3/2)	(3/2,2,5/2)	0.45
East Gojjam	(2/3,1,3/2)	(1,1,1)	(2/3,1,3/2)	0.32
North Shewa	(2/5,1/2,2/3)	(2/3,1,3/2)	(1,1,1)	0.23

Table 4. Fuzzy Comparison matrix of the alternatives with respect to SC criteria

	West Gojjam	East Gojjam	North Shewa	Priority Weight
West Gojjam	(1,1,1)	(1,1,1)	(3/2,2,5/2)	0.45
East Gojjam	(1,1,1)	(1,1,1)	(2/3,2,5/2)	0.45
North Shewa	(2/5,1/2,2/3)	(2/5,1/2,2/3)	(1,1,1)	0.1

Table 5. Fuzzy Comparison matrix of the alternatives with respect to MP criteria

	West Gojjam	East Gojjam	North Shewa	Priority Weight
West Gojjam	(1,1,1)	(2/3,1,3/2)	(2/3,1,3/2)	0.33
East Gojjam	(2/3,1,3/2)	(1,1,1)	(2/3,1,3/2)	0.33
North Shewa	(2/3,1,3/2)	(2/3,1,3/2)	(1,1,1)	0.33

Table 6. Fuzzy Comparison matrix of the alternatives with respect to SS criteria

	West Gojjam	East Gojjam	North Shewa	Priority Weight
West Gojjam	(1,1,1)	(1,1,1)	(2/3,1,3/2)	0.33
East Gojjam	(1,1,1)	(1,1,1)	(2/3,1,3/2)	0.33
North Shewa	(2/3,1,3/2)	(2/3,1,3/2)	(1,1,1)	0.33

Table 7. Fuzzy Comparison matrix of the alternatives with respect to RP criteria

	West Gojjam	East Gojjam	North Shewa	Priority Weight
West Gojjam	(1,1,1)	(2/3,1,3/2)	(2/3,1,3/2)	0.08
East Gojjam	(2/3,1,3/2)	(1,1,1)	(2/3,1,3/2)	0.39
North Shewa	(3/2,2,5/2)	(2/3,1,3/2)	(1,1,1)	0.53

Table 8. Fuzzy Comparison matrix of the alternatives with respect to PE criteria

	West Gojjam	East Gojjam	North Shewa	Priority Weight
West Gojjam	(1,1,1)	(2/3,1,3/2)	(2/5,1/2,2/3)	0.08
East Gojjam	(2/3,1,3/2)	(1,1,1)	(2/3,1,3/2)	0.39
North Shewa	(3/2,2,5/2)	(2/3,1,3/2)	(1,1,1)	0.53

The priority weights of the alternatives with respect to the criteria were combined and the priority weights of the alternatives were determined. As shown in Table 9, each column of the matrix was multiplied by the priority weight at the top of the column and then those values were added up for each row. At the end, the priority weights of the alternatives with respect cluster selection criteria were calculated.

The priority weights for the alternatives were found to be (0.30, 0.34, 0.36). According to the final score, North Shewa is the most preferred sheep husbandry cluster area as it has the highest priority weight, and East Gojjam is the next recommended alternative for sheep husbandry cluster.

Table 9. The priority weight of the alternatives

	GP	SC	MP	SS	RP	PE	Priority Weight
<i>Weight Alternative</i>	0.28	0.2	0.13	0.13	0.16	0.1	
West Gojjam	0.45	0.45	0.33	0.33	0.08	0.08	0.30
East Gojjam	0.32	0.45	0.33	0.33	0.39	0.39	0.34
North Shewa	0.23	0.1	0.33	0.33	0.53	0.53	0.36

4.5 Sensitivity Analysis

Sensitivity analysis is the study of how the uncertainty in the output of a mathematical model or system (numerical or otherwise) can be apportioned to different sources of uncertainty in its inputs. In this study, it is conducted in order to monitor the robustness of the preference ranking among the alternative areas by changing the priority weights of the criteria. Fifteen trails have been done to justify the result.

In Fig.2, the current situation indicates North Shewa is the most preferable area for sheep husbandry cluster. In case 1, when all the criteria weights are equal, in case 2, when Geographical Proximity (GP) and Sectorial Concentration (SC) criteria increase, and in case 5 when Market Potential (MP) and Support Services (SS) criteria increase, the ranking between East Gojjam and North Shewa is exchanged. On the other hand, when the weights of Resource Potential (RP) and Potential Entrepreneurs (PE) are significantly higher than other criteria, North Shewa became the most preferable area for sheep husbandry cluster. Generally, for most trial cases, North Shewa is the most preferable area for sheep husbandry cluster.

Sensitivity analysis shows that the ranking among the alternatives is quite sensitive to the changes in the weights of cluster selection criteria.

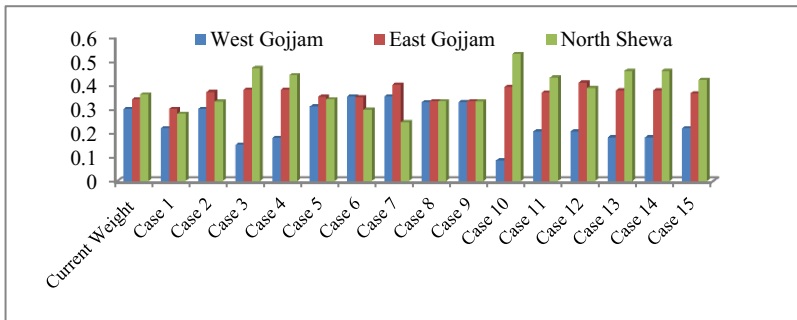


Fig. 2. Sensitivity Analysis

5 Conclusion

In this paper, to minimize defects of hide and skin at animal husbandry stage, livestock husbandry cluster was proposed. Fuzzy-AHP based methodology is selected to identify best origin/area for livestock husbandry cluster. Fuzzy-AHP has been playing an increasingly important role in multiple criteria decision-making under uncertainty. The proposed methodology was tested on a real-world data and was found that it functions satisfactorily. This study identified six criteria (Geographical Proximity (GP), Sectorial Concentration (SC), Market Potential (MP), Support Services (SS), Resource Potential (RP) and Potential Entrepreneurs (PE) for cluster selection and three best origin/areas (West Gojjam, East Gojjam and North Shewa) for livestock husbandry cluster. From these areas North Shewa is selected for sheep husbandry cluster. Here, sensitivity analysis is also performed to discuss and explain the results. As a future study we plan to use FUZZY-AHP to identify industrial as well as Micro and Small Enterprises (MSEs) cluster. We plan also to use other methods for cluster identification and to compare the results with Fuzzy-AHP results.

References

1. Bisrat, G.: Defect Assessment of Ethiopian Hide and Skin: The Case of Tanneries in Addis Ababa and Modjo, Ethiopia. *Global Veterinarian* 11, 395–398 (2013)
2. Calabrese, A., Costa, R., Menichini, T.: Using Fuzzy AHP to manage intellectual capital assets: an application to the ICT service industry. *Expert Systems with Applications* 40, 3747–3755 (2013)
3. Chang, Y.D.: Applications of the extent analysis method on fuzzy AHP. *European Journal of Operational Research* 95, 649–655 (1996)
4. Choua, Y.C., Sunb, C.C., Yenc, H.Y.: Evaluating the criteria for human resource for science and technology (HRST) based on an integrated fuzzy AHP and fuzzy DEMATEL approach. *Applied Soft Computing* 12, 64–71 (2011)
5. Choudhary, D., Shankar, R.: An STEEP-fuzzy AHP-TOPSIS framework for evaluation and selection of thermal power plant location: A case study from India. *Energy* 42, 510–521 (2012)
6. Durán, O.: Computer-aided maintenance management systems selection based on a fuzzy AHP approach. *Advances in Engineering Software* 42, 821–829 (2011)
7. Isaai, M.T., Kanani, A., Tootoonchi, M., Afzali, H.R.: Intelligent timetable evaluation using fuzzy AHP. *Expert Systems with Applications* 38, 3718–3723 (2011)
8. Kahraman, C., Cebeci, U., Ruan, D.: Multi-attribute comparison of catering service companies using fuzzy AHP: The case of Turkey. *International Journal of Production Economics* 87, 171–184 (2004)
9. Kilincci, O., Onal, S.A.: Fuzzy AHP approach for supplier selection in a washing machine company. *Expert systems with Applications* 38, 9656–9664 (2011)
10. Lee, K.S., Mogi, G., Zhuolin, L., Hui, S.K., Lee, K.S., Hui, N.K., Park, Y.S., Ha, J.Y., Kim, W.J.: Measuring the relative efficiency of hydrogen energy technologies for implementing the hydrogen economy: An integrated fuzzy AHP/DEA approach. *International Journal of Hydrogen Energy* 36, 12655–12663 (2010)
11. Lee, S.K., Mogi, G., Hui, K.S.: A fuzzy analytic hierarchy process (AHP)/data envelopment analysis (DEA) hybrid model for efficiently allocating energy R&D resources: In the case of energy technologies against high oil prices. *Renewable and Sustainable Energy Reviews* 21, 347–355 (2013)
12. LIDI: Profile of the Ethiopian Leather Industry Development Institute, Addis Ababa (2010)
13. Netsanet, J., Birhanu, B., Daniel, K., Abraham, A.: AHP-Based Micro and Small Enterprises' Cluster Identification. In: *Fifth International Conference on Soft Computing and Pattern Recognition* (2013)
14. Netsanet, J., Daniel, K., Jakub, S., Svatopluk, S., Vaclav, S.: Application of Fuzzy-AHP for Industrial Cluster Identification. In: *IBICA*, pp. 323–332 (2014)
15. Pedro, C.O., Hécio, M.T., Márcio, L.P.: Relationships, cooperation and development in a Brazilian industrial cluster. *International Journal of Productivity and Performance Management* 60, 115–131 (2011)
16. Porter, M.: Clusters and the new economics of competition. *Harvard Business Review* 76, 77–90 (1998)
17. Porter, M.: *The Competitive Advantage of Nations*. The Free Press, New York (1990)
18. Shamsuzzaman, M., Ullah, A.M.M.S., Bohez, L.J.: Applying linguistic criteria in FMS selection: fuzzy-set-AHP approach 3, 247–254 (2003)

19. Tetsushi, S., Keijiro, O.: Strategy for cluster-based industrial development in developing countries. Foundation for advanced studies on international development and national graduate institute for policy studies (2006)
20. Wang, Y.M., Chin, K.S.: Fuzzy analytic hierarchy process: A logarithmic fuzzy preference programming methodology. *International Journal of Approximate Reasoning* 52, 541–553 (2010)
21. Zheng, G., Zhu, N., Tian, Z., Chen, Y., Sun, B.: Application of a trapezoidal fuzzy AHP method for work safety evaluation and early warning rating of hot and humid environments. *Safety Science* 50, 228–239 (2011)

Thermal Analysis, Design and Experimental Investigation of Parabolic Trough Solar Collector

Yidnekachew Messele and Abebayehu Assefa

School of Mechanical and Industrial Engineering, Addis Ababa Institute of Technology,
Addis Ababa University, Ethiopia
yidn21@gmail.com,
abebayehu_assefa@yahoo.com

Abstract. Energy is one of the building blocks of modern society. Solar energy is a form of renewable energy which is available abundantly and collected unreservedly. In this paper, the application of solar energy using parabolic trough is analyzed. An experimental setup was developed to investigate the performance of the parabolic trough. Measurements of total direct radiation on the plane of the collector, ambient temperature, wind speed, water flow rate, and inlet and outlet temperatures of the water inside the absorber tube were measured and employed in studying the performance of the parabolic trough. A data logger and a computer were employed for data acquisition and the outputs of the experiment are illustrated with the help of graphs plotted using the data import-wizard data import-wizard of the data logger and MATLAB software. Finally, the efficiency which is used as a measure of performance is calculated and the experimental results are compared with the results obtained from the mathematical model.

Keywords: Thermal Analysis, Parabolic Trough, Solar Collector, Experimental Investigation, Solar Energy.

1 Introduction

Energy is one of the crucial inputs for socio-economic development. The rate at which energy is being consumed by a nation often reflects the level of prosperity that it could achieve. The total energy consumption increases with economic and population growth and, at the same time, various environmental problems associated with human activities become increasingly serious.

In addition to an increase in price of fossil fuel products, resources will be exhausted in a relatively short period of time. The current high price of fossil fuel resources is affecting economic and social development worldwide. The impact of energy crises is particularly felt in less developed countries where a high percentage of national budgets for development must be diverted to the purchase of fossil fuel products. To reduce the dependency on imported fuels with high price, most countries

have initiated programs to develop alternative energy sources based on domestic renewable resources. In order to achieve the goals of sustainable development, it is essential to minimize the consumption of finite natural resources and to mitigate the environmental burden to within nature's restorative capacity.

There is now a global consensus that the new sources of energy have to be renewable to satisfy the global energy demand in the long term. Solar thermal power plants are one of the most promising options for renewable electric power production. Unlike traditional power plants, concentrating solar power systems provide an environmentally friendly source of energy, producing virtually no emissions and consuming no fuel other than sunlight.

Ethiopia is one of the countries which is found around the equator in which a better solar radiation exists that creates favorable conditions for the exploitation of solar energy. This can make parabolic trough solar power generation system optional for power production in the country.

2 Statement of the Problem

Ethiopia, in addition to the persistent food insecurity, is suffering from energy supply. It is observed through studies and recent data, energy consumption increases proportionally to the gross national product. One of the possible methods of overcoming energy crisis is by increasing the use of freely available renewable energy sources such as solar energy.

Because of the proximity to the equator, Ethiopia receives adequate sunshine throughout the year. The annual average daily radiation in Ethiopia reaching the ground is about $5.2\text{kWh/m}^2/\text{day}$. The minimum annual average radiation for the country as a whole is estimated to be $4.5\text{ kWh/m}^2/\text{day}$ in July to a maximum of $5.55\text{ kWh/m}^2/\text{day}$ in February and March [1].

Many industries in the country use fuel for water heating process. However, energy costs for heating water is increasing at considerable rate by increasing operating costs and reducing profitability due to continuous escalating of fuel price. Hence, this research explores solar energy as a sustainable alternative for large scale water heating and power generation and it is a step forward to reduce dependency on imported oils.

3 Significance of the Research

Energy is one of the current issues of the country. There is a large energy demand in the country and to fulfill this demand the government is working on different sources of energy.

The result from the study explores solar energy as a sustainable alternative for large scale water heating and power generation and it is a step forward to reduce dependency on imported oils.

The study benefits different industries in reducing the cost related to fuel and other high cost energy sources. This in turn develops the energy alternatives for the industry as well as the country. In addition, this research can be used as a reference for further study in the area.

4 Solar Energy

Solar radiation, often called the solar resource, is a general term for the electromagnetic radiation emitted by the sun. Solar radiation can be captured and turned into useful forms of energy such as heat and electricity, using a variety of technologies. However, the technical feasibility and economical operation of these technologies at a specific location depend on the available solar resource.

From the rays of the sun, which pass through the earth's atmosphere to the ground, a portion is scattered by particles or clouds. The intensity of solar radiation outside the atmosphere is about 1.3 kW/m^2 . Even though only a fraction of this actually hits the earth's surface, the magnitude of the energy from this source is enormous. For example, utilizing only 1% of the earth's deserts and applying a conversion efficiency of 15% to produce electric energy would develop more electricity than is currently produced worldwide by fossil fuels [3]. This is not practical given the need to distribute the electricity to users around the world, but it does highlight the magnitude of this resource.

Global radiation is radiation energy incident on a surface, which is comprised of a diffuse (scattered) component and a direct normal component (the part coming undisturbed directly from the sun). Figure 1 illustrates the definitions of global, diffuse and Direct Normal Radiation (DNR).

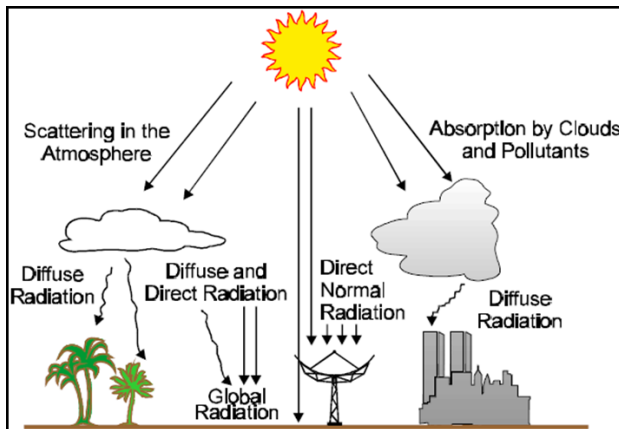


Fig. 1. Direct normal, diffuse and global radiation [3]

The main components of a CSP system are:

4.1 Integration of Parabolic Troughs with Other Energy Systems

As discussed earlier, the Solar Electric Generating System (SEGS) is fundamentally a steam turbine power plant in which the main fuel is solar radiation. Figure 2 shows a schematic diagram of a typical plant configuration [3].

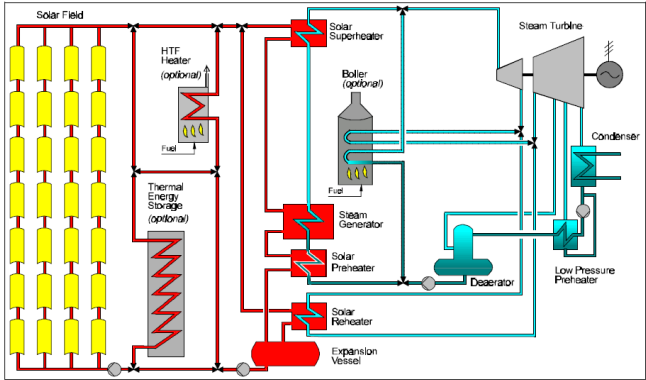


Fig. 2. Parabolic troughs integrated with steam power plants [12]

5 Mathematical Modeling of the Parabolic Trough

In analyzing the solar parabolic collector, it is important to identify each and every part of the collector and the terms used on the solar collector.

In the concept and design of the parabolic collector, the first definition is strictly geometric as ratio of aperture area to receiver area. The ratio of these two areas defines the concentration ratio of the parabolic trough as [24]:

$$C = \frac{A_a}{A_r} \quad , \quad A_a = W_a \times L \quad , \quad A_r = \pi DL \quad (1)$$

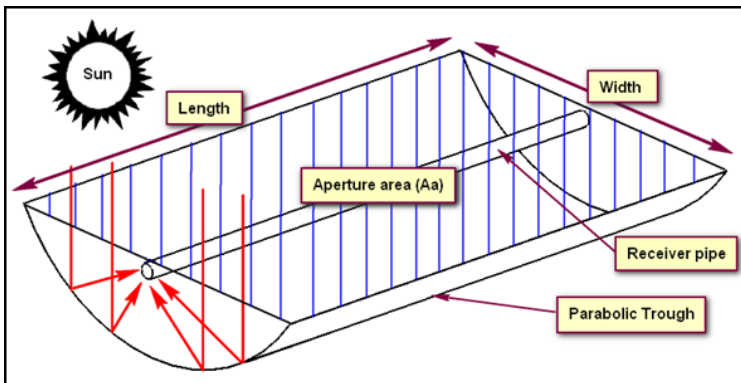


Fig. 3. Parabolic trough

5.1 Heat Collecting Element (HCE) Performance Model

The HCE performance model is based on an energy balance on the collector and the HCE. The energy balance includes the direct normal solar irradiation incident on the collector, optical losses from both the collector and HCE, thermal losses from the HCE and the heat gained by the HTF. Temperature gradient on the receiver can be accounted for by a flow factor FR to allow the use of inlet fluid temperature in energy balance equation. Thus, it is required to derive appropriate expressions for the collector efficiency factor F' , the loss coefficient U_L and the heat removal factor FR to numerically evaluate the outlet temperature. All the equations and relationships used in one-dimensional HCE performance models are described in the following sections.

5.2 One-Dimensional Energy Balance Model

The HCE performance model uses an energy balance between the HTF and the atmosphere, and includes all equations and correlations necessary to predict the terms in the energy balance, which depend on the collector type, HCE condition, optical properties and ambient conditions.

Figure 4 shows the one-dimensional steady-state energy balance for a cross-section of an HCE and Figure 8 shows the thermal resistance model and subscript definitions. For clarity, the incoming solar energy and optical losses have been omitted from the resistance model. The optical losses are due to imperfections in the collector mirrors, tracking errors, shading and mirror and HCE cleanliness. The effective incoming solar energy (solar energy minus optical losses) is absorbed by the selective coating ($\dot{q}'_{3SolAbs}$). Some of the energy that is absorbed by the selective coating is conducted through the absorber ($\dot{q}'_{23SolAbs}$) and transferred to the HTF by convection (\dot{q}'_{12Conv}); the remaining energy is transmitted back to the environment by convection (\dot{q}'_{35Conv}) and radiation (\dot{q}'_{34rad}).

Table 1. Heat flux definitions

Heat Flux Heat Transfer Path [W/m]	Heat Transfer Mode and Transfer Path
$\dot{q}'_{3SolAbs}$	Solar irradiation absorption from incident solar irradiation to outer absorber pipe surface
\dot{q}'_{23Cond}	Conduction heat flux from outer absorber pipe surface to inner absorber pipe surface
\dot{q}'_{12Conv}	Convection heat flux from inner absorber pipe surface to heat transfer fluid.
\dot{q}'_{35Conv}	Convection heat from outer absorber pipe surface to ambient
\dot{q}'_{34rad}	Heat radiation from outer absorber pipe surface to sky

The model assumes all temperatures, heat fluxes, and thermodynamic properties are uniform around the circumference of the HCE. All heat flux directions shown in Figure 7 are positive and all terms indicated in the above paragraph are defined in Table 1. Dotted variables indicate rates and the prime indicates per unit length of receiver and a double prime indicates per unit normal aperture area.

The energy balance equations are determined by conserving energy for each surface of the HCE cross section, referencing Figure 4.

$$\dot{q}'_{12conv} = \dot{q}'_{23cond} \tag{2}$$

$$\dot{q}'_{3SolAbs} = \dot{q}'_{23Cond} + \dot{q}'_{35Conv} + \dot{q}'_{34rad} \tag{3}$$

$$\dot{q}'_{HeatLoss} = \dot{q}'_{35Conv} + \dot{q}'_{34rad} \tag{4}$$

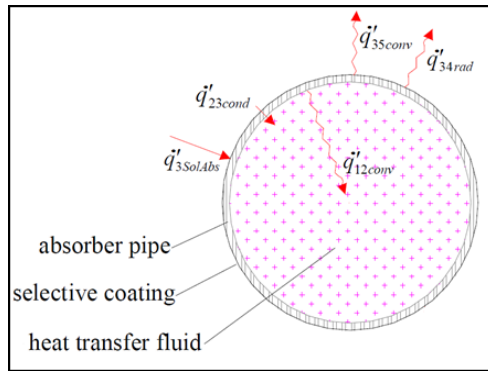


Fig. 4. One-dimensional steady-state energy balance

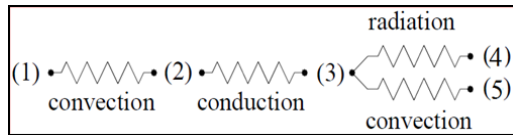


Fig. 5. Thermal resistance model for a cross-section of an HCE

In Figure 5, Point 1 is heat transfer fluid, Point 2 is absorber inner surface, Point 3 is absorber outer surface, Point 4 is sky and Point 5 is surrounding air.

5.3 Convection Heat Transfer between the HTF and the Absorber

From Newton’s law of cooling, the convection heat transfer from the inside surface of the absorber pipe to the HTF is:

$$\dot{q}'_{12Conv} = h_1 D_1 \pi (T_2 - T_1) \quad (\text{W/m})$$

$$h_1 = Nu_{D1} \frac{k_1}{D_1} \quad (5)$$

In these equations, both T_1 and T_2 are independent of angular and longitudinal HCE directions, as will be all temperatures and properties in the one-dimensional energy balance model.

5.4 Conduction Heat Transfer through the Absorber Wall

Fourier's law of conduction through a hollow cylinder describes the conduction heat transfer through the absorber wall [26].

$$\dot{q}'_{23Cond} = 2\pi k_{23} (T_2 - T_3) / \ln(D_2 / D_1) (W / m) \quad (6)$$

In this equation, the conduction heat transfer coefficient is constant and is evaluated at the average temperature between the inner and outer surfaces.

5.5 Heat Transfer from the Absorber Wall to the Atmosphere

The heat will transfer from the glass envelope to the atmosphere by convection and radiation. The convection will either be forced or natural, depending on whether there is wind. Radiation heat loss occurs due to the temperature difference between the glass envelope and sky.

Convection Heat Transfer

The convection heat transfer from the glass envelope to the atmosphere (q''_{35Conv}) is the largest source of heat loss, especially if there is wind. From Newton's law of cooling

$$\dot{q}'_{35Conv} = h_{35} D_2 \pi (T_3 - T_5) \quad (\text{W/m})$$

$$h_{35} = Nu_{D2} \frac{k_3}{D_2} \quad (7)$$

The Nusselt number depends on whether the convection heat transfer is natural or forced (i.e no wind or with wind). Since the experimental setup is in the open field, the convection heat transfer is assumed to be forced (with wind).

Radiation Heat Transfer

The useful incoming solar irradiation is included in the solar absorption terms. Therefore, the radiation transfer between the outer surface of the tube and sky is caused by the temperature difference between the outer surface of the tube and the sky. To approximate this, the outer surface of the tube is assumed to be a small convex gray

object in a large blackbody cavity (sky). The net radiation transfer between the glass envelope and sky becomes [26].

$$\dot{q}'_{34rad} = \sigma\pi D_2 \epsilon_3 (T_3^4 - T_4^4) \text{ (W/m)} \tag{8}$$

5.6 Solar Irradiation Absorption

Using basic energy balance equation, the useful energy gained per unit collector length expressed in terms of the local receiver tube temperature and the absorbed solar radiation per unit of the aperture area, which is the difference between the absorbed solar radiation and the thermal loss and is given by:

$$\dot{q}'_{used} = \frac{A_a \dot{q}_{ab} - A_i U_L (T_i - T_a)}{L} \tag{9}$$

$$\dot{q}_{ab} = \alpha_o I_t \tag{10}$$

$$A_i = \pi D_2 L \tag{11}$$

In terms of the energy transferred in to the working fluid at local fluid temperature:

$$\dot{q}'_{used} = \frac{\left(\frac{A_i}{L}\right)(T_i - T_f)}{\frac{D_2}{h_i D_1} + \frac{D_2}{2k} \ln\left(\frac{D_2}{D_1}\right)} \tag{12}$$

Rewriting Equation 7 in the form of Tt and substituting in Equation 10, the following equation is obtained:

$$\dot{q}'_{used} = F' \frac{A_a}{L} \left[\dot{q}_{ab} - \frac{A_i}{A_a} U_L (T_f - T_a) \right] \tag{13}$$

F' is collector efficiency factor which is given by:

$$F' = \frac{\frac{1}{U_L}}{\frac{1}{U_L} + \frac{D_2}{h_i D_1} + \frac{D_2}{2k} \ln\left(\frac{D_2}{D_1}\right)} \tag{14}$$

This can be rewritten in the form of:

$$F' = \frac{U_o}{U_L} \tag{15}$$

The actual useful energy collected by fluid is, therefore, given by:

$$\dot{q}'_{used} = F_R \left[\frac{A_a \dot{q}_{ab} - A_i U_L (T_{fi} - T_a)}{L} \right] \tag{16}$$

Where F_R is the collector heat removal factor, defined as the ratio of the actual useful energy gain to the useful energy gain, if the entire collector was at the fluid inlet temperature T_{fi} and it is expressed as:

$$F_R = \frac{\dot{m}C_{pf}(T_{fo} - T_{fi})}{\frac{A_a}{L} \left[\dot{q}_{ab} - \frac{A_a}{A_a} U_L (T_{fi} - T_a) \right]} \quad (17)$$

After rearranging the above equation, including the collector efficiency factor, it becomes:

$$F_R = \frac{\dot{m}C_{pf}}{A_a U_L} \left[1 - \exp \left(- \frac{A_a U_L F'}{\dot{m}_f C_{pf}} \right) \right] \quad (18)$$

Finally, rearranging the above equations in the form of T_{fo} , then the exit temperature of the water from the heat collecting tube can be calculated from the following equation:

$$T_{fo} = T_{fi} + \frac{\dot{q}'_{used}}{\dot{m}_f C_{pf}} \quad (19)$$

6 Setup of Experimental Components

6.1 Parabolic Trough Stand

The stand, four legged, holds all the components: the parabolic trough support, the ratchet mechanism and tracing mechanism, up right from the floor. It has three parts: the lower part of the stand is connected to the concrete foundation using anchoring bolts. The stand is linked to the upper part of the trough by a mechanism that allows 360° (the actual angle of rotation need to trace the sun's position 47.3°) rotation of the trough to trace the sun's monthly position.



Fig. 6. Stand of the parabolic trough

6.2 Trough Support

This part of the parabolic trough connects the lower part of the support to the upper part of the trough. The parabolic troughs are connected to this support using bearings so that it is free to rotate from east to west to trace the solar position. This part also gives a rigid structural support of the trough with the stand. The connection between the stand and trough support is contact connection to let the trough support rotate through 360° over the stand.

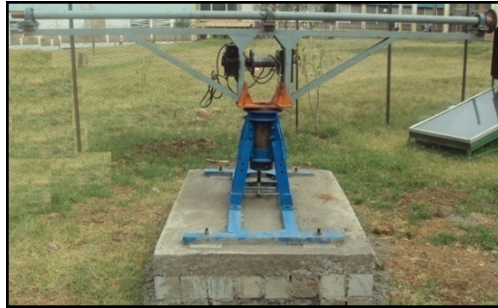


Fig. 7. Trough support

6.3 The Parabolic Trough

The parabolic trough is the most important part of the assembly. The solar radiation strikes the surface of the trough and is reflected to the focal point. The parabolic trough structure is made from RHS metal and angle iron and the reflecting material is aluminum sheet.



Fig. 8. The parabolic trough

6.4 Heat Exchanger

The heat exchanger is one of the components of this experimental setup. After the working fluid is heated by the solar radiation, the heat should be rejected at some point in the experiment setup because the working fluid will circulate again through the collector tube. The heat exchanger that is used for this purpose is shell and tube heat exchanger.



Fig. 9. Heat exchanger

6.5 HTF Pump

The pump is the driving force to circulate the HTF in this experimental setup. The selection of the pump is done considering the layout of the setup, the total pressure loss and the availability of the pump in the market.

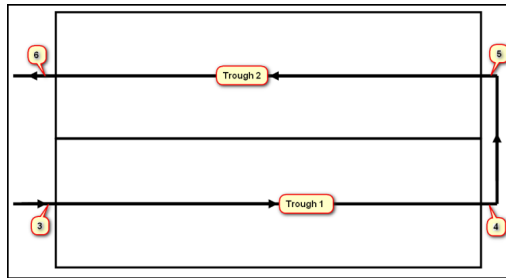


Fig. 10. Top view of the experimental setup layout

Where: 3-Inlet point of the first trough, 4-outlet of first trough, 5-inlet of second trough and 6-outlet of second trough.

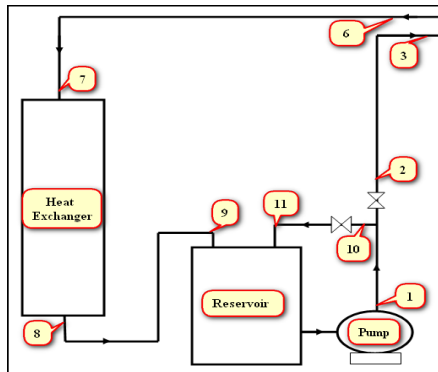


Fig. 11. Front view of the experimental setup layout

Where: 1-outlet of pump, 2-outlet of first flow controlling valve, 7-inlet of Heat exchanger, 8-exit of heat exchanger, 9-inlet of reservoir, 11-outlet of second flow controlling valve.

The parabolic trough is designed to track the sun in any direction. There are two rotational axes to make the tracking system easy. The first rotational axis allows the trough to track the sun from east to west during the day and the second rotational axis allows the trough to rotate in the vertical rotational axis to track the seasonal sun direction change from east to northeast and east to southeast.



Fig. 12. Rotational directions from east to west

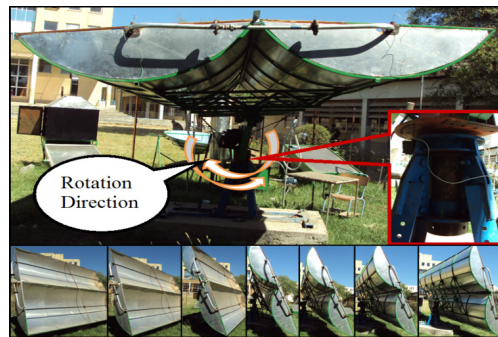


Fig. 13. Rotational Direction with vertical rotational axis

7 Results

In this section, the thermal performance of parabolic trough solar collector is investigated using the measured data of the inlet and outlet temperatures of the working fluid, ambient temperature, wind speed and global radiation around the experimental setup. Plots of the performance parameters are also done using MATLAB software from the data logger import wizard to Microsoft excel worksheet.

Generally, in the test setup, water is circulated through the absorber tube, then to the heat exchanger, finally to the reservoir and again it is pumped back to the absorber tube. A regulating mechanism which is a gate valve is used to alter the flow rate in the system. Temperature of water incoming into and outgoing from the absorber tube are logged using data acquisition system. The pump circulates the water throughout the day starting from sunrise up to sunset. This is in line with the general test procedure in all standards.

The ambient temperature variation in the vicinity of the trough during the day is shown in Figure 16.

The test was started at 8:40 in the morning and ended 16:30 in the afternoon. The maximum temperature measured was 22.99°C between 14:15 and 15:03 in the afternoon and the minimum temperature measured was 17.9°C at 8:40 in the morning.

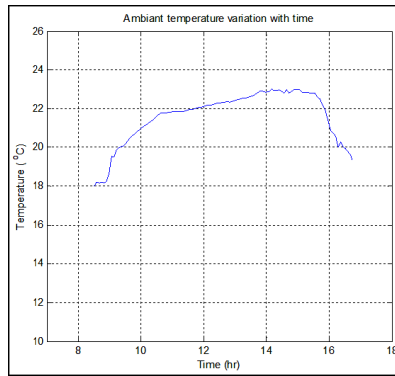


Fig. 14. Variation of ambient temperature

The direct solar radiation variation is shown in Figure 17. The maximum radiation measured was $1049.519\text{W}/\text{m}^2$ and the minimum value $302.32\text{ W}/\text{m}^2$.

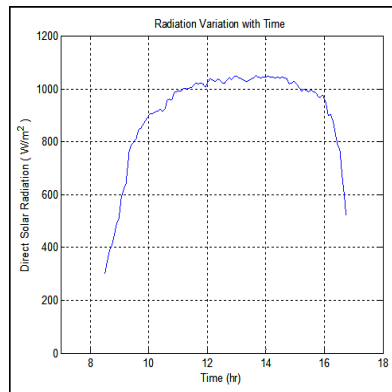


Fig. 15. Variation of Direct Solar Radiation

The temperature variations of the two troughs are shown in Figure 18. The first trough increased the water temperature to a maximum of 55.95°C , which remained above 55°C for two working hours, starting from 12:15 up to 14:23, and above 50°C for five and a half working hours starting from 10:17 up to 15:42. The second trough increased the water temperature to a maximum of 73°C , raising the temperature above 70°C for about five working hours starting from 10:42 up to 14:30.

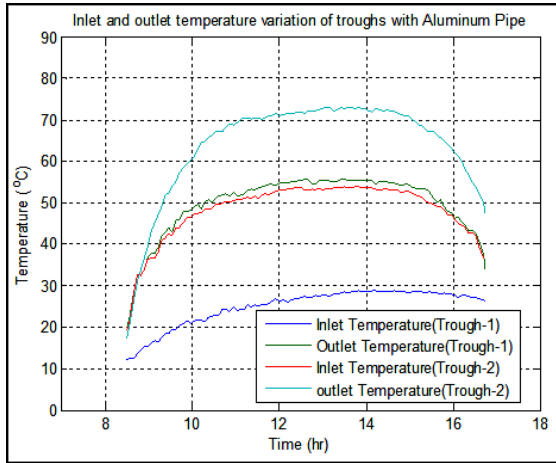


Fig. 16. Inlet and Outlet Temperature Variation of the water

The efficiency of the aluminum heat collector pipe fluctuates between 50% - 60% starting from 10:00 up to 16:00, more than six working hours.

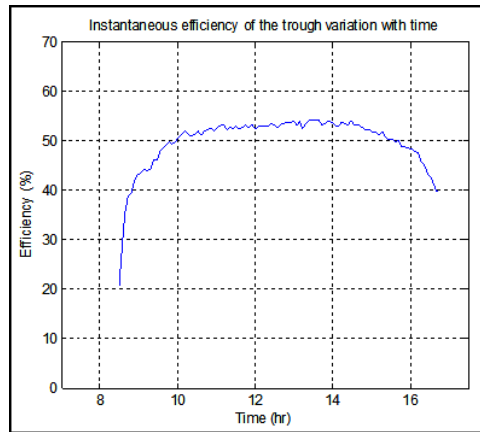


Fig. 17. Instantaneous efficiency of the aluminum heat collecting pipe

The test results are also compared with the mathematical model analysis. From graphs shown in Figure 20, the temperature difference between the analytical model and the actual test results are seen to lie between 0 and 10⁰C. These deviations can be accounted for by the manufacturing and test procedures errors.

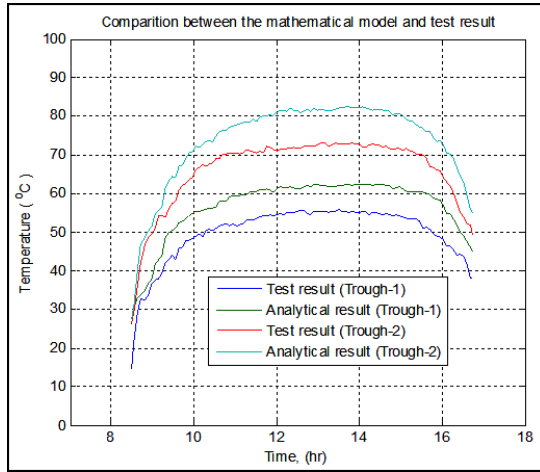


Fig. 18. Comparisons between the mathematical model and test result

8 Conclusion

The major aim of this project work was to design, manufacture and conduct an experimental investigation on the performance of parabolic trough and prepare a mathematical model to verify the results obtained during the test period

This paper has attempted to highlight the following issues

- Based on the current status of the country a new method of a solar energy application has been tested and new technical and technological opportunities of a solar energy application in water heating and steam generation has been established.
- Aluminum sheet material bends easily into the required parabolic trough shape. Black painted aluminum pipe is used as absorber tube. Temperature sensor thermocouples measure the changing water temperature at inlet and outlet of the central receiver. Daily data were collected from each material used as the absorber. Water temperature does increase in the absorber.
- On a clear sky day, a maximum of 73 OC and an average of 70OC of water temperature were recorded using aluminum pipe absorber tube.
- From the result, it can be observed that the parabolic solar trough is a very efficient high temperature water generating system for about five and a half working hours, from 10:00 to 15:30.
- The experimental and the analytical results are very comparable with some acceptable differences.
- The environmental factor plays a major role in the performance analysis of the solar collector. Environmental or weather conditions such as wind and scattered clouds conditions are factors that bring down the efficiency of the solar collector.
- The results of this study give guidance for the possible application of parabolic trough for energy generation.

References

1. Development, Ethiopian Rural Energy. Solar and Wind Energy Utilization and Project Development Scenarios (October 2007)
2. Energy, US department of. Energy basics, <http://www.eere.energy.gov>, http://www.eere.energy.gov/basics/renewable_energy/solar_resources.html#
3. International, Pilkington Solar. Status Report on Solar Trough Power Plants. Pilkington Solar International GmbH, Cologne (1996) ISBN 3-9804901-0-6
4. Prairie, M.: Overview of Solar Thermal Technology
5. Ruiz, P.F.: European Research on, Luxembourg, Belgium (2004)
6. Pytlinski, J.T.: Solar Energy Installations for Pumping Irrigation Water (1978)
7. Kreider, J.F., Kreith, F.: Solar Energy Handbook. McGraw Hill, New York (1981)
8. Spencer, L.C.: A comprehensive review of small solar-powered heat engines (1989)
9. Romero, M., Martinez, D., Zarza, E.: Terrestrial Solar Thermal Power Plants: On the verge of Commercialization (2004)
10. Kakac, S., Liu, H.: Heat Exchangers: Selection, Rating and Thermal Design. CRC Press, s.l. (2002) ISBN 0849309026
11. Perry, R.H., Green, D.W.: Perry's Chemical Engineers' Handbook. McGraw-Hill, s.l. (1984) ISBN 0-07-049479-7
12. http://www.engineeringpage.com/technology/thermal/fouling_factors.html, <http://www.engineeringpage.com/>
13. American Society of Heating, Refrigerating and Air Conditioning Engineers. Method of Testing to Determine the Thermal Performance of Solar Collector. Tullier Circle, Atlanta (1991) Issue 1041-2336
14. White, F.M.: Fluid Mechanics. McGraw-Hill, Boston (1991)
15. ASHRAE Handbook (2001)
16. Devices, Delta-T. User Manual for Temperature Probes. s.n., Cambridge (November 1996)
17. Corporation, National Instruments. User Guide and Specifications NI cDAQ-9172. s.n., Texas (2008) 371747F-01
18. Corporation, National Instrument. LabVIEW Fundamental. s.n., Texas (2007)
19. Valan Arasu, A., Sornakumar, T.: Design, Manufacture and Testing of Fiberglass Reinforced, vol. 81. Science Direct, Tamilnadu (2007) ISSN: 0038092X
20. Forristall, R.: Heat Transfer Analysis and Modeling of a Parabolic Trough Solar Receiver Implemented in Engineering Equation Solver. s.n., Colorado (2003) DE-AC36-99-GO10337
21. Incropera, F.P., DeWitt, D.P.: Fundamental of Heat and Mass Transfer
22. Velautham, S.: Renewable Energy Powered Organic Rankine Cycle (2006)

Performance Improvement by Scheduling Techniques: A Case of Leather Industry Development Institute

Abduletif Habib, Kassu Jilcha, and Eshetie Berhan

Addis Ababa University, Addis Ababa Institute of Technology, Addis Ababa, Ethiopia
{abduinda, jkassu}@gmail.com, eshetie_ethio@yahoo.com

Abstract. The model leather products manufacturing factory of leather industry development institute (LIDI) suffers from poor performances due to various problems. The purpose of this study is, therefore, to improve the performance of the case company using scheduling techniques. Proper scheduling technique can result in dramatic improvements in layout, utilization, idle time, make span and tardiness reduction. The existing company performance and various another scenarios were analyzed by using different sequencing rules plus Johnson's and Campbell's algorithms. The analysis and discussion showed that the feasible scheduling was of flow shop and while product layout was seen most preferable that result in reduction of machine idle time & make span by 3.00 & 4.33 hours respectively. Total flow time was reduced by 82.9% and machine utilization was improved by 16.15% when compared with existing layout. Through production lines 1 or 2 of scenario-2 with the sequence of J₁, J₂, J₃, J₄ and J₅, the company should make possible arrangements for such improvements.

Keywords: performance, scheduling, manufacturing, make-span, leather industry.

1 Introduction

Many companies fail to improve their work performances due to lack of knowledge how to improve it. In turn it results in failure to meet customer requirements in their product on time delivery and products quality as many studies finding proved. Therefore, this research paper focuses on performance improvement technique for the model leather products manufacturing factory of leather industry development institute using scheduling techniques optimize performance so that customers get satisfaction and increase their confidence level. To support these problems, it is important to summarize the previous results and conclusion.

Performance refers to the way people or machineries do their jobs and the results of their works [1]. Performance improvement is methodology to find the root causes of a performance problem and make improvements [2]. Performance analysis is one of the promising tools employed for assessing performance and then taking the necessary measures. The major performance improvement tools in an industry, manufacturing or service, are business process reengineering(BPR), business score card (BSC), benchmarking, ISO 9001:2000, scheduling, assembly line balancing, queue system, total quality management (TQM) and so on [3]. Before selecting particular

performance improvement tool, enterprises need to be clear in the aim to improve and expected outcomes and types of improvements required, holistic or specific area.

Table 1. Performance improvement tools comparison

Tools	Scope	Main im- provements	Time in months	Resources
BPR	Process & people	Process	6-12	high investment
BSC	Holistic model.	Cascading	4-6	Low investment
ISO 9000:2000	Processes and documentation	Process control	6-9	Assessment cost
TQM	Processes and products	Continuity	24-36	Consultant costs
Kaizen	Processes or func- tions	Attitude	24-36	Consultant costs
Scheduling	Processes and people	Make span	24-36	Scheduler costs
Line ba- lancing	Processes and facilities	Flow	36-60	Consultant costs
Queue model	Processes, system and people	Service	48-60	Technology costs

Scheduling is the process of allocating available production resources to complete a certain set of tasks in a given time period while satisfying one or more objectives. It can also decide the work flow, utilization and the layout. Hence, it has been decided to improve the performance of the model factory using relevant scheduling techniques with time parameter. Some of these techniques include job shop (a firm specializes in low to medium volume production and utilizes job or batch processes) the focus of which is to control scheduling bottleneck machines (the n on 1 problem), scheduling parallel machines (the n on m equal problem), scheduling serial processes (n on 2 or n on m problem), and scheduling goals. Another type, flow shop scheduling, is when a firm specializes in medium to high volume production and utilizes line or continuous processes. The rules that specify the job processing sequence when several jobs are waiting in line at a workstation are total shop time, earliest due date (EDD), first-come, first-served (FCFS), shortest processing time (SPT) and critical ratio (CR), a ratio that is calculated by dividing the time remaining until a job's due date by the total shop time remaining for the job [4]. Oliver Braun (2013), also worked on single-processor scheduling with time restrictions but still considered only single parameter and showed that the integrated model can result in better strategic planning decisions in terms of expected profit and conditional value at risk compared to traditional modeling approaches [5]. Camilo et al. (2009) developed an efficient numerical procedure to calculate mean job flow times and then solved for the optimal scheduled starting times using non-linear programming [6].

In specific sub-sectors such as leather products manufacturing, there were few study results that considered the improvements though the sub-sector has high contribution to GDP (Gross Domestic Product) of Ethiopia [7]. Furthermore, most of the

researchers consider only single parameter or constraint which will not come out with significant improvements. In the institute, problems with the performance of model factory can be tackled through performance-based scheduling techniques [8].

The discussion in this paper is focused on proper selection of scheduling techniques that was brought about improvements and significantly benefit the workshop to handle continuously increasing new design and product demand and come up with multiple improvements in layout, utilization, makespan, and efficiency. This solves the makespan time to optimum.

2 Material and Methods

The used methodology was considering the existing production layout type to identify processing time at each workstation. Then, using sequencing rules the sequence that minimizes makespan and utilization of resources at production lines were identified by direct observation followed by self-investigation such as interviewing and discussion to the relevant persons, especially designers, pattern makers, supervisors and other staff in the workstations. Five major jobs were selected and average process time and order due date and ten order times sample data in each workstation was collected randomly to determine activity delays and the number of delayed jobs.

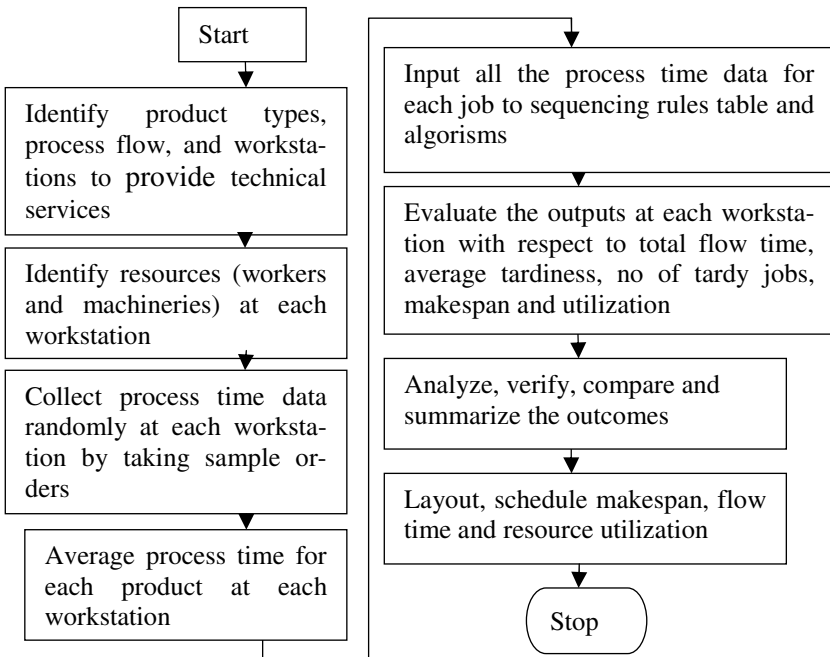


Fig. 1. Flow chart of methodology

Sequencing rules and Johnson’s algorithms were used and results were compared through critical ratio (CR). Then makespan was verified by using the Gantt chart. The validation process of the model was conducted by comparing the output of proposed layouts and that of the existing system. Finally, based on the output comparison, proper scheduling type, layout and resource utilization could be defined.

3 Discussions and Results

3.1 Existing Factory Layout and Production Process Flow

The existing condition contains two workstations. These are CAD/CAM center (design room) in workstation-1 and workstation-2 (leather products production section) in which cutting, skiving and stitching operations are being performed. The design outcomes from the CAD/CAM center are then moved to workstation-2 for further processing and or transferred to customers (producers industries) directly so that they can reproduce in mass for the market. By considering the former case where the designed patterns are transferred to workstation-2, the pattern and design specification sheet is distributed for the respective cutting, then to skiving & stitching operations. Let’s evaluate the existing method (Fig 2).

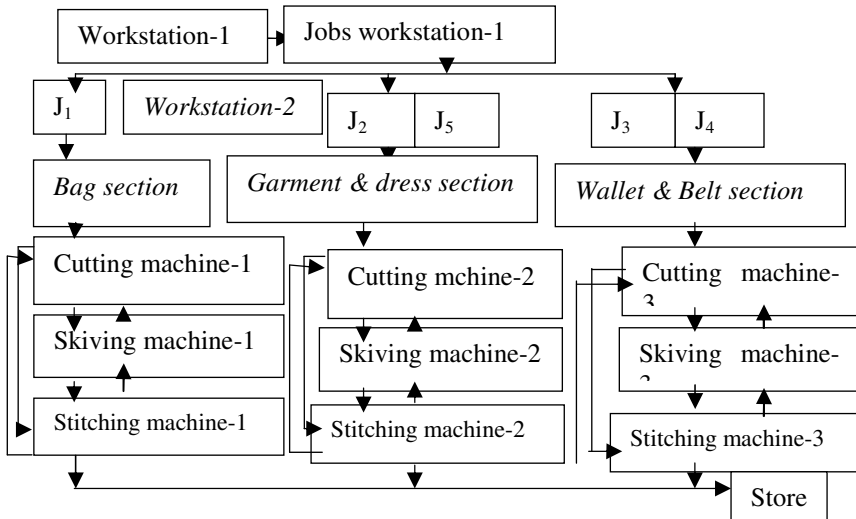


Fig. 2. Existing arrangements of sections and machines in workstations

3.2 Workstation-1 (CAD/CAM Center) – Scheduling Jobs on CAD/CAM Machine

The CAD/CAM center is always busy and has at least 5 unprocessed jobs (leather bag, leather garment, leather wallet, leather belt and dress) design & pattern making, respectively assigned as: J₁, J₂, J₃, J₄ & J₅. The scheduling is done separately using

dispatching rules and just to begin, data for process time and due date for two weeks (10 days or order intervals) is collected from product development center and an average is taken for evaluation (table 2). It takes only a couple of minutes to insert specification and trace the pattern for a single product on CAD/CAM machine. Here given that, (the job, process time, due date) as :(1,35,56),(2,49,38),(3,31,45),(4,18,30) and(5,35,33) respectively, by sequencing rules (FCFS, SPT, EDD and CR) the evaluation summary was obtained as follows (table 2).

Table 2. CAD/CAM Machine against dispatching rules for existing condition

Rules	CAD/CAM Machine			
	Mean flow time	Aver. Tardiness	No of tardy jobs	Utilization
FCFS	113	76.8	4	29.7%
SPT	87.0	50	4	38.36%
EDD	117.8.	80.3	4	28.5%
CR	94.6	56.5	5	35.5%

3.3 Workstation-2 and Jobs

Considering workstation-2 (Fig 2) in the m-machine problem with ‘n’ different products, we have potentially (n!)^m different schedules though full enumeration is usually impossible. Also, jobs processed in the section are unrelated and each job is processed on respective sections sequentially i.e. cutting, skiving and then stitching respectively but some components may repeatedly visit the same machines. It also violates the Johnson’s rule that the job time must be known and constant and job times must be independent of sequence (but in our case, the sequence is a must) and jobs must follow same two-step sequence plus job priorities cannot be used.

3.4 Proposed Product Layout Type and Production Flow

Scenario 1: Considering job shop flow where 5 jobs are done on identical parallel machines in batches in workstation-2, major assumptions were drawn as identical machines performance may be constant. Similar machines are arranged in similar subsections. Jobs should visit all substations before completion. The same sample is used. In this case the entire batch is completed before they move to the next machine (Fig.3).

By using algorithms to minimize mean flow time and/or make span in each substation (one for minimizing mean flow), sequence jobs by SPT and take jobs from list and assign to each 1st machine in workstation-2, substation-1 with least amount of assigned time and continue through all jobs. The other is to minimize make span while controlling mean flow time and do the same reversing scheduled tasks’ processing order. Students who are training data from workstation-2, substation-1 for 5 jobs in the months of April-May, 2014 was taken and arranged (job, process time in hrs per unit product) as (1,5),(2,6),(3,3),(4,8),(5,7). Using flow time and make span with FCFS, SPT, EDD, LPT and CR rules plus Gantt chart for all substations with the same procedure, the result is shown (table 3).

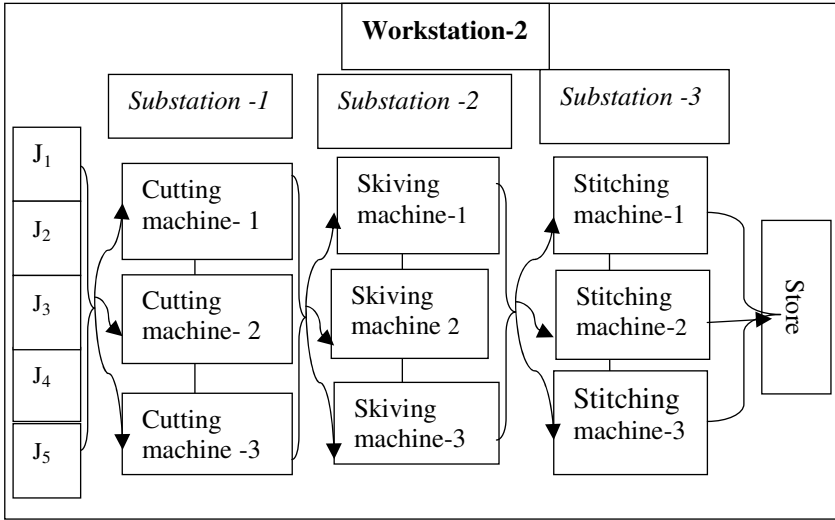


Fig. 3. Parallel identical machines arranged in production sequences

Table 3. Scenario-1 proposed method out comes

Rules	Substation-1			Substation-2			Substation-3		
	Mean flow time	Aver. Tardiness	No of tardy jobs	Mean flow time	Aver. Tardiness	No of tardy jobs	Mean flow time	Average Tardiness	No of tardy jobs
FCFS	17.9	11.2	4	22	12	1	25	13	3
SPT	15	0.5	1	11	0.7	1	25	2	1
EDD	58.5	0	0	14	7	4	27	19	2
CR	16	2.1	3	13	4	3	14	3	4

Scenario-2: Considering line production where n jobs are done in separate lines on production sequence of serial machines, taking into account assumptions such as: Each line contains non identical machines arranged in production flow sequence. A product once entered into a production can never jump to another line. The same sample is used See (Fig 4).

In this serial problem, choosing arbitrarily line-1, average process time for the four consecutive machines are considered. Campbell algorithm reduces $(4!)^5$ machines schedules to 3. CAD/CAM can feed all the serial lines. Time is in hours. Process time data was collected for 5 jobs for the months of April to May, 2014 (table 4).

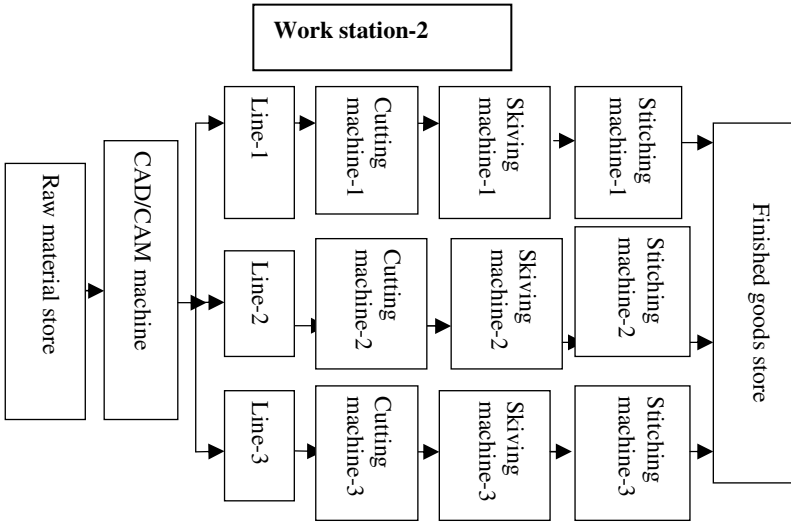


Fig. 4. Serial machine problem/flow shop process

Table 4. Scenario-2 proposed method out comes

Job Order	CAD/CAM Machine	Cutting Machine-1	Skiving Machine-1	Stitching Machine-1
J ₁	5	8	4	3
J ₂	7	9	5	8
J ₃	2	3	9	7
J ₄	6	1	6	4
J ₅	3	4	5	2

For production line-1, Johnson’s rule 1st sequence is: J₃-J₂-J₄-J₁-J₅ with makespan of 40 units. The 2nd and 3rd sequences and makespan are: J₃-J₄-J₂-J₅-J₁ and J₃-J₂-J₁-J₄-J₅, 42 and 45 respectively. Gantt chart also shows the same result. So, makespan is 40 hrs. For production line-2 and line-3, the same procedure followed and resulted in the summary (table 5).

Table 5. Scenario-2 proposed method out comes

Production lines	Performance-based evaluation parameters			
	Mean Flow time	Mean idle time	Makespan	Mean average utilization
Production line-1	12.22	4.3	40	30.29%
Production line-2	15	2.5	40	49.45%
Production line-3	17.5	2.2	42	41.2%
Average	14.90	3.0	40.67	40.31%

To analyze the existing scenario's (table 2) results show that SPT is better to schedule with as it has minimum flow time, tardiness and better utilization. But in workstation-2 the schedule also violates both sequencing and Johnson's rules. Scenario-1 considers all machines in workstation-2 categorized in section-wise (Fig.3) and the result summarized in (table3) shows that SPT schedule has a minimum flow time while LPT schedule has minimum lateness in all subsections which reduces service delays. Makespan for SPT and LPT is same. In scenario-2 (Fig.5) the outcome summary shows that the mean flow time is 14.90 hours and the mean idle time is 3.4 hours (table 5). The mean flow time and mean idle time is less than mean flow time of scenario-1 (table 2). Makespan for scenario-2 is less than that of scenario-1 ($40.67 < 45$) hours. For scenario-1 and scenario-2, SPT is reduced by $(87-51)/87=41.4\%$ and $(87-14.90)/87=82.9\%$ of its hours respectively. Machine utilization for scenario-2 is increased by 16.15% ($(40.31\%-33.8\%)/40.31\%$). Makespan is independent of sequencing rules in existing method but reduced by $(45-40.67)/45=9.62\%$ while idle time is reduced by $(6-3)/3=100\%$ from scenario-1 to scenario-2. Average tardiness and tardy jobs are reduced in SPT of scenario-1 than in existing method of CAD/CAM center by 93.6% $(50-3.2)/50$ and $(4-3)=1$ respectively.

4 Conclusion

This study investigated the appropriate scheduling techniques in the model leather products factory of leather industry development institute by using different dispatching rules and Johnson's algorithms across various scenarios which showed that the minimum flow time, makespan and idle time needed to produce a batch of products were found in flow shop scheduling (scenario-2) as 14.90 hours and the mean idle time was 3.4 hours (table 5). Makespan and flow time were reduced by 9.6% and 82.9% respectively while machine utilization was increased by 16.15%. Product layout was seen best in the selected scenario. Hence scenario-2 is the best alternative to schedule with as it leads to better improvements. Furthermore, the improvements will be good news for leather products manufacturing companies that are consulted under leather industry development institute those are suffering from performance (productivity) under 30% (LIDI Report, 2013) due to different problems like underutilization, and improper layout.

The study by other than Johnson's algorithm and dispatching rules, many factors affecting the scheduling of the model leather products factory will be focused on as future study area. These factors include scheduling with economic-based performance improvement parameters that minimize costs such as set-up costs, inventory holding costs, shortage costs, overhead and labor costs.

References

1. Intra health, n.d., performance improvement stages, steps and tools, <http://www.intrahealth.org/sst/stage2.html> (viewed August 13, 2014)

2. Gruenberg, T.: Performance Improvement: Towards a method for finding and Prioritizing potential performance improvement areas in manufacturing operations. *International Journal of Productivity and Performance Management* 53(1), 52–71 (2004)
3. Cherkos, T.: Performance Analysis and Improvement of Ethiopian Leather Footwear Factories: With Special Reference to Anbessa Shoe S.C. AAU (October 2011)
4. Sosimi, A., Ogunwolu, F.O., Adegbola, T.: A Makespan Optimization Scheme for NP-Hard Gari Processing Job Scheduling Using Improved Genetic Algorithm. *Journal of Industrial Engineering* 1(1), 1–5 (2014)
5. Braun, O.: Single-processor scheduling with time restrictions, University of Applied Sciences Trier (August 2013)
6. Mancilla, C., Storer, R.H.: Stochastic Sequencing and Scheduling of an Operating Rooms. Theses and Dissertations, Lehigh University, Department of Industrial and Systems Engineering (November 14, 2009)
7. Leather Industry Development Institute (LIDI), Leather Industry sector development plan, Performance Report 2010–20014, GTP1, Addis Ababa
8. Benchmark Implementation Plan for the Ethiopian leather products Sector, Pilot Project on Anbessa Shoe Share Company and Peacock Shoe Factory, prepared by ASSC (February 2009)

Response Time Reduction in the Leather Products Manufacturing Industry Using Arena Simulation Method

Haftu Hailu, Kassu Jilcha, and Eshetie Birhan

Addis Ababa University, Addis Ababa Institute of Technology, Addis Ababa, Ethiopia
{hhea192741, jkassu}@gmail.com, eshetie_ethio@yahoo.com

Abstract. The purpose of this study is to find the best strategy for products distribution network in terms of response time reduction. The decision making is based on the carriage capacity of the cars relative to the transportation time and cost records. The response time for the service delivered in the manufacturing company is not optimum and customers are always in complaint. Thus, due to high transportation cost and long time delivery, the institutional customers of the company were also influenced by response time to get the products on time. Optimum transportation time and cost has improved company's competitive strategy during last decades undoubtedly. As premises of this, transportation cost and times are the two variable used to determine the best strategy among the alternatives. A lot of mathematical models have been applied to do optimization on supply chain networks, but supply chain dynamics, such as uncertainty in production demand and transportation, are not present in most of them. As a result, to handle this variation, simulation can be a powerful tool. In this study, an arena simulation tool is used to analyze system performance using transportation cost and time as the performance parameters. In addition, based on the generated transportation time and cost, a simple mixed integer linear programming (MILP) was developed and Win QSB software has been used to solve and select the best network configuration. Finally, the study resulted in strategy 2 has a minimum time and cost. By considering current situation this strategy is able to save the time by nearly 6.5 hours in 30 days.

Keywords: Leather products, response time, arena simulation, production, MILP.

1 Introduction

It is very important to reduce the cycle time of product distribution that may results in reduction of cost of transportation. The main reason to this study motivation is that in the leather manufacturing industry there is high cycle time for product distribution to different centers. Therefore, this study focuses to answer this high cycle time reduction for the leather products distribution considering all possible strategies described in this study.

Response time is the total amount of time it takes to respond to request for service which is the sum of the service and wait time. Thus, factors that affecting response time for achieving predetermined goals is assessing by performance measurement.

Evaluation of supply chain performance is essential for operating and managing of supply chain [1]. Production distribution network as a key driver of the total profitability of companies has considerable impact on the supply chain performance directly [2]. Therefore, evaluation of performance in terms of response time is vital task for companies. Simply, simulation software can evaluate network performance. A few years ago, this method was bounded to predicting outcomes of changes made to equipment, processes or systems within a distribution centers or manufacturing facility. Today, simulation can be much broader encompassing of whole supply chain networks and their numerous facilities, transportation modes and inventory [3].

Configuration of production-distribution network has great effects on a supply chain's long-term performances. In the sense of dynamics inherent of supply chain such as transportation uncertainty and demand variation, they are not presented in nearly all mathematical models because of tractability [4]. A wide range of techniques such as linear programming, integer programming, mixed-inter linear programming, heuristic methods and genetic algorithms have been applied for modeling of production-distribution network. The principles of modeling methods such as simulation, mathematical optimization and heuristic models have been described in detail by reference [5]. In heuristic modeling methods, although acceptable answers are founded, but the answers are not guaranteed as an optimal. These methods are able to create an acceptable solution in a logical time. In decision making for production distribution center (DC), decision criteria is stated considered within two dimensions as Customer needs that are met and cost of meeting customer needs [6]. Response time is one of the customer service criteria which influenced by network design and configuration. Transportation is also supply chain costs which influenced by production-distribution configuration [6]. In this regard, simulation has appeared as a powerful method to investigate and design global strategies for the companies. According to Sara Hewitt, stochastic system use random number generators, so the output of the simulation is an estimate of the true system behavior [10].

Vieira had presented first ideas of a research project proposing to use computer simulation, in this case is Arena, in modeling and evaluation of supply chain performance [11]. The characteristic for network management seen as a complex task. But Arena simulation software could become a useful tool to understand the basics of any net-work management [12].

One of scenarios to improve transportation time and cost in leather products manufacturing industry was by adjusting the number of resources and distance between product distribution centers. And then Arena software was used to analyze the effects of that adjustment [13]. The relationships between key indicators of manufacturing product distribution system performance, such as transportation time and cost to be defined shortly in this work, are complicated and difficult to quantify [20]. For many companies, the key aspects of current competitiveness in order to raise customer satisfaction focus on effective delivery time [21].

A simulation model is a descriptive model of a process a system and usually includes parameters that allow the model to be configurable. The model will then be used to estimate the effects of various actions. Finally simulation draws conclusions and makes action decision based on the results of the simulation. In this way the real life system is not touched until the advantages and disadvantages of what may be major policy decision are first measured on the systems model [23, 24, 25]. Different

authors finding showed that arena simulation is the powerful tool to solve such a stochastic system of the companies [7, 8, 9]. This model also customized to come up with good result.

This study was conducted in Universal leather products manufacturing company producing ladies hand bag (LHB), leather luggage (LLG) and small leather articles (SLA). The products are shipping to distribution center by cars. The factory operates 6 days a week starting 8:00 am to 4:30 pm from Monday to Friday and till 12:30 pm on Saturday. So, now day’s product varieties are rapidly being introduced into the market. In this situation, the company survival is determined by constantly transporting products with optimum time and cost, or by improving an existing product distribution system [22]. By considering to the company’s ability and limited capacity of DC, direct shipping of products will be added to their transportation modes. The amounts of each product which can be sent directly were given by company’s sales experts. The ability of the company for direct shipping of products is 20% of LLG, 16 % of LHB and 10% of SLA.

In this study, the company’s actual cycle time for product distribution is 13.29 hrs. This high transporting time definitely identifies reasonable cost on the company. So, finding the best strategy which minimizes the transportation cost and time is enabling the company to solve the stated high cycle time problem. A simulation model is developed to evaluate system performance based on transportation cost and time for different DC configurations within the company in line with developing Mixed-integer linear programming (MILP) and solving by winQSB software. The existing company transportation time and cost was in random operational way. The different strategies were considered in this study to solve problems of high cost and time for the responsiveness of the company in reducing wastes.

2 Material and Methods

The study was conducted considering different materials and methods to achieve the goal of this study.

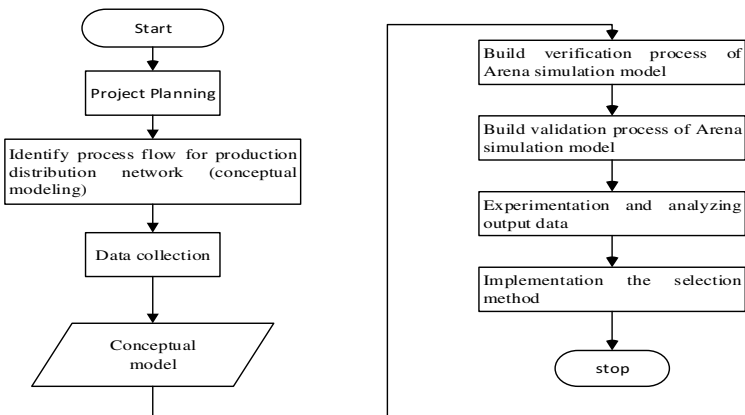


Fig. 1. Flowchart of methodology

The literature reviews and software applications discussed were utilized to come up with good solution regarding the objective of this study. Generally, the first methodology was direct contact to experts to identify transportation time and cost at each production distribution network, and then the second methodology was to use Arena simulation for identifying the minimum transportation time and cost needed for production distribution network. Applied methodology in this study divided in separated nine steps model. It is necessary to complete all these steps, before the simulation study was completed. The application of methodology makes certain a valid simulation outcome and helps the analyzer in the development of the model [17]. These steps are described clearly (Fig 1).

3 Discussion and Results

3.1 Production Process Flow

In general, there are 6 main sections in this company (Fig.2). Raw materials first come to store in the form of bundles for storing & quality checking. Next, the processes of cutting, numbering, skiving and splitting are done. After, preparation stitching and completed products process flow followed. Finally finished products go to inspection & packing to delivered to customers (Fig2).

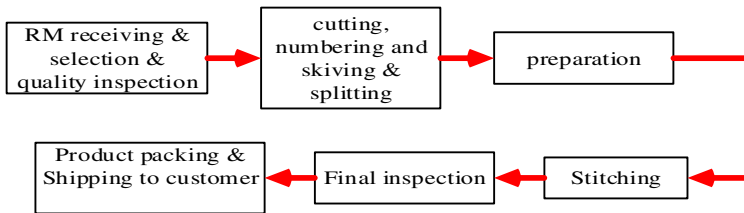


Fig. 2. Process flow of studied leather Product Company

3.2 Conceptual Model

First step in each simulation modeling is building conceptual model.

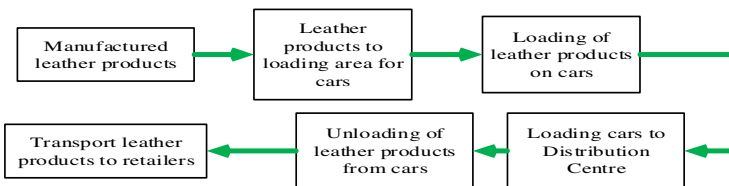


Fig. 3. Process flow of studied company product distribution network

Conceptual model included transportation cars as inputs from manufacturer to retailers. Out puts consisted of transportation time and cost for each type of cars among

the company's distribution network. Products are transported from manufacturer to distribution center by cars. Then, based on retailer/customers demand products are collected and shipped to retailers by their own vehicles. Hence, this study aim is to find the best strategy for transporting products from the manufacturing to customers. In order to create each product (entity) in the simulation model, the Create module was used to generate arrivals of each product. Then this Create module was connected with the Process module.

3.3 Decision Variables

Variables are classified as controllable and uncontrollable. The controllable variables include number of different types of cars that are transported from manufacturer to DC and from DC to different destination of retailers. Uncontrollable variables consist of variables which are not under the control of simulation analyzer such as: transportation cycle time. Both types of variables are important for making decision in terms of the fitted and best strategy of supply chain networks

3.4 Data Analysis

Relevant data's such as loading and unloading time, transportation time from an origin to its corresponding destination and products demand has been collected by interviewing, observing the report of the company with direct contact to the company sales experts and identified to which tool it fits.

In this study chi square test was done to find fitted distribution by using easy fit software. Normally Chi Square test is used to identify fitted statistical distribution. In general, the chi-square test statistic used is of the form:-

$$x^2 = \sum \left(\frac{\text{observed} - \text{expected}}{\text{expected}} \right)^2$$

3.5 Model Assumption

This study represents discrete-event modeling and the factory works for 450 minutes in a day. The actual working time for whole system is 7.30 hours per day. Break time include lunch and tea break are totally 1 hour per day. The manufacturer working time is 8.30 hours per day (Monday up to Friday) and 4.30 hours on Saturday. There is no maintenance process performed during the transportation period. Each car can transport the products and for running in simulation model, it should be waited until its capacity is being filled. The cars never starved for loading and trucks breakdown is not considered in this model.

3.6 Model Verification

Model verification is an important step in simulation modeling. It is applied to ensure the model whether it is running properly or not. In fact verification is a performance appraisal system. Several verification techniques can be applied to the model to make

sure that the model has been conducted correctly and it obtained various purposes of simulation in organizations or companies. These techniques are [15] used to conduct model code reviews, check the output for reasonableness, watch the animation for correct behavior, and use the trace and debug facilities provided in the simulation software. In this simulation modeling, each part of model was run with different set of inputs and the obtained outputs were compared with actual outputs.

3.7 Model Validation

Several techniques for validation of model have been offered by [14]. Model validation required abilities for drawing a conclusion map about the accuracy of the model based on the available evidence. In this study, validation of the model was done by comparing the actual data with the results that were generated by simulation software. In the other words, for validating of the model, the transportation cycle time from manufacturer to DC for a day was simulated and compared with actual time.

Table 1. Validation data

Replication	Simulated cycle time (hours)	Actual cycle time (hours)
1	14.07	13.29
2	13.93	
3	12.62	
4	13.70	
5	12.95	
Average	13.45	
Variation (%)	1.57	

3.8 Warming up Period Determination

In this case study, the system behavior represents steady state so that warm up period should be identified, because entities can be in different parts of the system while it is not working. While the system restarts, entities continue their routes to leave the system. Warm up period model was run for 1 shift with 5 replications. The system needs

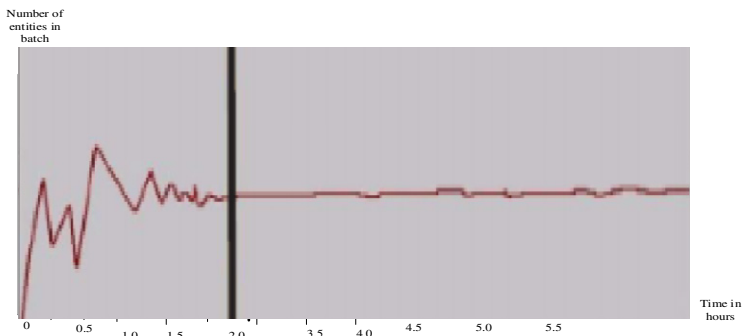


Fig. 4. Warm up period determination

2 hours to be worm up. While warm up period was identified, it is the time to find the number of replications for simulation running. To estimate the number of replications, the model initially was run for 5 times with a run length of 7.30 hours (1 working day).

3.9 Number of Replications

One of the important steps in any validation procedure is to determine the number of replications, because it has directly effect on result accuracy. To determine number of simulation runs, cycle time can be a good estimator [18]. To do so, the formula expressed below was proposed by [19] for estimating number of replications.

$$\sqrt{N(m)} = \left(\frac{S(m)t_{m-1, \frac{1-\alpha}{2}}}{x(m)\epsilon} \right)$$

Where N(m)= number of simulation runs to achieve the desired level of accuracy; X (m) = the mean estimate of an initial m number of runs; S(m)= the standard deviation estimate of m number of runs; α = level of confidence; ϵ = allowable percentage of error; and $t_{(m-1, (1-\alpha)/2)}$ = critical value of the two-tailed t-distribution at a level of significance, given m-1 degrees of freedom.

Mean, $x = (14.07 + 13.93 + 12.62 + 13.70 + 12.95)/5 = 13.45$ and variance, $s = \sqrt{1.57/(5 - 1)} = 0.6$. For 95% Confidence Level and error, $k = \pm 5\%$, $t_{0.025,4}=2.776$.

Number of replication, $n = \left(\frac{ts}{kx} \right)^2 = \left(2.776 * \frac{0.62}{0.05 * 13.45} \right)^2 = 3.82 \approx 4$

Table 2. Estimate of mean and standard deviation

		Simulation Cycle Time (Minutes)	$(xi - x)^2$
Replication	1	14.07	0.34
	2	13.93	0.23
	3	12.62	0.69
	4	13.70	0.06
	5	12.95	0.25
Total		67.27	1.57

Therefore, the minimum number of replications is 4. To have more reliable and valid result we examine our models under 5 time’s replications. Table 3 and table 4 show the generated result which obtained from ARENA. Figure 7 illustrates the average transportation time which has been obtained from simulation model in 5 replications on 30 days of simulation run. Strategies 1, 2 and 3 are related to the direct shipping of product 1, 2 and 3. From strategy 4 up to strategy 6, two types of products are shipped in a direct manner, and finally in the strategy 7 all the products are shipped directly with determined amounts.

Table 3. Transportation time generated by Arena

Replication \ Strategy	1	2	3	4	5
1	140.4	140.5	140.5	140.6	130.5
2	133.9	134.8	135.0	134.2	134.0
3	132.7	132.6	132.9	132.6	132.6
4	138.5	138.4	138.6	138.4	136.6
5	119.8	120.0	120.2	119.9	120.0
6	147.6	146.9	147.9	147.1	147.0
7	154.9	154.7	155.3	154.9	155.0

Table 4. Transportation cost generated by Arena

Replication \ Strategy	1	2	3	4	5
1	720353	727330	718051	721242	720733
2	811438	822053	819396	817652	824231
3	845836	848795	840176	847726	847726
4	751192	754016	745995	753155	751167
5	590692	585234	589665	591223	591223
6	622187	551899	621039	625002	621821
7	557165	553327	554787	553298	552876

Total transportation time reduced in single product direct shipping, and strategy 1 and strategy 3 have minimum and maximum total transportation time respectively (Fig 6). Strategy 4 to strategy 6 multi products direct shipping has been suggested for logistic network. Transportation time dramatically dropped, from strategy 3 to strategy 4. It can be concluded that by increasing the amount of direct shipping, a considerable time reduction has been obtained.

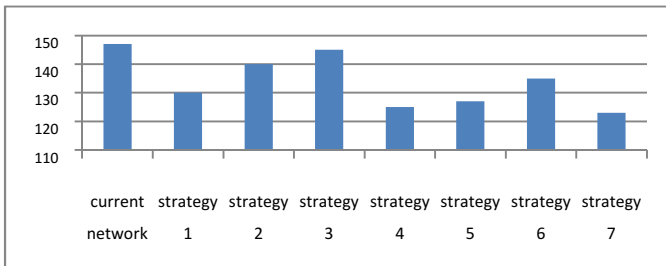


Fig. 5. Comparison of average transportation time between different strategies

Figure 7 illustrates fluctuations in terms of transportation cost among different strategies. Transportation cost has a direct relationship with amount of products which are transported directly. The cost of transportation increases while percentages of direct shipping have additive trend. Since cost and time have different value for decision making to find fit strategy, optimizing model and sensitivity analysis are required.

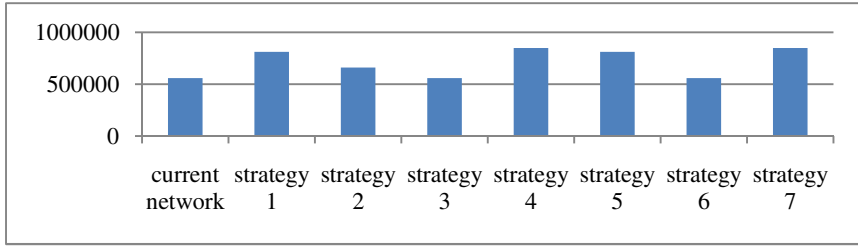


Fig. 6. Comparison of average transportation cost between different strategies

3.10 Selection Method

For selection of the best configuration based on mentioned criteria a simple MILP was developed and solved by WINQSB software. The description of a mixed integer linear programming MILP is as follow:

$$\text{Min } Z = h (162253.96 C_1 + 149421.27 C_2 + 138012.39 C_3 + 182744.34 C_4 + 166676.87 C_5 + 80474.72 C_6 + 188797.16 C_7) + k (89.17 T_1 + 83.63 T_2 + 85.25 T_3 + 25.79 T_4 + 26.24 T_5 + 27.96 T_6 + 26.93 T_7)$$

$$\begin{aligned} \text{s.t.} \\ \sum_{i=1}^n C_i = 1, \quad \sum_{i=1}^n T_i = 1 \\ C_1 - T_1 = 0, \quad C_2 - T_2 = 0, \quad C_3 - T_3 = 0 \\ C_4 - T_4 = 0, \quad C_5 - T_5 = 0, \quad C_6 - T_6 = 0, \quad C_7 - T_7 = 0 \\ C_i, T_i = 0 \text{ or } 1 \text{ (where 0 means not selected and 1 means selected)} \end{aligned}$$

$$h, k, i = \begin{cases} 0 \leq h, k \leq 1 \\ h + k = 1 \\ 1 \leq i \leq 7 \end{cases}$$

According to Win QSB results, strategy 2 has the minimum time and cost based on different weights of transportation time and transportation cost. By considering current situation this strategy is able to save the time by nearly 6.5 hours in 30 days. The percentage of transportation time improvement between the simulation result and the current situation result can be calculated as follows:

$$\begin{aligned} \text{The percentage of } TT \text{ improvement} \\ = \frac{\text{the current situation } TT - \text{the average } TT \text{ generated by arena}}{\text{the current situation } TT} \times 100\% \end{aligned}$$

The percentage of transportation time improvement resulted in this study is supported by the previous research [26].

$$\text{The percentage of } TT \text{ improvement} = \frac{146.5 - 140}{146.5} \times 100\% = 4.43\%$$

It improves the transportation time by 4.43%. On the other hand, this reduction in time is caused to increase transportation cost by: 3.59% during 30 days.

$$\begin{aligned} & \text{The percentage of TC improvement} \\ &= \frac{\text{the average TC generated by arena} - \text{the current situation TC}}{\text{the average TC generated by arena}} \times 100\% \\ \text{The percentage of TC improvement} &= \frac{612000 - 590000}{612000} \times 100\% = 3.59\% \end{aligned}$$

4 Conclusion

This study is aimed to reduce response time to the market. Different alternatives were suggested by the experts to reduce response time. Despite of reducing response time by 6.5 hours and specially transportation time improvement by 4.43%, definitely a reasonable transportation cost which is 3.59% has been identified within 30 days. Conducted simulation model in this study investigated the best strategy configuration which is strategy 2 based on minimum transportation time and cost by using ARENA software. Finally, among the strategies possible alternatives were evaluated by mixed integer linear programming model. The results are suggested to the company for further improvement in its production-distribution network. For future study in this regard, it is recommended to do more attention on sensitivity analysis for different volume of direct shipping. Also, it is recommended to investigate the effect of outsourcing on transportation cost and time. Finally the effect of various transportation modes on productivity of supply networks and distribution networks can be considered in the future studies.

References

1. Persson, F., Olhager, J.: Performance simulation of supply chain designs. *International Journal of Production Economics*, 231–245 (2002)
2. Persson, F.: Supply chain simulation: experience from two case study. In: Verbraeck, A., Hlupic, V. (eds.) *Proceedings of the 15th European Simulation Symposium* (2003)
3. Xia, X., Roland, L., Yu, W.F.: Re-Designing A Distribution Network in A Supply Chain: A Case Study. In: 2009 7th IEEE International Conference on Industrial Informatics, INDIN 2009 (2009)
4. Ding, H., Benyoucef, L., Xie, X.: A Simulation-based Optimization Method for Production-distribution Network Design. In: *IEEE International Conference on Systems, Man and Cybernetics*, pp. 4521–4526 (2004)
5. Chopra, S.: *Designing the distribution network in a supply chain*, pp. 123–140. Elsevier Science Ltd. (2003)
6. Schunk, D., Plott, B.: Using simulation to analyze supply chains. In: *Winter Simulation Conference* (2000)
7. Bruniaux, R., Pierreval, H., Caux, C.: A continuous simulation approach for supply chains in the automotive industry. *Simulation Modelling Practice and Theory*, 185–198 (2007)
8. Banerjee, A., Burton, J., Banerjee, S.: A simulation study of lateral shipments in single supplier, multiple buyers supply chain networks. *International Journal of Production Economics*, 103–114 (2003)
9. Coyle, C.: *The Management of Business Logistics*. South-Western/Thomson Learning, Mason (2003)

10. Hewitt, S.: Comparing analytical and discrete-event simulation models of manufacturing system. Thesis Master, Institute for System Research, University of Maryland (2002)
11. Vieira, G.E.: Ideas for modeling and simulation of supply chains with Arena. In: Proceeding of the 2004 Winter Simulation Conference (2004)
12. Nagarajan, K.V., Vial, P., Awyzio, G.: Simulation of SNMPV3 traffic flow meter mib using Arena simulation modeling software. In: Proceeding of the Modeling and Simulation, Marina Del Rey, USA (2002)
13. Abed, S.Y.: A simulation study to increase the capacity of a rusk production line. *International Journal of Mathematics and Computers in Simulation* 2(3), 228–237 (2008)
14. Sargent, R.: Validation and Verification of Simulation Models. In: Winter Simulation Conference (1999)
15. Harrell, R.: Simulation using ProModel. McGraw-Hill, New York (2003)
16. Benjamin, P.: Toolkit For Enabling Adaptive Modeling And Simulation (Teams). In: Winter Simulation Conference (2002)
17. Burnner, T., Werker, C.: A Practical Guide to Inference in Simulation Models. In: Economics and Evolution, vol. 32. Max Planck Institute of Economics, Germany (2003)
18. Henri Pierreval, C.: A continuous simulation approach for supply chains in the automotive industry. *Simulation Modelling Practice and Theory*
19. Ahmed, K.: Modeling Drivers' Acceleration and Land Changing Behavior. PhD thesis, ITS Program. Massachusetts Institute of technology, Cambridge (1999)
20. Na Li, N., Zhang, M.T., Deng, S., Lee, Z.H., Zhang, L., Zheng, L.: Single-station performance evaluation and improvement in semiconductor manufacturing: A graphical approach. *Int. J. Production Economics* 107, 397–403 (2007)
21. Cuatrecasas, A.L., Santos, F.J., Sanchez, C.V.: The Operations-Time Chart: A graphical tool to evaluate the performance of production systems – From batch-and-queue to lean manufacturing. *Computers & Industrial Engineering* 61, 663–675 (2011)
22. Heizer, J., Render, B.: Production and operations management, 2nd edn., pp. 229–243. Allyn and Bacon, Massachusetts (1998)
23. Robinson, S.: Simulation: The practice of model development and use. John Wiley & Sons (2003)
24. Banks, J.: Hand Book of Simulation. Wiley Inter Science Publications (2007)
25. Kelton, W.D., Sadowski, R.P., Sturrock, D.T.: Simulation with Arena, 4th edn. McGraw Hill, New York (2007)
26. Yang, F.: Neural network metamodeling for cycle time-throughput profiles in manufacturing. *European Journal of Operational Research* 205, 172–185 (2010)

Lead Time Prediction Using Simulation in Leather Shoe Manufacturing

Hermela Solomon, Kassu Jilcha, and Eshetie Berhan

Addis Ababa University, Addis Ababa Institute of Technology, Addis Ababa, Ethiopia
{hermu.y, jkassu}@gmail.com, eshetie_ethio@yahoo.com

Abstract. The purpose of this study is to acquire lead time prediction on the basis of actual data to optimize it. Most Manufacturing industries experiences problems of lead time prediction. The determination of planning values for manufacturing lead times has been viewed both as a problem of estimating independent, uncontrollable variables, and as a control problem, where emphasis is placed on managing the average lead times to match predetermined norms. Once a new order with specific and known processing requirements enters to manufacturing system, an exact lead time estimate is assigned to it that is based on the manufacturing current status. Thus, the customers of the case company was being influenced by unpredictable lead time. As a result of this, the customers were unable to get the product of shoe on time due to indeterminate and fluctuating lead time. To come +up with acceptable lead time, Monte Carlo and Arena simulation tools were used to analyze optimal lead time. The finding of the tools simulation resulted in fixed lead time and order cycle numbers for the products order over specified periods (what is the final result reported by simulation?)

Keywords: Production, Lead time, Leather shoe, Distribution, Monte Carlo.

1 Introduction

Many companies fail to meet customer requirements in their product on time delivery and companies fail to predict their lead time as many studies finding results proved. Therefore, this research paper focuses on lead time prediction for the leather shoe company using simulation techniques to optimize lead time so that customers get satisfaction. To support these problems, it is important to summarize the previous results and conclusion.

More and more companies on the global market are today capable of manufacturing individual or small series orders at comparable prices and quality. The main difference between these companies is the expected production order development time and the observance of delivery deadlines [21].

Since planned lead times represent the amount of time allowed for orders to flow through the production facility, they have a very significant impact on the system performance: tight lead times may lead to past orders and expediting, while loose lead times aggravate the overall material requirement planning (MRP) performance significantly [7].

Lead Time is also regarded an essential element in negotiating with the customer [9]. Lead Time estimation is critical as it exclusively affects customer relations and shop floor management practices. Due date quoting, which means commitment to meeting customer orders on time, is a direct outcome of lead time estimation. Short lead times improve a manufacturer's image and future sales potential. However, not only short but also accurate and precise lead time estimates are desirable.

Leather product manufacturing industries are one of the major industries in Ethiopia. Competitive advantage of the selected case is providing a large variety of products with reasonable price in a short time. In order to be more responsiveness, meet customer needs, reduces response time considering the company's ability.

This study was conducted in one of the oldest leather shoe company which is Anbessa Leather Shoe manufacturer company launched in 1927 E.C and producing Men's , women's and Children's shoe. The company has a capacity to produce 4000 pairs per day and operates 6 days a week and 70% of their products are sold in local market and the rest are exported to different countries.

In this study the problem of indeterminate and fluctuate production lead time was found a big problem in the company. To solve this problem of the company, this study was conducted in selecting the best method of simulation which are Monte Carlo and ARENA Simulation. These tools helped to optimize lead time of Anbessa leather Shoe Company to meet customer orders to make delivery on time and to have fixed production lead time period.

From the late 1960s, the lead time management problem has been consistently addressed in the literature. In an early works by [20] the determination of planning values for manufacturing lead times has been viewed both as a problem of estimating independent, uncontrollable variables, and controllable problems. Monte Carlo Simulation is one of techniques using random variables to estimate possible activity duration in a probability distribution [1].

There are published studies aimed at solving lead time by [18, 24] were modeled three lead time estimation models. Each model requires the estimation of two or more parameters. The parameters in each model were estimated via linear regression from steady state simulation data. These were Jobs in queue (JIQ), Jobs in bottleneck queue (JIBQ) and Combined model (COMB) and the experience obtained during many tests of practical implementation of these procedures, led us to the conclusion that it would be possible to predict lead times of planned orders on the basis of actual operational and assembly order lead times achieved in the past. These predictions (on the basis of enterprise requirement planning (ERP)-system data or on the basis of manually acquired past data) are accurate enough for individual production. Using these data, it is possible to predict lead times even for fairly complex products with several machining operations and individual order features. Based on [21]. There are six procedure for predicting lead times for future production orders.

Lead time estimation are closely related to job due dates. In practice, job due dates are determined either internally by production personnel or externally as a result of negotiation between marketing personnel and the customer. Internally set due dates are often used in make to stock operations, assemble-to-order operations may set due dates internally for the production of components while setting the due dates for finished products externally, while externally set due dates are more prevalent in make to order operations. Whether due dates are set externally or internally, the degree to

which they are met is largely dependent on the ability to accurately predict the time necessary to complete product processing requirements, i.e. lead time [15].

Job and shop characteristics are the two factors affecting flow times. Earlier studies have taken job related parameters (such as total work content and number of operations) into account in rules like CON, TWK, SLK, and NOP [10]. Shop related measures such as total shop load at the time of order arrival or shop congestion on an orders prospective route have been considered in later rules JIQ and JIS are two well-known examples [5]. Using job and shop status data in combination has proven to be more effective [25,4]. Experiments have shown that shop related measures specific to an order are more effective than aggregate shop congestion indicators in estimating flow times [14]. This work makes use of queuing theoretic findings and assumes that utilization rate and process time are factors defining the flow time. The proposed model is quadratic in both factors. Higher order terms also appear in another approach [16]. Studies bases their OFS (operation flow time sampling) and COPS (congestion and operation flow time sampling) rules on the fact that flow times are correlated [23]. A sample of recently completed orders is used to infer ongoing average flow time per operation to reflect near term shop congestion. Their results support the previous findings about superiority of using shop conditions. However; significant differences between the rules are found depending on shop balance. Probability distribution based approaches explicitly consider the impact of variance either of the workload (e.g. JIQ) or of the error in estimation [6] refers to an earlier work making use of the normal distribution in quoting LT estimation subject to uncertainty. Threat flow time as a normal distributed random variable and base LT estimation on a service level requirement [8].

In general, the study focused on lead time prediction and optimized leather shoes company products delivery time using simulation techniques.

2 Material and Methods

The study was conducted considering different materials and methods to achieve the goal of this study. The literature reviews and software applications discussed in this paper were utilized to come up with good optimized lead time prediction solution.

Before starting with the procedure for predicting production order, the right data were collected from the case company. At the beginning of lead time prediction, it is necessary to define the interval for data acquisition of past operational and assembly orders and so that it was done. This interval can be a month, a quarter, a year or several years. As mentioned, a company wishing to predict lead times must have an ERP/PPC system as a basis for all further steps, because this is the database of orders processed in the past. Operational or assembly order codes, type and sequence of operations in manufacturing and assembly orders, actual execution times of operational or assembly orders, date of completing a particular operational or assembly order in the previous workplace, and date of completing a particular operational or assembly order in the observed workplace.

The input data for this predicting method were the actual lead times of operational and assembly orders processed in the past in the case company's workplaces.

The principles of modeling methods such as simulation, mathematical optimization and heuristic models have been described in detail by reference [17]. Some customer service criteria which influenced by delivery time and configuration were: customer experience, response time, return ability, order visibility, product variety and product availability. Hence, simulation has appeared as a powerful method to investigate and design method to predict lead time of the product for the case company.

Different authors finding showed that simulation arena is the best tool to solve such a stochastic system of the companies. This model also customized with Monte Carlo simulation in order to compromise random customer order because today's customer order is not the same with tomorrow order's by simulating those random customer order distribution move towards up with good result.

2.1 Data Analysis

Data was collected from the case company and identified to which tool it fits. Normally Chi Square test is used to identify fitted statistical distribution. Process of leather shoe product including loading and unloading time, transportation time from one station to another station and demand probability distribution with in past two years has been collected using primary and secondary data collection techniques. Chi Square test was done in this study to find fitted distribution for data by using easy fit software. One statistical test that addresses this issue is the chi-square goodness of fit test. This test is commonly used to test association of variables in two-way tables, where the assumed model of independence is evaluated against the observed data.

3 Discussion and Results

3.1 Production Process Flow

In Anbessa shoe factory, there are 3 main sections (Fig 1). Raw materials in the form of leather are issued from store. They will be sent to cutting section and wait for cutting processes. After cutting processes, numbering, skiving the bundle, it will goes

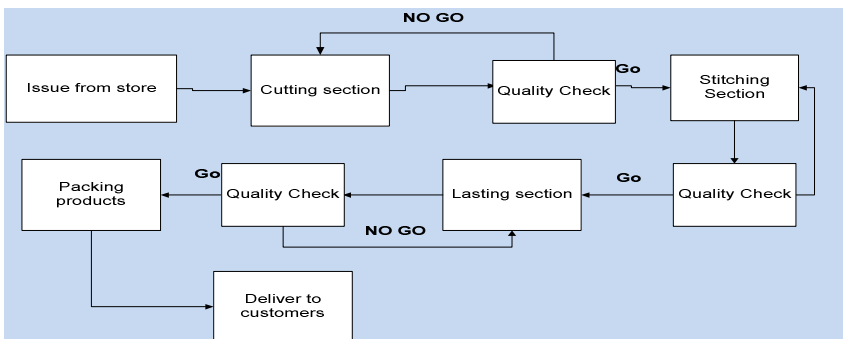


Fig. 1. Product Process flow chart

to stitching section to assemble all components of the leather shoe. After stitching complete luggage products are then going for finishing to lasting section. Finally, inspected and packaging will be done. The product will be kept in finishing storage before delivery to customer.

3.2 Conceptual Model

Simulation modeling is all about building the conceptual model of the case. In fact, modeler should identify and gather all the details and formulations that are necessary for the model in mind map then convert all of them into simulation software. Conceptual model includes inputs .Inputs to Arena models were order per arrival, time between arrivals, assembly time in each section, waiting jobs in queue in each section and time for inspection until deliver to customers.

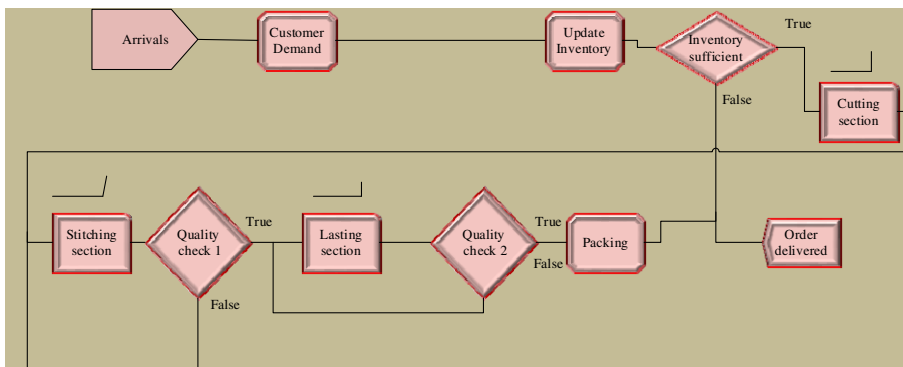


Fig. 2. Conceptual model using Arena Software

3.3 Model Assumption

This study represents continuous event modeling with below assumptions

- The manufacturer working time is eight hours per day
- There is no maintenance performed during the production period.
- Machines are available all the time and no breakdowns
- There is always sufficient raw materials in storage, so the production process never starves.
- All Workers are assumed to produce equal product with in equivalent period of time.
- Transportation time is included in processing time.
- No power interruption

We Are Interested in

- Simulating the system for 10,000 hours
- estimating process utilizations, average job waiting times and average job flow times (the lead time for a job from start to finish)

3.4 Model Formulation

Considering the given definitions and assumptions in the prior section, the schematic model of the system presented in equation (1) helps in calculating delivery time using Monte Carlo tool as:

$$T_i^s = \sum_{j=1}^{j=m_i} t_{ij}; i = 1, 2, \dots, n. \tag{1}$$

Where: t_{ij} = delivery time from station i to station j , T = total delivery time of the whole processes,

A set of uniformly distributed random numbers is needed to generate the arrivals at the checkout counter. Random numbers have the following properties:

1. The set of random time is uniformly distributed between 5 and 15 days.
2. The set of random numbers of orders are uniformly distributed between 500 and 20,000.
3. Successive random numbers are independent.

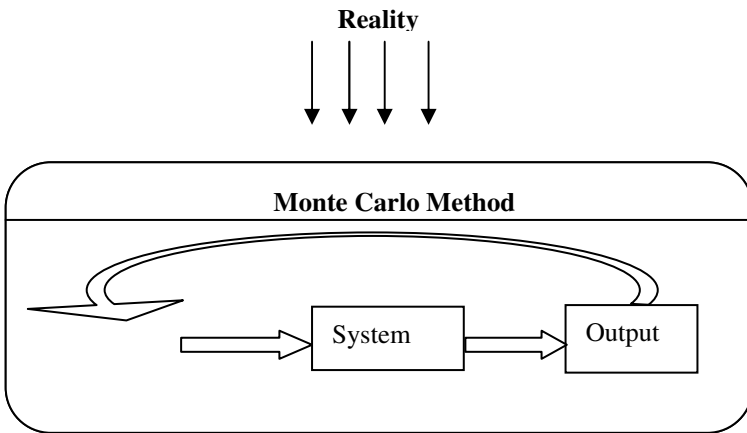


Fig. 3. Monte Carlo process Analysis

Anbesa Shoe Factory demand or order arrival is classified into four intervals to make the study easy in over the two years as shown Table 1. For each order arrival random numbers also generated in line with the probability distribution of orders.

Table 1. Possibility of Order Arrival

Order Arrive	Probability	Random number assigned
0 - 500	0.1	Less than 0.1
600 - 5000	0.15	Between 0.1 and 0.15
5001 - 10000	0.40	Between 0.15 and 0.40
10001 - 20000	0.35	Between 0.35 and 0.40

The key to this simulation is to use a random number to key a lookup from the table range 500 up to 20,000 orders (named *lookup*). Random numbers greater than or equal to 0 and less than 0.10 yielded a demand of 500 and random numbers greater than or equal to 0.35 yielded a demand of 20,000. The study generated 200 random numbers using Monte Carlo simulation tool and would get final optimal result.

Table 2. Final average results with Monte Carlo Tool

Order Arrive	Processing time in all section(hr)	Jobs waiting in the queue(hr)
500	26.2	9.01
5000	246.5	90.1
10000	491	180.2
20000	980	380.4
Sum	1743.7	659.71
	435.925	164.9275

A Monte Carlo based simulation method is designed for running the model to estimate the quantity(order) arrival and the lead time of product in order to help managers to make decisions regarding product lead time and order arrival cycle to customers with alternative scenario that they can see results for many possible variants. This approach offers not just one outcome, but a distribution of possible outcomes.

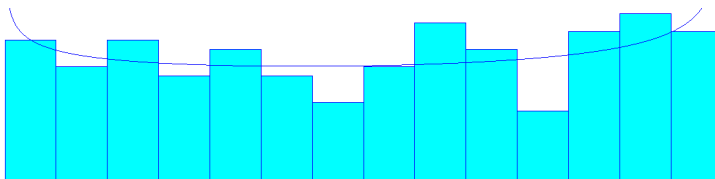


Fig. 4. Beta distribution of demand

Since order arrival distribution was between 500 and 20,000 the study test their distribution using input analyzer and select minimum corresponding p-value which is Beta distribution with the values of BETA (0.882, 0.831)(Fig 4.).

4 Model Validation

Several techniques for validation of model have been offered [13]. In this study, validation of the model was done by comparing the actual data with the results that were generated by simulation software. In addition, for validating of the model, estimating lead time with in company when producing for one day was simulated and it was compared with actual lead time. Table 3 shows the validation data of Monte Carlo and Arena Simulation tools.

Table 3. Validation of Simulation Tools

Order Arrive	Actual Lead Time of the Company Min(hr)	Lead Time using Monte Carlo and Arena(hr)
0 - 500	48	35.21
501 - 5000	352	336.6
5001- 10000	704	671.2
10000 -20000	1056	980

The Monte Carlo and arena software simulation results shown is less than the actual lead time of the company (table 3). This implies that the simulation result is the predicted lead time for the company.

5 Conclusion

Due to ever-fiercer market competition, companies must predict lead times and delivery times with ever greater accuracy to handle its customers and satisfy them. If they give incorrect deadlines, they may not get a request from a particular company next time, which can lead the company into crisis. This article proposes a technique for predicting production order lead times on the basis of actual lead times of past operational or assembly orders. Using the proposed procedure, the company can predict the lead time required for delivery of any new order to any customer and make variations of delivery lead time calculations on the basis of an acceptable risk level by selecting the confidence interval with respect to the size and complexity of an order, and taking into account the company's policy towards its customers.

The study purpose alternatives simulation tools were suggested by the experts to predict lead time of products by using MONTE CARLO and ARENA software. Finally results are suggested to the company for further improvement in its production lead time network within different customers' orders. For example, when the customer orders is (501-5000) pieces, the predicted lead time has been reduced to 336.6 hrs from the current company's lead time of 352hrs. On this basis of the tests, it was found that the procedure for predicting lead times of production orders was well designed and provided very useful data for sales, as well as for production planning and control. Signing a supply contract on the basis of reliable statistical data is completely different from signing a contract on the basis of uncertain, experience-based guesswork.

For future study in this regard, it is recommended to improve lead time by taking into account the sequence of operations required to complete an order, the influence of the number of operations per order, and the influence of the processing time of operations and machine breakdown.

References

1. Yahia, A.A.E.F.: Time Schedule Preparation By Predicting Production Rate Using Simulation
2. Arena Professional Reference Guide. Rockwell Software Inc. (2000)
3. Benkó, J.: Modeling Kanban Systemic Production with Arena Simulator. *GépGyártás XLIX(4)* (2009)
4. Bertrand, J.W.M.: The use of workload information to control job lateness in controlled and uncontrolled release production systems. *Journal of Operations Management* 3(2), 79–92 (1983)
5. Eilon, S., Chowdhury, I.G.: Due dates in job shop scheduling. *International Journal of Production Research* 14(2), 223–237 (1976)
6. Enns, S.T.: A dynamic forecasting model for job shop flow time prediction and tardiness control. *International Journal of Production Research* 33(5), 1295–1312 (1995)
7. Ho, J., Chang, Y.: An integrated MRP and JIT framework. *Computers and Industrial Engineering* 41, 173–185 (2001)
8. Hopp, W.J., Sturgis, M.L.R.: Quoting manufacturing due dates subject to a service level constraint. *IIE Transactions* 32(9), 771–784 (2000)
9. Moodie, D.: Demand management: The evaluation of price and due date negotiation strategies using simulation. *Journal of Production Operations Management Society* 8(2), 151–162 (1999)
10. Nyhuis, P., Wiendahl, H.P.: *Logistische Kennlinien*, pp. 81–94. Springer, Heidelberg (1999)
11. Percentile (2006), <http://en.wikipedia.org/wiki/Percentile>
12. Harrell, R.: *Simulation using ProModel*. McGraw-Hill, New York (2003)
13. Sargent, R.: *Validation and Verification of Simulation Models*. In: *Winter Simulation Conference* (1999)
14. Ragatz, G., Mabert, V.: A framework for the study of due date management in job shops. *International Journal of Production Research* 22(4), 685–695 (1984)
15. Ruben, R.A., Mahmoodi, F.: Lead time prediction in unbalanced production systems (November 14, 2010)
16. Ruben, R.A., Mahmoodi, F.: Lead time prediction in unbalanced production systems. *International Journal of Production Systems* 38(7), 1711–1729 (2000)
17. Chopra, S.: *Designing the distribution network in a supply chain*, pp. 123–140. Elsevier Science Ltd. (2003)
18. Smith, M.L., Siedman, A.: Due date selection procedures for job-shop simulation. *Computers and Industrial Engineering* 7(3), 199–207 (1983)
19. Smith, M.L., Siedman, A.: Due date selection procedures for job-shop simulation. *Computers and Industrial Engineering* 7(3), 199–207 (1983)
20. Tatsiopoulos, I.P., Kingsman, B.G.: Lead time management. *European Journal of Operational Research* 14(4), 351–358 (1983)
21. Berlec, T., Starbek, M.: Predicting Order Due Date, doi:10.1007/s13369-012-02791

22. Veral, E.: Computer simulation of due-date setting in multi-machine job shops. *Computers and Industrial Engineering* 41(1), 77–94 (2001)
23. Vig, M.M., Dooley, J.K.: Dynamic rules for due-date assignment. *International Journal of Production Research* 29(7), 1361–1377 (1991)
24. Vig, M.M., Dooley, K.J.: Mixing static and dynamic flow time estimates for due date assignment. *Journal of Operations Management* 11, 67–79 (1993)
25. Weeks, J.K.: A simulation study of predictable due-dates. *Management Science* 25(4), 363–373 (1979)
26. Wiendahl, H.P.: *Load-oriented Manufacturing Control*, pp. 37–199. Springer, Berlin (1995)

Ensemble Neurocomputing Based Oil Price Prediction

Lubna A. Gabralla¹, Hela Mahersia², and Ajith Abraham^{1,3,4}

¹ Faculty of Computer Science and Information Technology,
Sudan University of Science & Technology, Khartoum, Sudan
lubnagabralla@gmail.com

² Image Signal and Information Processing Research Laboratory,
National Engineering School of Tunis (ENIT), University of Tunis El Manar, Tunisia
helamahersia@yahoo.fr

³ IT4Innovations - Center of excellence, VSB - Technical University of Ostrava,
Czech Republic

⁴ Machine Intelligence Research Labs,
Scientific Network for Innovation and Research Excellence, WA, USA
ajith.abraham@ieee.org

Abstract. In this paper, we investigated an ensemble neural network for the prediction of oil prices. Daily data from 1999 to 2012 were used to predict the West Taxes, Intermediate. Data were separated into four phases of training and testing using different percentages and obtained seven sub-datasets after implementing different attribute selection algorithms. We used three types of neural networks: Feed forward, Recurrent and Radial Basis Function networks. Finally a good ensemble neural network model is formulated by the weighted average method. Empirical results illustrated that the ensemble neural network outperformed other models.

Keywords: Oil price prediction, ensemble neural network, computational intelligence.

1 Introduction

Oil is one of the most important topics in the contemporary world and it will remain as a keyword in the world politically, and economically. Oil has unique properties that can be used to control and conquer the world successfully. The history of oil discovery goes back to the 1859, when the first oil was drilled in Pennsylvania, United States [1]. After that it became very useful in the manufacturing engines and cars, planes and machinery. In 1914 during the first World War the head of the French government at the time described that every drop of oil is equal to a drop of blood in an orientation to its importance [2].

At present, oil is the most important source of energy and one of the elements of modern civilization for humans. It is used as a fuel for cars, airplanes, factories and agricultural equipment, trucks, commercial and military ships and electric power generation for homes, workplaces and other places. Oil prices have undergone many changes and instabilities over the years. It was known as oil shocks, and the first

shock was in the October war in 1973, where the price rose from 2.29\$ to 10.73\$ for a barrel until 1974, and these prices continued to rise, and even achieved strong jump to 32.51\$ per barrel in 1981, this what is known as the second oil shock, and the third oil shock was after Iraq's invasion of Kuwait, when the price of oil rose from 17.31\$ for the year 1989 to 22.26\$, in 1990. Oil prices continued volatility until prices of oil was collapsed in 1998 and the average barrel of oil was around 9.69\$ for OPEC. This was as a result of decline in global oil demand after the financial crisis [3]. Several researchers and scientists were interested in studying the factors that influence the oil prices, like climate [4], politics [5], and stock market [6] etc. Owners of the economic sector, such as commercial institutions and companies operating in the field of oil were very much concerned and wanted to know the prices of oil in the coming years in order to determine their economic policies and building plans for the future and make informed decisions, which will help them to avoid the problems of inflation and economic stagnation, losses and financial crises. Recently several artificial intelligence algorithms were used for oil price prediction. Artificial neural networks have many characteristics and does not need any hypotheses (a priori) to be introduced and is able to deal with incomplete information and with the large number of variables and generally it is flexible in modeling [7]. Therefore, this study aims to employ several types of neural network algorithms to develop a computational model that is able to predict oil prices with high accuracy and high performance, which can contribute to the development of the local and global economy. The rest of this paper is organized as follows: After a short literature review in Section 2, Section 3 describes the research methodology in detail. The data used and their divisions are found in Section 4, and experimental results are reported in Section 5 followed by concluding remarks.

2 Related Works

Artificial neural networks (ANN) [8] are designed to represent data by simulating the work of the human brain. ANN's emerged in different areas such as industrial, medical and business, and achieved successful results therefore many researchers also used ANN in the oil industry. Kaboudan [9] selected multilayer perceptron (MLP) and Genetic programming (GP) to forecast crude oil price using monthly data, such as world crude production, OECD consumption, world stocks and lagged crude FOB crude oil price of US imports. Two methods are compared to a random walk and their results proved that GA has an advantage over random walk predictions. Yu et al. [10] constructed an empirical mode decomposition (EMD) based on neural network ensemble learning. They used daily West Texas Intermediate (WTI) data from 1/1/1986 to 30/9/2006 as training and Brent from 20/5/1987 to 30/9/2006 as testing data. Results proved that EMD based neural network ensemble can be used for oil price prediction. Haidar et al. [11] suggested a network to predict the oil prices using two groups of inputs, crude oil futures data, and Dollar index, S&P500, gold price and heating oil price. The authors measured performance by hit rate, root mean square error, correlation coefficient, mean squared error and mean absolute error. The authors concluded that heating oil spot price support forecast crude oil spot price in numerous steps prediction. Alizadeh and Mafinezhad [12] proposed General

Regression Neural network (GRNN) using six factors monthly data to predicting Brent crude oil price. Experiment results show that the model achieved high accuracy in normal and crisis situations. Mingming and Jinliang [13] collected data covering Brent and West Texas Intermediate (WTI) from 1946 to 2010 and adopted multiple wavelet recurrent neural networks (MWRNNs) to forecast crude oil prices. The study showed that the model has high prediction accuracy. Yu et al. [14] introduced a fuzzy ensemble prediction model, support vector machine, radial basis function networks and back-propagation neural networks to predict crude oil prices. They used the data covering a period from January 2000 to December 2007 using West Texas Intermediate and Brent crude oil spot. Results showed that the agent-based fuzzy ensemble prediction model outperformed other individual methods in accuracy. Most of the studies in the literature focused on constructing a new model using one percentage of training and testing on the other hand, few researchers were interested in using different inputs for testing. So the objective of this paper was to provide a variety of the training and testing percentages with a set of different inputs using several kinds of neural networks to get high accuracy for the model.

3 Research Methodology

3.1 Feed Forward Neural Networks (FFN)

Back propagation [15][16] method is a supervised learning scheme and the most popular technique in multilayer networks when a set of input produces its own actual output and then compare it with the target value by calculating the error, after that error is fed back through the network. The weights of each connection are adjusted to reduce the error by several ways, such as gradient descent etc. until sufficient performance is achieved. To improve the generalization, there are several learning methods such as Levenberg – Marquardt (LM), Bayesian regularization (BR) and BFGS quasi-Newton (BFG-QN) back propagation algorithm [17].

3.2 Recurrent Neural Network (RCN)

RCN is the state of the art in nonlinear time series prediction, system identification, and temporal pattern classification. As the output of the network at time t is used along with a new input to compute the output of the network at time $t + 1$, the response of the network is dynamic [8].

3.3 Radial Basis Function (RBF)

Radial basis function network [8] is an artificial neural network that uses radial basis functions as activation functions. The output of the network is a linear combination of radial basis functions of the inputs and neuron parameters. RBF is successful in numerous fields especially for system control, time series and prediction.

3.4 Ensemble Neural Network [18]

The generalized ensemble method find weights for each output that minimizes the MAE of the ensemble. The general ensemble model (GEM) is defined by:

$$F_{GEM} = \sum_{i=1}^n \alpha F_i(X) \quad (1)$$

Where $\alpha F_i(x)$ are chosen to minimize the MAE between the outputs and the desired values. Figure 1 illustrates the ensemble neural network approach.

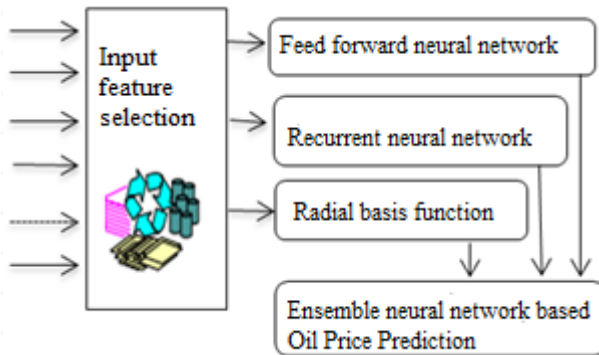


Fig. 1. Ensemble neural network for oil price prediction

4 Experiments

The daily data [19, 20] (from 1999 to 2012) were used to predict the West Taxes Intermediate (Output). The dataset consists of 14 variables as following:

- Date (DT).
- West Texas Intermediate (WTI).
- Federal Fund rate (FFR).
- Volatility Implied Equity Index (VIX).
- The regional Standard & Poor's equity index US, Europe and Asia (SPX).
- Gasoline prices New York & US Gulf Coast (GPNY) & (GPUS).
- Heating oil spot prices (HP).
- Future contracts 1,2,3,4 for WTI (FC1, FC2, FC3, FC4).
- EUR/USD exchange rate (ER).
- Gold prices (GP).

We also examined the effect of training and testing data by randomly splitting them as follows:

- 90% - 10% (A)
- 80% - 20% (B)
- 70% - 30% (C)
- 60% - 40% (D)

We used WEKA for pre-processing experiments [21] and formulated 7 different sub datasets, which were derived from the original dataset after implementing several attribute selection algorithms, such as:

- Attributes ranking principal: ranked list of attributes based on evaluated individually each attribute [22].
- Wrapper attributes Selection: It depends on an induction algorithm to estimate the merit of feature subsets [23].
- Relief for regression: Evaluates quality of attributes according to value of the given attribute for the near instance to each other and different predicted (class) value [24].
- Correlation based Feature Selection (CFS): assesses the value of group of attributes by concerning the individual predictive ability of each features as well with the possibility of repetition among the features. Selecting a subset of the original attributes to reduce the dimensionality of the data and then constructing a model from these reduced number of features in some cases could improve the prediction accuracy and performance, and a simpler model that is easier to interpret [19][22]. Table 1 summarizes the results of attribute selection and the 7 sub-data sets (SDS1-SDS7) obtained.

Table 1. Attributes selection methods and their features

Sub dataset	Method	Features
SDS1	Correlation based Feature Selection subset evaluator	WTI; SPX; FG1
SDS2	Correlation based Feature Selection subset evaluator	DT; VIX; WTI; SPX; GPNY; GPUS; HP; ER; FC1; FC2; FC3; FC4
SDS3	Correlation based Feature Selection subset evaluator	VIX; WTI; GPNY; ER; FC1
SDS4	Correlation based Feature Selection subset evaluator	WTI; GPNY; FC1
SDS5	Correlation based Feature Selection subset evaluator	VIX; WTI; GPNY; FC1
SDS6	Wrapper subset evaluator	WTI; FC1
SDS7	Wrapper subset evaluator	WTI; GPUS

5 Experimental Results

5.1 Feed Forward Neural Network

Neural network experiments are accomplished in MATLAB. We used one hidden layer exploring 40-45-50-55-60 neurons and used tan-sigmoidal transfer function for

the hidden layer and pure linear function in the output layer. We measured the performance using mean absolute error (MAE).

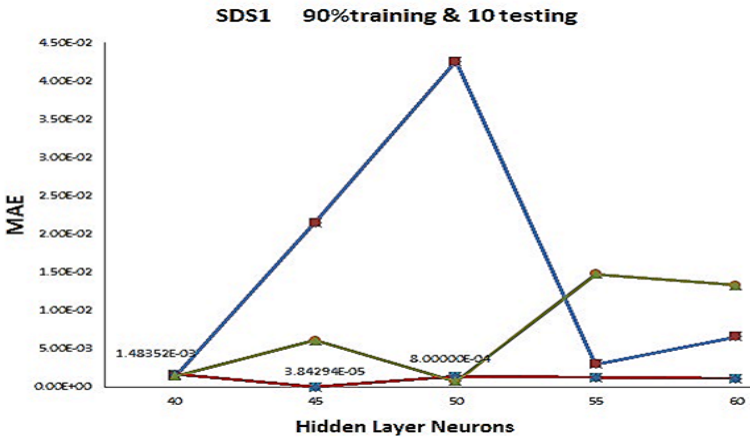


Fig. 2. Comparison of feed-forward networks using three different training algorithms

According to Table 2, in most of the sub-datasets the best results were obtained when using the Bayesian regularization back propagation method with 80% training and 20% testing, and sub-dataset1 (SDS1) achieved MAE= 3.843E-05 with 90% training and 10% testing using 45 neurons. Figure 2 shows the best results in SDS1 using feed-forward networks comparing Levenberg –Marquardt (LM), Bayesian regularization (BR) and BFGS Quasi-Newton (BFG-QN) algorithms.

Recurrent Neural Network

We used a hidden layer with 10 neurons and used three training algorithms Levenberg –Marquardt (LM), Bayesian regularization (BR) and BFGS Quasi-Newton (BFG-QN). Bayesian regularization method outperformed other algorithms by 51.85%. It is noted from Table 3 that for all the sub-datasets in the percentage 80% training and 20% testing is the best (shaded area), except in sub-dataset (SDS5) 90% training and 10% testing is the best. On the other hand the lowest value of the error is 3.941 E-05 when using 90% training and 10% testing with sub-dataset (SDS5).

Radial Basis Function Network

We constructed the network until it reached a maximum number of neurons or the sum-squared error falls beneath an error goal. Table 4 shows the results obtained using the seven sub-datasets and different number of neurons: 40, 45, 50, 55 and 60. The shaded area indicates the best results when using 60 neurons in few sub-datasets then followed by 55 neurons. According to the percentage of training and testing sub-dataset (SDS4 - SDS5-SDS6 - SDS7) achieved the best results with 80% training &20% testing. The best results over all sub-datasets is an MAE of 2.206 E-05 in SDS6 with 80% training & 20 % testing using 45 neurons.

Table 2. Performance of FFN

Sub-datasets	Data	Mean Absolute Error			Hidden layer neurons
		LM	BR	BFG-QN	
SDS1	A	1.48352E-03	3.84294E-05	8.00000E-04	45
	B	4.81400E-03	1.35000E-04	4.00000E-04	45
	C	5.77900E-03	3.55000E-04	1.50000E-03	45
	D	8.15200E-03	4.38000E-04	6.30000E-03	40
SDS2	A	9.17000E-04	2.83700E-03	9.00000E-04	40
	B	4.48000E-04	2.18100E-03	4.00000E-04	40
	C	2.74700E-03	1.02200E-03	1.90000E-03	45
	D	3.15700E-03	7.69000E-04	1.00000E-03	60
SDS3	A	1.83600E-03	1.26594E-04	1.10000E-03	50
	B	3.00400E-03	1.24897E-04	9.00000E-04	50
	C	1.21500E-02	5.31100E-03	3.10000E-03	45
	D	1.69940E-02	4.06200E-03	1.80000E-03	45
SDS4	A	5.58850E-02	5.74155E-05	7.80000E-03	50
	B	2.86700E-03	6.45000E-05	1.60000E-03	50
	C	2.48740E-02	3.04484E-04	7.70000E-03	50
	D	1.67079E-02	1.94000E-04	2.40000E-03	50
SDS5	A	3.50300E-03	9.65000E-04	2.00000E-03	55
	B	1.40900E-03	6.80000E-05	8.00000E-04	60
	C	2.19900E-02	3.73327E-04	2.80000E-03	40
	D	4.28700E-02	2.10900E-03	2.60000E-03	40
SDS6	A	1.44560E-02	6.23000E-05	5.70000E-03	40
	B	1.94854E-02	6.04640E-05	3.10000E-03	60
	C	3.23723E-01	3.49000E-04	4.95000E-02	55
	D	1.07220E-01	1.62000E-04	2.25000E-02	55
SDS7	A	5.69780E-02	1.90200E-03	1.92300E-01	40
	B	1.39500E-02	1.97000E-04	9.70000E-03	55
	C	9.65800E-02	2.25700E-03	3.45000E-02	40
	D	2.83030E-01	4.50000E-04	3.32000E-02	55

Experiments Using Ensemble Method

We compared the results of three different types of neural networks and observed that the RBF network outperformed other methods in obtaining the lowest error (MAE= 2.206 E-05). Also the data set using training 80% and testing 20% accomplished the best results in all the neural network methods. In the feed-forward and radial basis networks, the best results were obtained when using 45 neurons. RBF networks outperformed again in the time factor, as it was faster than feed-forward and recurrent neural network.

To further improve the results, we used the ensemble neural network. We selected the best results using the category of 80% training and 20% testing for three neural networks. Experimental results illustrate that the proposed ensemble neural network is superior to other methods by achieving the lowest MAE = 2.186E-05 as shown in Table 5.

Table 3. Performance of RCN

Sub-datasets	Data	Mean Absolute Error		
		LM	BR	BFG-QN
SDS1	A	1.14102E-03	1.27500E-03	2.52400E-03
	B	1.17400E-03	2.48800E-03	5.79000E-04
	C	1.03650E-02	6.98300E-03	1.25780E-02
	D	1.29940E-02	7.69800E-03	1.29150E-02
SDS2	A	4.60000E-04	3.77000E-04	2.18328E-02
	B	2.22000E-04	1.74000E-04	1.18390E-02
	C	1.36700E-03	1.44800E-03	7.03600E-02
	D	8.61000E-04	3.32000E-04	3.45920E-02
SDS3	A	5.22000E-04	4.43900E-03	2.55500E-03
	B	3.57762E-04	1.05000E-04	6.40100E-03
	C	4.21100E-03	2.82000E-04	4.10660E-02
	D	6.82000E-04	1.68000E-04	1.67200E-02
SDS4	A	7.68000E-04	4.80285E-05	1.58640E-02
	B	5.20102E-05	3.94799E-05	7.16800E-03
	C	4.73000E-04	2.17658E-04	2.00530E-02
	D	2.08400E-03	1.81824E-04	6.17200E-03
SDS5	A	3.94114E-05	9.52860E-05	4.38700E-03
	B	1.16000E-04	1.12000E-04	4.13600E-03
	C	4.41680E-04	3.77708E-01	1.42909E-01
	D	4.95000E-04	3.62000E-04	9.70200E-03
SDS6	A	1.83283E-04	1.51900E-03	3.33800E-03
	B	1.72000E-04	1.89800E-03	1.58400E-03
	C	1.82716E-03	6.89000E-03	6.86900E-03
	D	9.69000E-04	5.59500E-03	4.03800E-03
SDS7	A	5.80000E-04	2.48000E-04	1.01690E-02
	B	1.47000E-04	2.19000E-04	2.20500E-03
	C	1.30400E-03	1.01200E-03	4.03600E-03
	D	4.30824E-04	6.10500E-03	1.03750E-02

Table 4. Performance of Ensemble neural network

Prediction models	MAE
Feed forward	6.0464E-05
Recurrent	3.9479E-05
Radial Basis Function	2.2191E-05
Ensemble	2.1862E-05

Table 5. Performance of RBF

Sub-datasets	Data	Mean Absolute Error based on the number of neurons				
		40	45	50	55	60
SDS1	A	1.03146E-04	9.34926E-05	1.00340E-04	5.70070E-05	5.08729E-05
	B	1.53000E-04	1.48000E-04	1.09000E-04	2.58000E-04	2.40010E-05
	C	4.84160E-05	5.01626E-05	5.02828E-05	5.02436E-05	4.98570E-05
	D	3.88153E-05	3.81595E-05	3.81548E-05	3.81543E-05	3.81548E-05
SDS2	A	5.05680E-03	2.70885E-03	1.53678E-03	1.33255E-03	1.36009E-03
	B	4.09600E-03	3.23000E-03	3.02800E-03	1.85800E-03	1.68000E-03
	C	6.13000E-03	6.04800E-03	5.16800E-03	4.47700E-03	2.85700E-03
	D	8.42200E-03	8.58200E-03	7.28400E-03	5.07200E-03	4.30200E-03
SDS3	A	2.97000E-04	2.93000E-04	2.91000E-04	2.86000E-04	2.75000E-04
	B	7.31000E-04	7.51000E-04	7.21000E-04	3.81000E-04	2.24000E-04
	C	1.50700E-03	6.52000E-04	5.11000E-04	9.74410E-04	6.96864E-04
	D	1.60300E-03	1.56300E-03	1.51400E-03	1.52400E-03	1.54100E-03
SDS4	A	2.11138E-04	2.09055E-04	2.09022E-04	2.09002E-04	2.08946E-04
	B	6.08910E-05	6.23224E-05	6.23330E-05	6.23224E-05	6.35603E-05
	C	1.25481E-04	1.10437E-04	1.13763E-04	1.16843E-04	1.16840E-04
	D	4.06000E-04	3.84000E-04	3.81000E-04	3.80000E-04	3.80000E-04
SDS5	A	1.00800E-03	8.17000E-04	3.54000E-04	1.60000E-04	4.34000E-04
	B	1.53000E-04	1.03000E-04	1.26000E-04	1.04000E-04	1.21483E-04
	C	4.14000E-04	4.01671E-04	4.01344E-04	4.28000E-04	2.34880E-04
	D	2.54000E-04	2.16000E-04	1.97000E-04	1.85000E-04	8.74530E-05
SDS6	A	9.04430E-02	9.04440E-02	9.04440E-02	9.04440E-02	9.04440E-02
	B	2.21914E-05	2.20646E-05	2.21172E-05	2.21087E-05	2.21087E-05
	C	4.51944E-05	4.51615E-05	4.49053E-05	4.48180E-05	4.47060E-05
	D	3.51792E-05	3.41134E-05	3.40980E-05	3.35330E-05	3.35330E-05
SDS7	A	2.34000E-04	2.85000E-04	2.76000E-04	2.72000E-04	2.76000E-04
	B	4.03840E-05	4.03361E-05	3.98230E-05	3.98158E-05	3.98230E-05
	C	5.37473E-05	5.34101E-05	5.34009E-05	5.36361E-05	5.36107E-05
	D	1.21000E-04	1.20000E-04	1.20000E-04	1.16000E-04	1.16000E-04

6 Conclusions

In this paper, we presented an ensemble neural network model for prediction of oil price. The model is based on three different types of neural networks: feed-forward, recurrent and radial basis function networks. The network structure was selected after many experiments including a number of hidden neurons and several learning methods. In addition, four different groups of training and testing were experimented and many attribute selection algorithms were implemented, which led to 7 different sub-datasets. The results illustrate that the radial basis function achieved the best MAE and less time to run when compared to other individual methods. Ensemble methods were found to be superior when compared to the individual neural networks and learning methods.

References

1. <http://history1800s.about.com/od/oil/a/first-oil-well.htm>
2. http://www.theforbiddenknowledge.com/..the_rothschild_blood_line.htm
3. http://www.opec.org/opec_web/en/publications/457.htm
4. Pindyck, R.S.: The dynamics of commodity spot and futures markets: a primer. *The Energy Journal*, 1–29 (2001)
5. Griffin, J.M.: OPEC behavior: a test of alternative hypotheses. *The American Economic Review*, 954–963 (1985)
6. Soytaş, U., et al.: World oil prices, precious metal prices and macroeconomy in Turkey. *Energy Policy* 37(12), 5557–5566 (2009)
7. Maimon, O.Z., Rokach, L.: *Data mining and knowledge discovery handbook*, vol. 1. Springer (2005)
8. Abraham, A.: Artificial neural networks. In: *Handbook of Measuring System Design* (2005)
9. Kaboudan, M.: Compumetric forecasting of crude oil prices. In: *Proceedings of the 2001 Congress on Evolutionary Computation*. IEEE (2001)
10. Yu, L., Wang, S., Lai, K.K.: Forecasting crude oil price with an EMD-based neural network ensemble learning paradigm. *Energy Economics* 30(5), 2623–2635 (2008)
11. Haidar, I., Kulkarni, S., Pan, H.: Forecasting model for crude oil prices based on artificial neural networks. In: *International Conference on Intelligent Sensors, Sensor Networks and Information Processing*, ISSNIP 2008. IEEE (2008)
12. Alizadeh, A., Mafinezhad, K.: Monthly Brent oil price forecasting using artificial neural networks and a crisis index. In: *2010 International Conference on Electronics and Information Engineering (ICEIE)*. IEEE (2010)
13. Mingming, T., Jinliang, Z.: A multiple adaptive wavelet recurrent neural network model to analyze crude oil prices. *Journal of Economics and Business* 64(4), 275–286 (2012)
14. Yu, L., Wang, S., Lai, K.K.: A generalized intelligent-agent-based fuzzy group forecasting model for oil price prediction. In: *IEEE International Conference on Systems, Man and Cybernetics, SMC 2008, IEEE* (2008)
15. Reed, R.D., Marks, R.J.: *Neural smithing: supervised learning in feedforward artificial neural networks*. MIT Press (1998)
16. Chauvin, Y., Rumelhart, D.E.: *Backpropagation: theory, architectures, and applications*. Psychology Press (1995)
17. Demuth, H., Beale, M., Hagan, M.: *Neural network toolbox™ 6. User's guide* (2008)
18. Maqsood, I., Khan, M.R., Abraham, A.: An ensemble of neural networks for weather forecasting. *Neural Computing & Applications* 13(2), 112–122 (2004)
19. <http://www.eia.gov>
20. <http://www.usgold.com>
21. Garner, S.R.: Weka: The waikato environment for knowledge analysis. In: *Proceedings of the New Zealand Computer Science Research Students Conference*. Citeseer (1995)
22. Guyon, I., Elisseeff, A.: An introduction to variable and feature selection. *The Journal of Machine Learning Research* 3, 1157–1182 (2003)
23. Hall, M.A.: Correlation-based feature selection for machine learning. *The University of Waikato* (1999)
24. Robnik-Šikonja, M., Kononenko, I.: An adaptation of Relief for attribute estimation in regression. In: *Machine Learning: Proceedings of the Fourteenth International Conference, ICML1997* (1997)

Internet of Things Communication Reference Model and Traffic Engineer System (TES)

Adel H. Alhamedi, Hamoud M. Aldosari, Václav Snášel, and Ajith Abraham

VŠB-Technical University of Ostrava

17. listopadu 15/2172, 708 33 Ostrava - Poruba, Czech Republic

{ade0004,mub0002,vaclav.snasel}@vsb.cz, ajith.abraham@ieee.org

Abstract. One of the biggest challenges facing *Internet of Things (IoT)* is the existing infrastructure of Internet and its mechanism of action. This paper proposes a new system, which sends the *full* Internet best path (between source and destination objects) to source object on IoT. This will help data of source object to reach its final destination object faster. This system saves most of recalculation of the Internet best paths again and again in the Internet Routers during a data trip. The authors call this system *Traffic Engineer System (TES)*. The most important effect of this system is that it changes the form of "*Internet of Things Communication Reference Model*". This paper merges two addressing layers (IP/ID and Link) from this model in one new layer; where routers transition data through one address and the data have its full best path.

Keywords: Internet of Things (IoT), IoT Communication Reference Model, Traffic Engineer System (TES).

1 Introduction

Internet of the future is likely to be dramatically different from the Internet we use today. This development is opening up huge opportunities for each of the scientific research and the economy. However, it also involves risks and undoubtedly represents an immense technical and scientifically challenge [1]. The Internet Protocol (IP) is suited for networking devices with stringent requirements [2]. According to Internet of the future concepts, *data will be self-addressable and self-routable* [3]. In the world of Internet of Things (IoT), can we use IP as the data exchange protocol directly without any modifications with urgent need for more efficiency? Especially with this tremendous progress in communication techniques and the power of hardware. The ability to uniquely identify things (objects) is critical for the success of IoT. This will not only allow us to uniquely identify billions of devices, but also to control remote devices through the Internet [4]. On the other hand in IoT, routing packets and inter-nodal communication have received little attention; mainly due to the sheer reliance on the today's Internet as it is and as a backbone [5]. IBM's newest study reveals how new technologies support the development of the Internet of Things, and how the Internet of Things provides the foundational infrastructure for a smarter planet [6].

1.1 Previous Work

This paper represents a supplement to authors’ previous paper [7] related to the communication between objects in Internet of Things (IoT). Figure 1 illustrates the output model "IoT Communication Reference Model" in that paper. The model is built from 7 layers, bottom to top, in the following order: Physical, Quality of Service, Security, Link, IP/ID, End-to-End, and Data; where two new layers have been added to the original model (Security and Quality of Service). *This paper merges other two layers from the model in one new layer!*

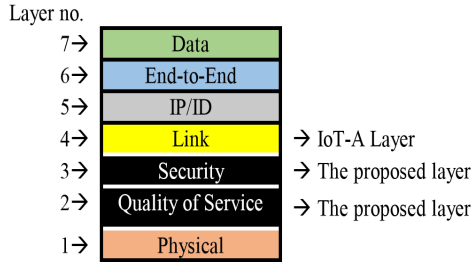


Fig. 1. IoT Communication Reference Model

1.2 Addressing in the Model

Why does today’s Internet require two types of addresses (like IPv4 and MAC) to achieve the process of communication between source and destination objects? And then transfer data between these objects through the Internet communication devices (Routers). *In the beginning*, we all agree that the issue of "Addressing" as it is now is a successful design that has been implemented and has achieved the communication goals. However, that does not mean this is the only way! There may be other more efficient ways, especially since we are on the verge of a big changes in everything through IoT. This leads us to think about everything related to today’s Internet. *Secondly*, to answer this question we need to know a detailed and profound answer to the following question: what is the job and benefits of IPv4 and MAC Addresses? Also, because this paper wants to develop "IoT Communication Reference Model", it has to find the relation between these addresses, model’s layers, and Internet devices. *Finally*, note that authors use IPv4 and MAC for simplicity the idea and as a guide for what they need to prove (although authors know that the Internet of Things is more complicated than that). For example, nowadays we are going to transition to a new generation of IP addresses (IPv6) [8]. As well as with respect to the enormous diversity in addressing when adding various and highly heterogeneous objects to the Internet to build IoT. At the end of this paper it will become clear how the authors implement this idea with IPv6 and new parallel addresses to MAC.

Table 1 shows a comparison between MAC and IPv4 Addresses [9, 10, 11, and 12] and contains the answer to all previous questions. IPv4 is an address for location and carries in data packet to help routing table in Routers to route the data to its final destination; Communication Address from End-to-End. MAC is a name for hardware

and carries in data frame to help CAM table in Switches to forward the data to the next neighbor on the track; Connection Address for Hop-by-Hop.

So we cannot make Internet work without one of these addresses, because each address has its characteristics, advantages, specific role, and final goal. *However, it is possible to change its mechanism of action!*

Table 1. Comparison between IPv4 address and MAC address

VS	Addressing	IPv4	MAC
1.	Name	Internet Protocol	Media Access Control
2.	Addressing	It's address for location	It is more comparable to a name than an address
3.	Assigned by	Network administrator or Internet Service Provider	During manufacturing the Network card
4.	Nature	Logical Address	Physical Address
5.	Static/Dynamic	Static/Dynamic and could be changed	Static and Permanent and could not be changed
6.	Unique	Unique; relation to the time or provider	Unique; relation to the hardware
7.	Location	Intranet/Internet - Per Card	LAN/Link - Per Card
8.	Size	32 bit (4 Byte)	48 bit (6 Byte)
9.	Notation	Dotted Binary/Decimal	Columned Hexadecimal
10.	Formula	Network ID + Host ID (As needed)	Manufactory ID + Card ID (fifty-fifty)
11.	Types	Class A,B, and C - Private and Public	Accordingly manufacturers classification
12.	Datagram	Packet carry First Source's IP and Final Destination's IP	Frame carry neighborly relationship Source MAC and Destination MAC
13.	Goals	Communication Address from End-to-End somewhat resembles Passport number in our life	Connection Address for Hop-by-Hop somewhat resembles the local ID in our life.
14.	Using for	Routing Packets by Routing Tables	Forwarding Frames by CAM Tables (MAC Tables)
15.	IoT CRM Layer	5 th IP/ID Layer	4 th Link Layer
16.	Devices	Layer 3 Devices; like Router	Layer 2 Devices; like Switch
17.	Tables	Routing Table support around 370,000 Network	CAM or MAC Table support around 8000 MAC
18.	General broadcast	255.255.255.255	ff:ff:ff:ff:ff:ff
19.	Translation	DNS: URL or Name → IP Address	ARP: IP Address → MAC Address
20.	Source	IANA	IEEE

1.3 Imaginary Simulation Scenario

An imaginary simulated scenario to simplify the problem is introduced: Suppose that there are two persons on a trip in a country for the first time. The first person has a map for this country while the second has no map. Therefore, if we assume that they will go to the same destination, who will be the faster? Definitely the first one, because the second will lose some time by asking people about the destination. The current situation today when data transition in Internet to its destination, it represents the other person: asking each Router during the trip about the best path to the destination. Note: Persons represent the Data, Destination represents destination object IP address, Map represents the Best path, and People represent Routers.

The remaining part of this paper is organized as follows: Section 2: discusses the issue of best path to the final destination with latest mechanisms and its problems from our point of view. Section 3: our proposed system "Traffic Engineer System". Section 4: results and discussions. Section 5: conclusions and future work. Section 6: acknowledgments.

2 Best Path for Final Destination

With the rapid growth of the Internet and the establishment of IP as the Layer 3 protocol of choice in most environments, the drawbacks of traditional IP routing

became more and more obvious [13]. One of these drawbacks is re-calculation of the Internet best paths to final destination again and again in the Internet Routers during data trip. In recent years, there have been several attempts to reduce the repetition of this process, improve its performance, and accelerate it. The best and most powerful of those attempts is Cisco Express Forwarding (CEF) from Cisco, Multiprotocol Label Switching (MPLS), and Traffic Engineer mechanism.

2.1 Cisco Express Forwarding (CEF)

The basic function of a router is to move packets through the Internet. For a router to forward packets, it needs to look up the destination IP address of the packet in a Routing Table and decide which route to use to switch the packet by using ARP Table [14]. Cisco developed Cisco Express Forwarding (CEF) for its line of routers, offering high-performance packet forwarding through the use of dynamic lookup tables. A CEF-based multilayer switch (switch/router) consists of two basic functional blocks. 1) The Layer 3 Engine is involved in building routing information (Routing Table and ARP Table). 2) The Layer 3 Forwarding Engine can be used to switch packets in hardware (Forwarding Information Base FIB and Adjacency Table). Forwarding Information Base (FIB) contains routing or forwarding information and the next-hop address for each entry. Adjacency Table consists of the MAC addresses of nodes that can be reached in a single Layer 2 next-hop. *In brief*: “CEF means route once from FIB and switches many in hardware”. At times, however, a packet cannot be switched in hardware, according to the FIB [15]. As a result, there is still a question and search for the best path for one-time in each session between the source (sender) and the destination (receiver) objects; which means that the process occurs only for one-time in Router. However, it is still repeated in all Routers from the Internet and along the trip until the data reaches its final destination. Figure 2 illustrates this issue and the time it takes to repeat this process.

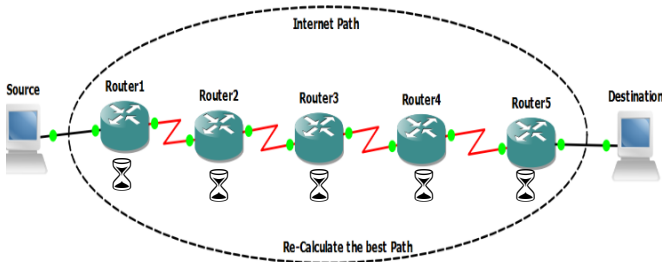


Fig. 2. Re-Calculate the Best Path on Routers

2.2 Multi-Protocol Label Switching (MPLS)

Multiprotocol Label Switching (MPLS) has been around for several years. It is a popular networking technology that advertises labels attached to IP packets between

routers to forward the packets across the Internet. Where routers build a label-to-label mapping so forward the traffic by looking at the label and not the destination IP address. MPLS is a great benefit to the service providers that deploy it and to their customers, because one of the reasons for a label-swapping protocol is the need for speed [14]. MPLS was created to combine the benefits of connectionless Layer 3 routing and forwarding with connection oriented Layer 2 forwarding [13].

2.3 Traffic Engineer

With IoT the delivery of Internet communications services has become very competitive and end-users are demanding very high quality service from their service providers. Consequently, performance optimization of large scale IP networks, especially Internet backbones, has become an important problem [16]. Traffic engineering (TE), or the ability to steer traffic through the Internet, is to get the traffic from edge to edge in the Internet in the most optimal way. For example, TE can bring a solution by steering the traffic or a portion of it away from the overloaded links, because the forwarding paradigm of IP is based on Routing Protocol mechanisms, which is least-cost path forwarding (best path). The IP forwarding paradigm does not take into account the available bandwidth capacity of the link, which might differ significantly from the cost that is assigned to the link [14].

2.4 Lessons Learned

It is noticeable that all these genius solutions are not exposed to the root of the problem! These solutions are made from the perspective of the devices themselves or a service provider not from the infrastructure of the Internet as a whole. The problem is to re-calculate the Internet best paths to final destination again and again in the Internet Routers during a data trip. This leads to excess consumption of the processor, memory and time in Internet devices. *The only benefit of a repetition of this process is to select the best path in an accurate way, which can be achieved without re-calculate this process.* This paper attempts to prevent the process of re-calculation of the best path using the proposed “Traffic Engineer System” (see Figure 2).

3 Traffic Engineer System

There is a massive growth of the Internet towards the Internet of Things. As a result of this, there are an enormous diversity and tremendous increase will happen in objects *Addressing*. We need a system able to help source object data to reach its final destination through these addresses. Where calculating and determining the full best path depends on final destination address. But for one-time in Routers along the trip to save time and thereby increases the speed. *Traffic Engineer System TES* will perform this function on Internet of Things. Like other TCP/IP-based services, TES is a protocol that works on servers. These servers maybe in 1) Internet Routers, 2) DNS servers (As additional service to DNS), or 3) *NEW* servers. This paper prefers the third option and will work to create it; this paper is the first step. The idea of working

for this system is derived from some of the systems and protocols, such as: Domain Name Services (DNS) as a system and Open Shortest Path First (OSPF) as a protocol [9, 10, 17, and 18].

3.1 TES Server

TES server provides the source object with the full Internet best path to final destination; instead of the source's data continues to ask Internet Routers repeatedly along the trip about it. Thus, the objects should have the TES server IP like DNS server IP.

- 1) The source object sends a unicast query containing the IP for destination object to a TES server.
- 2) The source object eventually receives a reply, which includes the full best path for the destination object through the internet.

Source object can then arrive to the destination faster; where transition will be with one address and switch all the time in hardware.

Since the best paths to all possible Internet final destinations are huge and a single server might not hold all these paths. As a result, there is a good suggestion that the source object's TES server (*Nearest server*) asks another TES server. At most calculated that process twice; the first at the TES Nearest server and the second at the *Service Provider TES server*. Please see Figure 3. We conclude from this that TES adds an additional delay -sometimes substantial- to the Internet applications that use it. Fortunately, it is possible to *cache* the desired best paths in TES Nearest server, which helps to reduce TES network traffic as well as the average TES delay. Also, it is possible to use a *local cache* for the best path in IoT objects themselves for 6 minutes for example; to get rid of TES-requests to any TES server as well as to get rid of TES reply time.

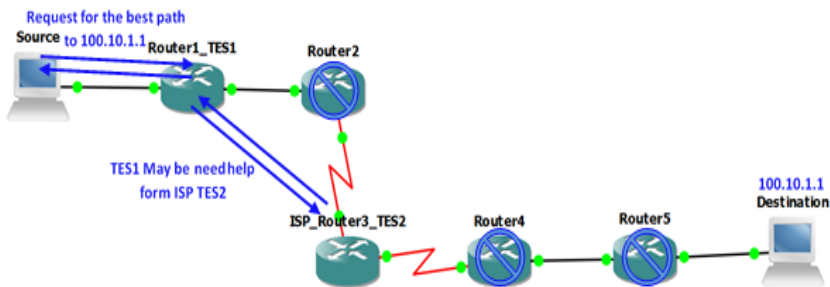


Fig. 3. Traffic Engineer System Servers

3.2 Database of Topology

TES aims to find the full best path to any destination in the Internet. In operation, each TES server sends an Internet network design (topology) to other servers on the Internet to build its database. After that database is constructed, the server filters the

database to the best paths only. The *Best-Path* table is populated with these resulting best paths per destinations in the Internet. From this table, TES servers respond to the query from source objects and give it the full best path for its data to travel to the final destination. In case of any changes in the Internet network TES servers update its databases and then the Best-Path table. TES servers for sure require more processor power, memory size, and high bandwidth, which is easy nowadays with the tremendous progress in Internet bandwidth and the hardware industry.

3.3 TES Best Path Format and Some Facts

This section represents the biggest challenges this paper faces and will help us to overcome them. The small size of the Maximum Transfer Unit MTU in Internet, which equals 1500 bytes, may represent a significant impediment to TES. TES may cause an increase in MTU size if it is designed freely. This paper did not want to change the MTU in this paper to demonstrate the principle feasibility of the idea at first, although there is tremendous progress in bandwidth and the hardware industry making it possible.

Firstly, we need to determine the size of the additions that are added by Link and IP/ID layers to know the available size to design the "best path" additions. Do we need to change this or not? The maximum size of additions in bytes (header and tail) of the link layer service provided by Ethernet to packet to build the frame is 38 bytes [20]. The maximum size of additions in bytes (header) of the IP/ID layer service provided by IPv4 to segment to build the packet is 60 bytes [11].

Secondly, we have an old technique that can help us understand what this paper wants to do in depth, which is *Source Routing*. In IP packet header options, source routing allows a sender of a packet to partially or completely specify the route the packet takes through the network. There are restrictions on a lot of Internet devices that do not support this feature for security reasons, but this issue was solved in the previous paper with Security Layer in IoT Communication Reference Model [7]. For example: Loose source routing option, in which a series of IP addresses for router interfaces is listed (up to nine addresses). The packet must pass through each of these addresses, although multiple hops may be taken between the addresses. This means that each router, instead of examining the destination IP in traveling packet and choose the next hop to forward the packet to, in source routing, the source takes some or all of these decisions by itself. This way, it removes the decision-making from the routers and puts it into the hands of the users [11, 20, 21, and 22].

Thirdly, the average AS-hops between any source and any destination in the Internet is 3 or just a bit above 3 and the maximum is around 13. The most frequent total Router-hops is less than 15 hops and the largest distance and worst case in total Router-hops is between 25 and 30 hops in over 95% of all possible traffic through the Internet [23].

4 Results and Discussions

4.1 New Frame Format

This paper will add the full best path to final destination instead of the source and destination IP addresses to data segment. However, unlike Source Routing, it will add a series of IP addresses for Routers themselves (Router ID) to represent the full best path. It chooses Router ID to make it easy to routers to select appropriate link to its neighbors. The series consists of 15 IP addresses as a maximum, because we have 60 bytes for this; where: {15 Router * 4 bytes IP address = 60 bytes} and the most frequent total Router-hops in Internet is less than 15 hops. However, what will happen if total Router-hops are more than 15? Furthermore, there is a good suggestion that the source object's data requests the rest of full best path again from the last TES server/router (*Transit* server), where it stopped due to depletion of the IP addresses series. Therefore, it must retain destination address in frame header. At most calculated that process twice; the first at the *Nearest* TES server and the second at the *Transit* TES server. Figure 4. Illustrates the new data frame format.

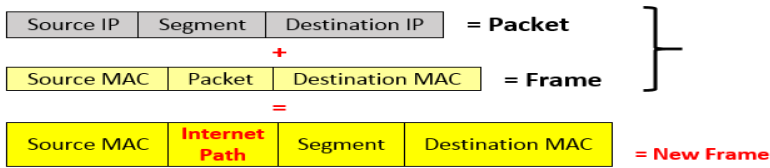


Fig. 4. New Data Frame Format

Thus in the worst case with TES, the question about the full best path will be repeated for only three times in Nearest, Service Provider, and Transit TES servers. In the perfect case, Data will be directed to their final destination without any question when we use a local cache in IoT objects.

All we need in the near future is to study the implementation of the idea of TES with all kinds of addresses in IoT, whether with IPv6 or other addresses of IoT objects. In fact, there is a simple solution to this issue. Like MPLS, routers can keep working with IPv4 in partial of Router's IDs and IPv6 for IoT objects. Thus, when any object in IoT wants to communicate with another object, the full best path will be written in IPv4 (series of Routers IDs).

4.2 IoT Communication Reference Model

Now we are getting to the most important results of this paper and TES, which is its impact on IoT Communication Reference Model. Based on the fact that routers transition data through one address and the data have its full best path, this paper merges two addresses layers from the model in one new layer. Figure 5 illustrates the new model with new proposed layer *Addressing Layer*, which will replace the two

layers IP/ID and Link. The IoT Communication Reference Model is built, bottom to top, in the following order: Physical, Quality of Service, Security, *Addressing*, End-to-End, and Data.

The new layer: *Addressing* layer is classified as Layer 4. The major function of *Addressing* layer is to provide objects in IoT with the full best path to their final destination.

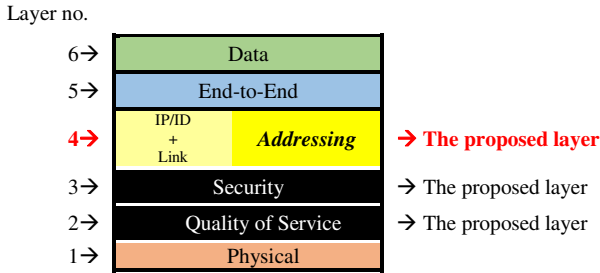


Fig. 5. New IoT Communication Reference Model

5 Conclusions

IoT will achieve '6A connectivity' (i.e., any time, any one, any thing, any place, any service, and any network) eventually as the vision of ITU and European project cluster (CERP-IoT). Building this infrastructure of any NETWORK remains the biggest challenge for driving future ubiquitous and pervasive computing [24]. The Traffic Engineer System (TES) is the proposed solution to this challenge. TES is a distributed database implemented in a hierarchy of TES servers and TCP/IP protocol that allows source objects in IoT to query the full best path for specific destination objects. With TES Source object can arrive to the destination faster; although transition will be with one address and switch all the time in hardware. In the worst case with TES, the question about the full best path will be repeated for only three times in Nearest, Service Provider, and Transit TES servers. In the perfect case, Data will be directed to their final destination without any question when it uses a local cache in IoT objects.

In the end as authors to this paper, we should recognize that there is a great challenge facing this paper. This paper requires a lot of changes and testing at the level of today's operating systems and devices in addition to creating a new system (TES), which requires support of one sponsors at implementation stage as a project and a specialized team work in many disciplines of Information Technology.

Acknowledgments. Our thanks to Eng. Mohammed Alissa, Senior Computer Engineer, Royal Saudi Air Force, Eng. Mohammed Alhoraiby, IT Infrastructure Director, King Saud bin Abdulaziz University for Health Sciences, and Dr. Nashwa Elbendary, Assistant Professor, Arab Academy for Science and Technology, for their valuable time they spent with us in a debate on TES and its repercussions.

References

1. Mattern, F., Floerkemeier, C.: From the Internet of Computers to the Internet of Things. Springer, Heidelberg (2010)
2. Ma, H.-D.: Internet of Things: Objectives and Scientific Challenges. Journal of Computer Science and Technology. Springer Science + Business Media, LLC & Science Press, China (2011)
3. Atzori, L., Iera, A., Morabito, G.: The Internet of Things: A survey. Computer Networks Journal 54, 2787–2805 (2010)
4. Gubbi, J., Buyya, R., Marusic, S., Palaniswami, M.: Internet of Things (IoT): A vision, architectural elements, and future directions. Future Generation Computer Systems Journal (2013)
5. Oteafy, S.M.A., Al-Turjman, F.M., Hassanein, H.S.: Pruned Adaptive Routing in the Heterogeneous Internet of Things. In: IEEE Global Communications Conference, GLOBECOM 2012 (2012)
6. van den Dam, R.: Internet of Things: The Foundational Infrastructure for a Smarter Planet. In: Balandin, S., Andreev, S., Koucheryavy, Y. (eds.) NEW2AN 2013 and ruSMART 2013. LNCS, vol. 8121, pp. 1–12. Springer, Heidelberg (2013)
7. Alhamedi, A.H., Aldosari, H.M., Snael, V., Abraham, A.: Internet of Things Communication Reference Model. In: The 6th CASoN Conference Porto, Portugal (2014)
8. Graziani, R. (ed.): IPv6 Fundamentals, A Straightforward Approach to Understanding IPv6. Cisco Press (2013)
9. Dulaney, E., Harwood, M. (eds.): CompTIA Network+ N10-005 Authorized, 4th edn. Pearson, Exam Cram (2012)
10. Lammle, T. (ed.): CompTIA Network+ Study Guide, 2nd edn., Exam N10-005. Sybex (2012)
11. Doyle, J. (ed.): CCIE Professional Development, Routing TCP/IP, Volume I, A detailed examination of interior routing protocols. Cisco Press (1998)
12. Doyle, J., Carroll, J.D.H. (eds.): CCIE Professional Development, Routing TCP/IP, Volume II, A detailed examination of exterior routing protocol and advanced IP routing issues. Cisco Press (2011)
13. Pepelnjak, I., Guichard, J. (eds.): MPLS and VPN Architectures, CCIP Edition, Prepare for CCIP certification as you learn to design and deploy MPLS-based VPNs. Cisco Press (2002)
14. De Ghein, L. (ed.): MPLS Fundamentals, A Comprehensive Introduction to MPLS Theory and Practice. Sybex (2007)
15. Hucaby, D. (ed.): CCNP Self-Study, CCNP BCMSN, Official Exam Certification Guide, 4th edn. Cisco Press (2007)
16. Awduche, D., Chiu, A., Elwalidl, A., Widjaja, I., Xiao, X. (eds.): Overview and Principles of Internet Traffic Engineering. RFC 3272. The Internet Society (2002)
17. Kurose, J.F., Ross, K.W. (eds.): Computer Network, A Top-Down Approach, 6th edn. Pearson (2013)
18. Stewart, B.D., Gough, C. (eds.): CCNP Self-Study, CCNP BSCI, Official Exam Certification Guide, 4th edn. Cisco Press (2007)
19. University of Aberdeen, School of Engineering, Research Information, Prof. Godred Fairhurst,
<http://www.erg.abdn.ac.uk/~gorry/eg3567/lan-pages/enet-calc.html>

20. Postel, J. (ed.): Internet Protocol, Darpa Internet Program, Protocol Specification. RFC 791 (1981)
21. Malkin, G. (ed.): Traceroute Using an IP Option. RFC 1393 (1993)
22. Soliman, M., Nandy, B., Lambadaris, I., Ashwood-Smith, P.: Source Routed Forwarding with Software Defined Control, Considerations and Implications. ACM (2012) 978-1-4503-1779-5/12/12
23. Ph.D. Internet researcher Geoff Huston, BGP Routing Table Analysis Reports, <http://bgp.potaroo.net/as6447/>
24. Park, S., Crespi, N., Park, H., Kim, S.-H.: IoT Routing Architecture with Autonomous Systems of Things. IEEE World Forum on Internet of Things, WF-IoT (2014)

Modeling Cloud Computing Risk Assessment Using Machine Learning

Nada Ahmed¹ and Ajith Abraham^{1,2,3}

¹ Faculty of Computer Science and Information Technology,
Sudan University of Science Technology, Khartoum, Sudan
naessa@pnu.edu.sa

² Machine Intelligence Research Labs (MIR Labs),
Scientific Network for Innovation and Research Excellence, WA, USA

³ IT4Innovations- Center of excellence, VSB -Technical University of Ostrava, Czech Republic
ajith.abraham@ieee.org

Abstract. Cloud computing emerged in recent years as the most significant developments in modern computing. However, there are several risks involved in using a cloud environment. To make the decision of migrating to cloud services there is a great need to assess the various risks involved. The main target of risk assessment is to define appropriate controls for reducing or eliminating those risks. We conducted a survey and formulated different associated risk factors to simulate the data from the experiments. We applied different feature selection algorithms such as Best-First, and random search algorithms methods to reduce the attributes to 3, 4, and 9 attributes, which enabled us to achieve better accuracy. Further, seven function approximation algorithms, namely Isotonic Regression, Randomizable Filter Classifier, Kstar, Extra Tree, IBK, multilayered perceptron, and SMOreg were selected after experimenting with more than thirty different algorithms. The experimental results reveal that feature reduction and prediction algorithms is very efficient and can achieve high accuracy.

Keywords: Cloud computing, classification algorithms, data mining, feature selection.

1 Introduction

Cloud computing is one of the most significant developments in modern computing, where computing resources such as: processing and storage are being offered as on demand services to individuals, companies, and government agencies, with users employing cloud computing for database management and mining, sharing and storing information, and deploying web services [1]. On an operational level, cloud computing free up the resources and refocusing them on core business activities, thereby, the potential for innovation is increased. A recent Gartner research report predicts that the global cloud market is expected to burst in the coming years [2]. The emergence of cloud computing represent a fundamental change in the way information technology service is invented, deployed, developed, maintained, scaled, updated, and paid

[3]. Consumers in cloud computing use services as needed, shared resources as a service that can rapidly and elastically scale up and down as needed, and pay only for what is used, all these things are characterizing the cloud computing [4][5].

It provides a level of abstraction between the physical infrastructure and the owner of the information being stored and processed because it stores the application software and databases in large data centers, where the management of the data and services are not trustworthy [6]. In recent years there is obvious migration to cloud computing with end users, quietly handling a growing number of personal data, such as photographs, music files, book marks, and much more, on remote servers accessible via a network [7]. The use of cloud computing services can cause risks to consumers. Before consumers start using cloud computing services they must confirm whether the product satisfies their needs and understand the risks involved in using this service [8].

In this paper, we present the application of data mining approach to assess the various risk factors. Feature selection techniques are used to reduce features and achieve better accuracy. We further applied different function approximation algorithms to find out which one is performed best over different datasets.

The rest of this paper is organized as follows. In Section 2 we give a brief summary about the cloud computing risk factors. Section 3 discusses feature selection methods followed by the different prediction algorithms in Section 4. Experiments are shown in Section 5 and finally conclusions are provided in the last Section.

2 Cloud Computing Risk Factors

We define the various risk factors associated with cloud computing as follows:

2.1 Authentication and Access Control (A&AC)

Organization's private and sensitive data must be secure and only authenticated users can access it. When using cloud, the data is processed and stored outside the premise of an enterprise, which brings a level of risk because outsourced services bypass the "physical, logical, and personnel controls", any outside or unwanted access is denied.

2.2 Data Loss (DL)

Data loss means that the valuable data disappear without a trace. Cloud customers need to make sure that this will never happen to their sensitive data.

2.3 Insecure Application Programming (IAP)

APIs is an important and necessary part to the security and availability for whole cloud services. Building interfaces, injecting services will increase risk, there for some organization may in force to relinquish their credentials to third party in order to enable their agency.

2.4 Data Transfer (DT)

Sensitive data are obtained from customers, processed and stored at the cloud provider end. All data flow over network needs to be secured in order to prevent seepage of customer's sensitive information. The application provided by the cloud provider to their customers has to be used and managed over the web. The risk comes from the security holes in the web applications.

2.5 Insufficient due Diligence (IDD)

Before using the cloud services, the organization needs to fully understand the cloud environment and its associated risk.

2.6 Shared Environment (ShE)

Multi-tenancy is a key factor of cloud computing services. To achieve scalability cloud provider provide shared infrastructure, platform, and application to deliver their services. This shared nature enables multiple users to share same computer resources, which may lead to leaking data to other tenants, also, if one tenant carried malicious activities the reputation of other tenants may be affected.

2.7 Regulatory Compliance (RC)

If the provider is unable or unwilling to subject to external audits and security certification, and they do not give their customers any information about the security controls that have been evaluated. It should only be considered for most trivial functions. Regardless of the location, the custodian is ultimately responsible for ensuring the security, protection, and integrity of the data, especially when they are passed to a third party.

2.8 Data Breaches (DB)

Breaching into a cloud environment will potentially attack all users' data. Those attackers can exploit a single flaw in one client application to get to all other client's data as well, if the cloud service databases are not designed properly.

2.9 Business Continuity and Service Availability (BC& SA)

The nature of the business environment, competitive pressure, and the changes happening in it leads to some events that may affect the cloud service provider, such as a merger, goes broke, bankruptcy, or its acquisition by another company. These things lead to loss or deterioration of service delivery performance, and quality of service. Another important thing to the cloud-computing provider is that their customers must be provided with service around the clock, but outages do occur and can be unexpected and costly to customers.

2.10 Data location and Investigative Support (DL&IS)

Most cloud service providers have many data centers around the globe. Regarding privacy regulation in different jurisdictions, in different countries where the government restricts the access to data in their borders, or if the data stored in high-risk countries, all these things make data location big concern issue. The investigation of an illegal activity may be impossible in a cloud computing environment, because multiple customer's data can be located in different data centers that are spread around the globe. If the enterprise relies on the cloud service for the processing of business records then it must take into account the factor of the inability or unwillingness of the provider to support it.

2.11 Data Segregation (DS)

The risks arise here come from the failure of the mechanisms to separate data in storage, and memory, from multiple tenants in the shared infrastructure.

2.12 Recovery (R)

Cloud users do not know where their data is hosted. Some events such as man-made, or natural disaster may happen; in such events customers need to know what happened to their data and how long the recovery process will take.

2.13 Virtualization Vulnerabilities (VV)

Virtualization is one of the fundamental components of the cloud service. However, it introduces major risks as every cloud provider uses it. Beside its own risks it holds every risk posed by physical machines.

2.14 Third Part Management (TPM)

There are many issues in cloud computing related to third party because the cloud service provider does not directly manage the client organizations. Some old concerns in information security appear with outsourcing such as integrity control and sustainability of supplier and all risks that clients may take if it relies on a third party.

2.15 Interoperability and Portability (I& P)

Interoperability and portability become crucial because if the organization locks to a specific cloud provider, then the organization will be at the mercy of the service level and pricing policies of that provider and it won't have the freedom to work with multiple cloud provider.

2.16 Resource Exhaustion (RE)

Cloud provider allocates resource according to statistical projections. Inaccurate modeling of resource usage can lead to many issues such as: service unavailability, access control compromised, economic and reputational losses, and infrastructure oversize.

2.17 Service Level Agreement (SLA)

The organization needs to ensure that the terms of (SLA) are being met. Risk may appear with service level application such as the data owner as some cloud provider include explicitly some terms state that the data stored is the provider's not the customer's. In the few cases where cloud vendors went out of business, their customer private data were sold as part of the asset to the next buyer. Also, SLA terms should include licensing conditions. There is the possibility for creating original work in the cloud, but if not protected by the appropriate contractual clauses, this original work may be at risk. One of the SLA terms must be for responsibilities of cloud provider for enabling governance.

2.18 Data Integrity (DI)

One of the most critical elements in all systems is data integrity. Cloud computing magnified the problem of data integrity and endangers the data integrity in transaction management, at the protocol level, which does not support transactions or guaranteed delivery. If data integrity is not guaranteed and there is a lack of integrity controls, this may result in deep problems.

3 Feature Selection

Data Mining is an iterative process within which the progress is defined by discovery of earlier, unidentified, valid patterns, and relationship in large dataset, through either automatic or manual tools [9, 10, 11]. Feature selection is the process of extracting subset instances from the original data set and it presents as an important technique in data preprocessing in data mining [12]. The goal of feature selection is to find the minimum number of related attributes from n attributes, which describe the concept as close as possible of the original set does, and perform the best [13]. Attribute selection reduces dataset size by removing irrelevant and redundant attributes. Applying feature selection technique involve both search algorithm and criterion function and the search algorithm applies the criterion function as a measure of the effectiveness of each feature subset, to compare and generate suitable solutions of the feature selection problem [14]. For feature selection we used Best-first, Random search, and ranker methods. Best-first [13] method search in the space of feature subset, and maintain a queue of possible solution. It deals with space of set as a graph called "feature selection lattice", then applies one of standard graph searching algorithms. In random search [12], first it randomly select subset, then continue in two different ways. One of them is to follow sequential search. The second is to continue randomly and generate the next subset randomly.

4 Risk Assessment Algorithms

Risk assessment involves building an accurate prediction model. Prediction comprises of two steps: in the first step called the training phase a predictor is built describing a predefined set of training data and in the second step the model is used for test data [15][16]. Following are the predictors used.

A) Extremely Randomized Decision Trees

It is tree-based ensemble method for supervised classifier (commonly known Extra-Trees). Extra-Trees based on randomization process, the splitting rules are randomly drawn at each node of the extra-tree, and the base of chosen one of this rules to be associated with that node is the best performance according to a score computation. This allows increasing the speed of training, weakening the correlation between the induced decision trees, and reducing the complexity of the induction process [17].

B) Instance-Based Knowledge (IBK)

IBK is an Instance-Based Learning method and is an implementation of K-nearest-neighbors classifier [18]. In its representation it does not derive a rule set or decision tree and storing it, instead it uses the instances themselves to represent what is learned [11]. IBK compare each new instance with existing ones using distance metric; most commonly Euclidean distance and the closest existing instance is used to assign the class for the test sample [15].

C) Sequential Minimal Optimization (SMOreg)

Sequential Minimal Optimization is an iterative algorithm for solving the regression problem using Support Vector Machine (SVM)[19]. This algorithm does it job by replacing all missing values and transforms nominal attributes in to binary ones [20]. The notable features of SMO algorithm is it does not require Quadratic programming solver, this make it differ from most SVM algorithms [20].

D) Multilayered Perceptron (Artificial Neural Networks)

A multilayer perceptron is a feed forward artificial neural network model that maps sets of input data onto a set of appropriate outputs. The fundamental aspects which a neural network depends upon are, input and activation function of the unit, network architecture and the weight of each input connection [21].

E) K- Nearest Neighbors (K-NN or K*)

K-NN is an instance-based learning algorithm that stores all training instances and does not build a model until a new instance need to be classified [9] and they use some domain specific distance function to retrieve single most similar instance from the training set [22].

F) Isotonic Regression

Is a regression method, it does its job by weighted least squares to evaluate linear regression models [23].

G) Randomizable Filter Classifier

Used for running an arbitrary classifier on data that has been passed through an arbitrary filter. Like the classifier, the structure of the filter is based exclusively on the training data and test instances will be processed by the filter without changing their structure [24].

5 Experimentation Results

We conducted a survey to finalize the risk factors. In the survey, we asked the participants to categorize risk factors to three levels according to the likelihood of happening and their effect on cloud computing. These categories are: Important, Neutral, and Not Important. 35 international experts responded to the survey from different countries and all of them agreed that the previously defined factors are important, which means that they have great effect over cloud computing. Next we give each risk factor a numeric range of values, and finally we formulated expert rules and use some statistical methods to generate the data based on the rules. The Dataset contains 18 input attributes and comprise of 1940 instances. The 18 attributes were labeled as DI, IDD, RC, BC&SA, TPM, I&P, DL, IAP, DL&IS, R, RE, SLA, A&AC, ShE, DB, VV and DI. Risk Factors are illustrated in Table 1 with their corresponding ranges for numeric values. The feature selection methods reduce the number of features to 3 (*first dataset*), 4 (*second dataset*), and 9 (*third dataset*). We used percentage split to test and evaluate the algorithms. In percentage split the dataset is randomly splitting to training and testing data as follows:

- 60% - 40% (A)
- 70% - 30% (B)
- 80% - 20% (C)
- 90% - 10% (D)

The experiments were performed in WEKA that provides a collection of machine learning algorithms and data preprocessing tools in a graphical user interface environment for data exploration and algorithm evaluation [18].

Table 1. Risk factors and their associated range values

Risk Factor	Range value	Risk Factor	Range value
<i>DT</i>	0 - 3	<i>R</i>	1 - 3
<i>IDD</i>	1 - 3	<i>RE</i>	0 - 2
<i>RC</i>	0 - 1	<i>SLA</i>	0 - 3
<i>BC & SA</i>	1 - 3	<i>A & AC</i>	0 - 3
<i>TPM</i>	0 - 2	<i>ShE</i>	1 - 3
<i>I & P</i>	0 - 1	<i>DB</i>	0 - 2
<i>DL</i>	0 - 3	<i>DS</i>	0 - 1
<i>IAP</i>	0 - 1	<i>VV</i>	1 - 3
<i>DL & IS</i>	0 - 3	<i>DI</i>	0 - 2

Performance statistics are calculated across all datasets using Root Mean Square Error (RMSE) and Correlation Coefficient (CC) but since CC is almost (0.9999 or 1) we did not include them in the tables. We apply attribute selection method to reduce the number of the attributes. In the preprocessing step, the data is filtered to remove the irrelevant data and improve the quality.

We investigated three feature reduction algorithms and finally managed to reduce to 4, 5, and 10 attributes. Table 2 summarizes the best test results of (RMSE) from each data percentage with each algorithm. Figure I represent the RMSE for all dataset percentages and shows that more percentage of training data (90% training and 10 % testing) produces the best results, which means better learning. Table 3 summarizes the best results of (RMSE) from each algorithm among all data sets. Figure 2 represents the RMSE we obtained using the different algorithms for all datasets. It shows that the K star algorithm did not perform well for first, and second datasets, and a little better with the third dataset. Multilayered Perceptron and SMOreg algorithms performed well with the first and second datasets and give the worst result with the third dataset. Randomizable Filter Classifier, Extra Tree, IBK, and Isotonic Regression algorithms performed well among all datasets.

Table 2. RMSE Performance for each dataset

Algorithm	A	B	C	D
	Root mean squared error (RMSE)			
Isoreg	0.0021	0.0019	0.0018	0.0017
MLP	0.0019	0.0015	0.0016	0.0006
SMOreg	0.0026	0.0019	0.0023	0.0024
IBk	0.0020	0.0018	0.0018	0.0017
K-Star	0.0045	0.0046	0.0037	0.0034
RFC	0.0021	0.0019	0.0018	0.0017
Extra Tree	0.0042	0.0035	0.0034	0.0031

Table 3. The best test results from different algorithms

Algorithm	First dataset	Second dataset	Third dataset
Isoreg	0.0017	0.0017	0.0017
MLP	0.0006	0.0016	0.0118
SMOreg	0.0019	0.0028	0.0366
IBk	0.0017	0.0017	0.0017
KStar	0.0177	0.0094	0.0034
RFC	0.0018	0.0018	0.0017
ExtraTree	0.0031	0.0032	0.0032

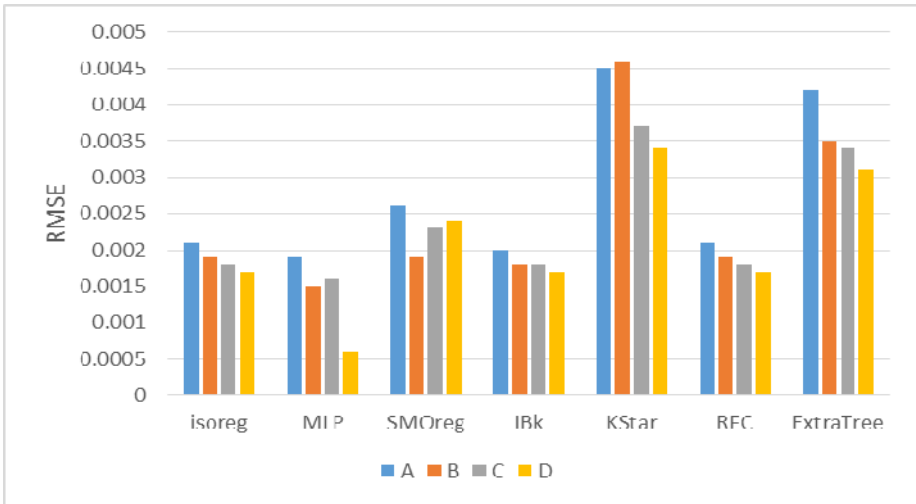


Fig. 1. Best RMSE for test data using different datasets (60, 70, 80, 90 %)

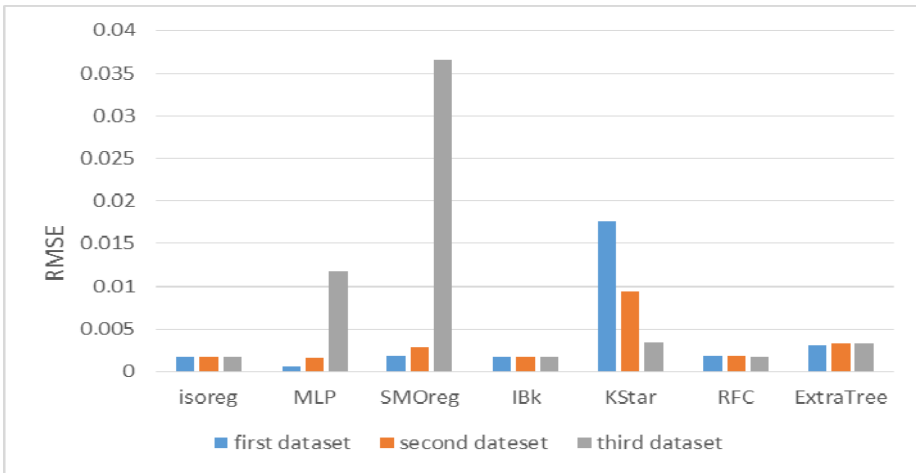


Fig. 2. Best RMSE for test data for different algorithms

6 Conclusions

The purpose of this experimental work is to achieve feature reduction, best accuracy on test data set and to find out the best available approach applied to our dataset. We examined the performance of several machine-learning algorithms for modeling the cloud computing risk environment.

The effect of subsets of training and testing data is also illustrated here by split the sub dataset randomly into four different groups. The empirical results show that the splitting of dataset to (90% - 10%) gives the best result among all other partitioning,

and also show that Randomizable Filter Classifier, Isotonic Regression, IBK, and Extra Tree algorithms gives the best results among all data sets.

References

1. Paquette, S.J., Wilson, P.T., Susan, C.: Identifying the security risks associated with governmental use of cloud computing. *Government Information Quarterly* 27, 245–253 (2010)
2. Brender, N., Markov, I.: Risk perception and risk management in cloud computing: Results from a case study of Swiss companies. *International Journal of Information Management* 33, 726–733 (2013)
3. Avram, M.: Advantages and Challenges of Adopting Cloud Computing from an Enterprise Perspective. *Procedia Technology* 12, 529–534 (2014)
4. Carroll, M., van der Merwe, A., Kotze, P.: Secure cloud computing: Benefits, risks and controls. In: *2011 Information Security South Africa (ISSA)*, pp. 1–9 (2011)
5. Sun, D., Chang, G., Sun, L., Wang, X.: Surveying and analyzing security, privacy and trust issues in cloud computing environments. *Procedia Engineering* 15, 2852–2856 (2011)
6. Subashini, S., Kavitha, V.: A survey on security issues in service delivery models of cloud computing. *Journal of Network and Computer Applications* 34, 1–11 (2011)
7. Zissis, D., Lekkas, D.: Addressing cloud computing security issues. *Future Generation Computer Systems* 28, 583–592 (2012)
8. Chandran, S.A., Mridula: Cloud Computing: Analyzing the risks involved in cloud computing environments. In: *Proceedings of Natural Sciences and Engineering*, pp. 2–4 (2010)
9. Phyu, T.N.: Survey of classification techniques in data mining. In: *Proceedings of the International MultiConference of Engineers and Computer Scientists*, pp. 18–20 (2009)
10. Kantardzic, M.: *Data mining: concepts, models, methods, and algorithms*. John Wiley & Sons (2011)
11. Chauhan, H.K., Pundir, V., Pilli, S., Emmanuel, S.: A Comparative Study of Classification Techniques for Intrusion Detection. In: *2013 International Symposium on Computational and Business Intelligence (ISCBI)*, pp. 40–43 (2013)
12. Liu, H.Y., Lei: Toward integrating feature selection algorithms for classification and clustering. *IEEE Transactions on Knowledge and Data Engineering* 17, 491–502 (2005)
13. Jain, A.Z., Douglas: Feature selection: Evaluation, application, and small sample performance. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 19, 153–158 (1997)
14. Serpico, S.B.B., Lorenzo: A new search algorithm for feature selection in hyperspectral remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing* 39, 1360–1367 (2001)
15. Ali, S.S., Kate, A.: On learning algorithm selection for classification. *Applied Soft Computing* 6, 119–138 (2006)
16. Han, J.K., Micheline: *Data Mining, Southeast Asia Edition: Concepts and Techniques*. Morgan kaufmann (2006)
17. Désir, C.P., Heutte, C., Salaun, L., Thiberville, M., Luc: Classification of endomicroscopic images of the lung based on random subwindows and extra-trees. *IEEE Transactions on Biomedical Engineering* 59, 2677–2683 (2012)
18. Witten, I.H.F., Trigg, E., Hall, L.E., Holmes, M.A., Cunningham, G., Jo, S.: *Weka: Practical machine learning tools and techniques with Java implementations* (1999)

19. Shevade, S.K.K., Bhattacharyya, S.S., Murthy, C., Krishna, K.R.: Improvements to the SMO algorithm for SVM regression. *IEEE Transactions on Neural Networks* 11, 1188–1193 (2000)
20. Karagiannopoulos, M., Anyfantis, D., Kotsiantis, S.B., Pintelas, P.E.: Feature selection for regression problems. In: *Proceedings of HERCMA 2007* (2007)
21. Kotsiantis, S.B., Zaharakis, I.D., Pintelas, P.E.: *Supervised machine learning: A review of classification techniques* (2007)
22. Cleary, J.G.T., Leonard, E.: K^* : An Instance-based Learner Using an Entropic Distance Measure. In: *ICML*, pp. 108–114 (1995)
23. Wu, C.-H., Su, W.-H., Ho, Y.-W.: A study on GPS GDOP approximation using support-vector machines. *IEEE Transactions on Instrumentation and Measurement* 60, 137–145 (2011)
24. <http://weka.sourceforge.net/doc.dev/weka/classifiers/meta/package-summary.html>

Adaptation of Turtle Graphics Method for Visualization of the Process Execution

Jakub Štolfa¹, Svatopluk Štolfa¹, Martin Kopka¹, and Václav Snášel^{1,2}

¹ Department of Computer Science, FEI, VŠB - Technical University of Ostrava,
17. listopadu 15, 708 33, Ostrava-Poruba, Czech Republic

² IT4Innovations, VŠB - Technical University of Ostrava,
17. listopadu 15, 708 33, Ostrava-Poruba, Czech Republic

{jakub.stolfa,svatopluk.stolfa,martin.kopka,vaclav.snasel}@vsb.cz

Abstract. Process mining is relatively young discipline that uses many methods to obtain a results that can be beneficial. These methods are used in different approaches that study the process instances from the statistical point of view, clustering methods are used, etc. Our intention is to present simple graphical method that can help to understand the visualized data directly. Adaptation of the turtle graphics is used to visualize the logs content - process instances. The main purpose of this paper is to present the usability of our method in the area of process mining and show its benefits for specific tasks.

Keywords: Turtle Graphics, Process mining, Process.

1 Introduction

Process mining was introduced by Aalst in 2004. This area of the research has been developing during the years and uses different methods and tools to analyze the data logs from performed processes [1, 2, 4, 5, 7]. The data are analyzed to obtain the overview whether the process was performed as it was designed, whether there are some deviations, we are trying to find interesting information about the process instances. The methods that are used varies from statistical methods throughout the different analytical methods, clustering etc. to the graphical methods. All approaches have their pros and cons. Sometimes these methods are too complex to use them for all situations. Also methods used for visualization provide complex process diagrams. Thus we were thinking about an easy to use method that can be used for the data analysis and is simple enough to be used as a first analytic tool to study the processes. Our idea is to use some graphical representation of the data that will be easy to interpret. We have decided to adopt turtle graphics approach to visualize the process - we are interested in process flow and process result too.

The paper is organized as follows: Section 2 introduces the state of the art; Section 3 describes the analyzed process that was used for the experiments, Section 4 depicts proposed approach, Section 5 presents the experiments that we have performed and explains obtained results; concluding Section 6 provides a summary and discusses the planned future research.

2 State of the Art

In a computer graphics the turtle graphics method is a method for creating vector graphics using a relative cursor. This cursor is called turtle. It is based on the turtle geometry that is a local, coordinate free version of computational geometry. Turtle geometry has been used to study many diverse subjects from simple polygons to complex fractals, from the Euler characteristic and the formula of Gauss-Bonnet to curved space-time and Einstein's general theory of relativity [3]. Turtle graphics method became popular method how to teach programming to young children [8].

Our virtual turtle is the creature than know only its position, direction it is facing and the step size. Turtle can respond to four basics commands:

- Forward x : the turtle moves forward x steps in the direction it is facing. During that move the turtle draws a line from initial position to its final position. It is like the turtle drag its tail in the sand and making line.
- Move x : the turtle moves forward x steps in the direction it is facing but do not draw a line.
- Turn α : the turtle changes counterclockwise its direction according the angle α
- Resize s : this command change the length of the turtle step

The turtle's location can be represented by a point P given by a pair of coordinates (p_1, p_2) ; similarly the turtle's heading can be represented by a vector w given by another pair of coordinates (w_1, w_2) . The pair (P, w) is called the turtle's state. The turtle's state (P, w) is a complete description of what the turtle knows, and the four turtle commands forward, move, turn, resize are the only way that a programmer can communicate with the turtle [6].

This was the classics definition of the turtle graphics. In our paper we present method for representing process execution based on the turtle graphics. We tell the turtle how to turn not according where the turtle is facing, but according the sides of the area. It means we tell the turtle go to the top, left, right or bottom side. How we use this setting is described in the following sections of the paper.

3 Process Context

We have used data logs of the SAP system. This section describes the details about the data that we use for the experiment. Current SAP system runs in the company that operates in four European countries. We chose business process of the invoice verification that is implemented in SAP system, user activities are controlled by SAP workflow system. Users participate in the invoice verification workflow in several different roles (creator, completer and verifier, approver, and accountant decision maker and poster). Generally, it is process where the creator should create the invoice, verifier should verify it, send to the approvers and finally when accountant gets the invoice, he does invoice posting.

Detailed description of obtained log and data preprocessing is described in our previous work [10]. We know the architecture of the process model because user activities in the SAP are controlled by SAP business workflow. It means that process execution should follow the process model. On the other hand the process has several degrees of freedom and moreover we can find out some deviations given by processed data. This model is depicted in the Fig. 1. Process starts with event Creation. Next one is Verification. Next is Approval event. Approval event can be done repeatedly. Last event is Posting.

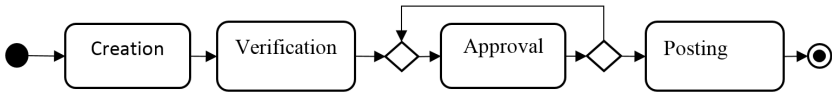


Fig. 1. Process model - modeled by UML activity diagram

4 Proposed Approach

This section presents our ideas of usage Turtle Graphics method for visualization of process execution.

We have four types of events - creation, verification, approval, posting in the examined process. We set up aliases for the type of the events in the process. It means that Verification event is V, Creation event is C, Approval is A, and, finally, Posting is S.

We define following turtle's moves according the executed processes that we have obtained from the process log - see the section 3. It means that the sequence that contains information about the process topology and duration of the particular event is the input for the turtle.

In our approach we control the turtle's behavior using fixed directions as- signed to specific events in process log. It means we tell the turtle how to turn not according where the turtle is facing, but according the sides of the area. We assigned following directions to specific events:

- event C - turtle goes up (North),
- event V - turtle goes right (East),
- event A - turtle goes down (South),
- event S - turtle goes left (West).

The turtle's state is defined as a pair (P, w) , where $P = (px, py)$ represents the coordinates of the turtle's position, vector $w = (wx, wy)$ represents the turtle's heading. The turtle's state can be written also as (px, py, wx, wy) . Initial state of the turtle is $P_i = (0, 0, 0, 1)$.

The turtle can be controlled by commands (the turtle knows its actual state $P = (px, py, wx, wy)$ before command execution):

- $move(x, y)$ - moves turtle to state (x, y, w_x, w_y) without keeping track,
- $turn(x, y)$ - moves turtle to state $(p_x, p_y, w_x + x, w_y + y)$,
- $forward(D)$ - moves turtle D steps in direction of its heading while the turtle keeps track to final state (p'_x, p'_y, w_x, w_y) , where $p'_x = p_x + D * \cos\alpha$, $p'_y = p_y + D * \sin\alpha$, $\alpha = \arctan \frac{w_y}{w_x}$

Following Table 1 contains the interpretation of specific process steps by the turtle activity.

Table 1. Interpretation of the process steps

Process step	Turtle move	Start state	1. Turn	2. Forward
C	up (North)	(p_x, p_y, w_x, w_y)	$turn(-w_x, 1 - w_y)$	$forward(1)$
V	right (East)	(p_x, p_y, w_x, w_y)	$turn(1 - w_x, -w_y)$	$forward(1)$
A	down (South)	(p_x, p_y, w_x, w_y)	$turn(-w_x, -1 - w_y)$	$forward(1)$
S	left (West)	(p_x, p_y, w_x, w_y)	$turn(-1 - w_x, -w_y)$	$forward(1)$

Explanation of the definition of turn command in our rectangular model of moves is easy. We know the direction the turtle is facing and we know to which direction we want to send the turtle based on event. Then we can count the degree of the turn like delta of the degree from the x axis which is turtle facing and the degree of the side where we would like to send the turtle. For example turtle is facing 270 degree from the x axis (last turn was A event) and we would like to send the turtle to the top (C event) - it is 90 degree. So $270 - 90 = 180$. It means that turtle have to turn 180 degree to the right.

Length (D) of the step is determined by the specific event duration. Event duration is divided to the three categories [9]:

1. Category 1: If the event lasts less than 32 hours it fits to the first category
2. Category 2: If the event lasts more than 32 hours and less than 168 hours it fits to the second category
3. Category 3: If the event lasts more than 168 hours it fits to the third category

It means if the event fits to the first category the turtle has to do one step, if it fits to the second category the turtle has to do 2 steps, and the same principle is used for the third category.

This graphical approach can show us some interesting information - the loops of the events or loops of the sequences of the events, deviation in the process execution, the resultant position of the process, etc. 1 Motivation of our research in this area and usage of the turtle graphics to visualization of the process execution is that we can find and compare interesting information in the process execution more easy way than compare it on the sequence level. For example we can compare how close is the execution of the sequences, or where these sequences have common point and the turtle graphics can tell it to us more easy way.

5 Experiments and Results

The first experiment shows four most used sequences without time parameter - CVAAS, CVAAAS, CVAS, CVAAAAS, CVAAAAS as you can see in the Fig. 2. Starting point is surrounded by circle. End points are red. When we look at figure, we can see that all sequences end in points that are on the line from the starting point to the bottom. This is the behavior that we expect from the correct sequence. The difference is only the distance from the starting point. The other finding is the shape of the picture. We can see that the picture is simple, there is only one step right - it means verification in our case. Accordingly, there is only one posting activity for all sequences. The difference is the bottom direction. Every sequence its specific number of approval activities. We can see, that there is no deviation found for the most used sequences.



Fig. 2. Most used sequences of the process visualized by Turtle Graphics

Next experiments show the deviations according to the proper process execution - CVAAS (fig. 3.). On the left side we can see that the sequence visually breaks the rule of one C and V directions. Thus there is a creation and verification made two times for this type of sequence. Otherwise the sequence is correct. The drawing in the middle shows that there is a problem at the end of the sequence, A and S were performed in addition even when the sequence was correctly ended. The third drawing on the right side shows that there is unexpected S activity directly after first A and then two A and final ending is made. These three results shows that the deviations at the beginning, at the end and in the middle of the sequence can be easily recognized by the visualization.

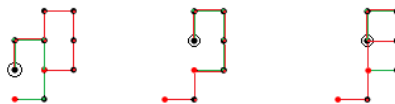


Fig. 3. Topology deviations of the sequences

These experiments led us to the simple hypothesis about resulting point of the process towards its start points. The positions are pictured in fig. 4. and explained in following list:

- A: Simple correct ideal process with one approval step and all purchase documents pre-approved. Typically raw material invoices under amount limit or periodic invoices.
- B: Complex correct process with several/many approval steps. Typically raw material invoices above amount limit or investment invoices or invoices touching more cost centers (plants).
- C, D: Less or more complex process with several/many verification steps. Typically invoices that have not completed purchasing / receipt process or for which receipt disagrees with the order. Combination with variant B is possible when the invoice has "B" parameters.
- E: Repeated posting workflow step means either discrepancy in accounting setting (unlikely event, probably some new posting case) or internal task in clerk office (responsibility, ...). These invoices should be analyzed.

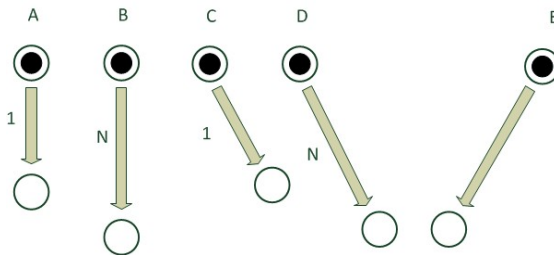


Fig. 4. Hypothesis about start and final positions of the process

Result: we checked this hypothesis through the sample of 10720 invoices with following result. Invoices, resulting in a relative position according to the described hypotheses had expected characteristics only in case when their approval process contain just one C step. As was analyzed based on this result, more C process steps means some specific mistake in the process and these cases must be analyzed separately in future (all examined processes with more C steps had a problem which coerced rerun the process from the begin).

The processes with more C steps cannot be simply identified by analysis of the final state in given model. We should use the comparison of specific subtask of the process (up-right direction) for this purpose.

Final experiment on this paper shows the visualization of the topology with the time parameter added - Fig. 5. There are three sequences shown. The topology is similar - CVAAS. The difference is in the length of each activity. We can see three different sizes of the shape. From the topological point of view, there is no problem or deviation. This is quite visible for our type of data, but we plan to study this by graphical methods in the future. From the time point of view there is a difference that is nicely visible as well. We have to say again, that even from time point of view the example and our sequences are quite simple, thus visible. The graphical methods will be used for the comparison in the future. Anyway, the visualization of such data is useful and another approach how to study the data.

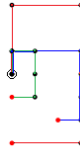


Fig. 5. Topology and time parameter visualization

6 Conclusion and Future Work

Performed experiments showed the potential of the visualization using turtle graphic. We discussed few experiments in this paper, but even these experiments show that the approach can be used for the visualization and is useful. There are many ways how to use and customize the settings to draw the paths.

We have adjusted the approach to our data, similar adjustment can be made for different data as well. The approach and the result visualization can be extended in many ways.

Following tasks are opened for the future work:

- The definition of turtle movements determines the shape of turtle's line and also the final state position. Change of the definition can make the visualization more clearer. Actually two variants of the change is prepared (1. change of the azimuth definition of specific steps and 2. definition of the step by the angle of the current direction - not fixed azimuth).
- Comparing the graphical properties of the turtle's lines with aim of comparing the similarities of whole process / subprocess and clustering the lines according to found parameters of similarity. This approach could open new access to clustering of the processes.

For example graphics comparison methods can be used to compare similar path shapes that differs e.g. only in the time manner, the lines between nodes can show the information about the usage of this by the thickness of the line, there can be information about the direction of each used line, also in the wholesale view etc. Our approach seems to be interesting and showed some potential, we would like to explore, use and describe more possible extension to it to support the readiness and interpretation of the results.

Acknowledgments. This work was supported by the internal grant agency of VŠB – Technical University of Ostrava, Czech Republic, under the projects no. SP2014/157 "Knowledge modeling, simulation and design of processes".

References

1. Van Der Aalst, W.M.P., Van Dongen, B.F., Gnther, C.W., Mans, R.S., Alves De Medeiros, A.K., Rozinat, A., Song, M., Verbeek, H.M.W., Weijters, A.J.M.M.: Process Mining with ProM. In: Belgian/Netherlands Artificial Intelligence Conference, p. 453 (2007)
2. Van Der Aalst, W., Adriansyah, A., Van Dongen, B.: Causal nets: A modeling language tailored towards process discovery (2011)
3. Abelson, H., di Sessa, A.: Turtle geometry: the computer as a medium for exploring mathematics. MIT Press, Cambridge (1986)
4. van Dongen, B.F., de Medeiros, A.K.A., Verbeek, H.M.W(E.), Weijters, A.J.M.M.T., van der Aalst, W.M.P.: The proM framework: A new era in process mining tool support. In: Ciardo, G., Darondeau, P. (eds.) ICATPN 2005. LNCS, vol. 3536, pp. 444–454. Springer, Heidelberg (2005)
5. Dumas, M., Van der Aalst, W.M.P., Hofstede, A.H.M.: Process Aware Information Systems: Bridging People and Software Through Process Technology. Wiley-Interscience (2005)
6. Goldman, R., Schaefer, S., Ju, T.: Turtle geometry in computer graphics and computer-aided design. *Computer-Aided Design* 36(14), 1471–1482 (2004) ISSN 0010-4485
7. Grigori, D., Casati, F., Castellanos, M., Dayal, U., Sayal, M., Shan, M.C.: Business process intelligence. *Computers in Industry* 3, 321–343 (2004)
8. Mindstorms, P.S.: Children, computers and powerful ideas. Basic Books, New York (1980)
9. Štolfa, J., Štolfa, S., Slaninová, K., Martinovič, J.: Searching Time Series Based on Pattern Extraction Using Dynamic Time Warping. *Dateso, Roudnice nad Labem*, 81–90 (2014)
10. Štolfa, J., Kopka, M., Štolfa, S., Koberský, O., Snášel, V.: An Application of Process Mining to Invoice Verification Process in SAP. In: 4th International Conference on Innovations in Bio-Inspired Computing and Applications, IBICA 2013, pp. 61–74 (2014)

NFL Results Predictor as a Smart Mobile Application

Petr Kakrda, Ondrej Berger, and Ondrej Krejcar

University of Hradec Kralove, Faculty of Informatics and Management, Center for Basic and Applied Research, Rokitanskeho 62, Hradec Kralove, 500 03, Czech Republic
{petr.kakrda, Ondrej.Berger, Ondrej.Krejcar}@uhk.cz

Abstract. This paper introduces the chosen area and problems which are connected to the forecasting of the results of American football games in National Football League (NFL). We cover the existing mobile applications for forecasting results which are analysed as a non-complex solutions with only limited prediction success ratio. We found that is possible to cover many other information sources to make a more complex analysis of each game and have a prediction based on a deep knowledge from match history. The suggested solution consists of a separate mobile application whose draft and implementation is described in the paper. Our algorithms implemented in developed solution provide a much better prediction success ratio than other mobile application.

Keywords: Smart, Mobile, Prediction, Algorithms, NFL.

1 Introduction

An effort to correctly forecast the result or the winner of a specific game has a long tradition and can be found in almost each sport. One of them is American football in the USA which has nearly a hundred-year history. Moreover, due to the high popularity and the possibility to bet for the winner a large amount of money is involved in the professional league, called NFL. NFL experts and fans try to estimate which team wins in a specific game. However, without certain knowledge, it would be only guessing. In order to obtain a relatively reliable forecast, it is possible to use various sources of data. For example, the forecasts from experts and specialists for the given sport who mainly use their experience for estimation of the result. With the expansion of PC, a lot of computer systems for results forecasting were developed. The forecasts are not 100% reliable as they have only a certain probability. In the case of teams that perform similarly it is very hard to predict the winner and each forecast is different. The reason for this is the diverse approach to estimating the winner of the game, the dependence of other factors, etc. There many available forecasts from various sources for suggesting the winner. Therefore, it is complicated for the fan to decide which source to trust. Finally, the following of forecasts from different sources on regular basis is time-consuming and inconvenient.

In regards to the expansion of smartphones, the mentioned problem can be solved with a mobile application which uses so called smart access [1-5]. The application

itself would prefer certain sources in a suitable way according to the success of their forecasts.

2 Results Forecasting of NFL Games

American football is one of the most watched and popular sports in the USA. Its popularity also grows in the world, especially in the Western Europe. An increasing number of fans and players can be also found in the Czech Republic. Currently, there are around 25 active clubs. NFL (National Football League) is the highest professional league of the American football with almost hundred-year old tradition. Similarly as in each sport, the fans of American football and experts are trying to estimate which team will win the game. A fan, who tries to win a bet or just wants to be precise when estimating the result of the game, has a variety of sources and data to use that can help him/her to assess who is the most probable winner.

Note: Almost always the match ends with the victory of one of the competing teams. It can also finish as a tie, but in this sport it is very rare (in the last 40 years, there were only 19 ties in NFL).

2.1 Experts

The most common source of the forecasts are the opinions of experts in the form of estimating a concrete winner or in the form of a ladder of teams according to the current form (so called Power Rankings). Another common source, and in the case of bidding the main source, are the fixed-odds bets of the betting shop for specific teams. The process itself of the winner of the game estimating is not simple. There are many factors which play a role in the game. Of course the most essential is the current form of the competing teams. The favoured team with many wins in the past has a much larger chance to win another game as has a weak competitor with many losses in a row. However, the bigger problem is to estimate the winner of a game where the two teams have similar or equal strengths. In the case of changing circumstances, it is sometimes impossible to estimate the current form of the team. Lastly but not least, it is also important to take into account the injuries of the important players. The American football matches are also characteristic for the possibility to obtain large number of points in a short time. Even the loss of a favoured team against a significantly weaker team is relatively common in NFL. All of these factors make the estimation of a game result particularly difficult. Therefore, it is only possible to choose the winner with a certain probability.

The NFL experts are trying to take into account the mentioned factors in their forecasts. Also the experiences of individual experts play a role in the estimations. The longer the expert works in the world of NFL, the more likely he/she is to estimate the winner of the upcoming games.

2.2 Systems for Forecasting

Another source of forecasts are computer programmes and systems for estimating the results of NFL games. These sources offer a different approach towards winner forecasting than the human counterparts. In the case of an expert, where there is a specific level of judgement that is easy to define, the programme has usually a clearly defined algorithm or a number of rules for the selection of the favoured team in the given match. Moreover, the entry data, most commonly in the form of statistic, are fixed and do not depend on the experiences of the expert. The programmes for results forecasting were first introduced in the 80s in relation to the growth of computational machines and PC. Even though, the method for winner estimating in the way of statistic and other data analysis may look at first as more reliable, the history suggests that it is not the case. The main problem is primarily the selection of relevant data. Each of the played American football games offers a wide range of data. Some of the basic statistics of one of the NFL teams follows:

- Number of points
- Quarterback's statistics (CP/AT, YDS, TD, INT)
- Statistics of the players with running yards (ATT, YDS, TD, LG)
- Statistics of players with receiving yards (REC, YDS, TD, LG)
- Statistics of defence players (T-A, SCK, INT, FF)
- Overall number of yards gained through throw/running
- Average number of gained yards in one attempt
- Number of penalties or penalty yards
- Statistics of fumbles
- Statistics of kicks
- Statistics of punts
- Statistics of kick-offs and punt returns
- Effectiveness of offense in the red zone and running of third and fourth's downs
- Duration which the offense spent on the field
- etc.

Regarding the above mentioned examples, it is clear to see that overall statistics of each NFL team offer high number of data. However, not all the data suggest the current form of the given team. This fact is together with the unpredictability of NFL matches the reason for having on average a similar success rate for the winner forecasting as experts. Therefore, it is useful to also follow this source of forecasts.

2.3 Success and Probability

In order to forecast the result from a specific match, it could be useful for a person to gather the estimations from multiple sources. In the case of different predictions for the winner (for a high number of sources this is a common phenomenon) the problem regarding to which source to favour arises. The most reliable guidance could be the

current success of the given source, namely the ratio of correct versus incorrect tips. Using this guidance, it is possible to avoid the unreliable sources (with the success rate of below 50%) and favour the ones which are relatively reliable (with success rate of above 65%). Still, there is a relatively high number of different tips. By simply taking the forecast from the most successful source does not bring the greatest hope for the correct tip for a concrete match. The success of sources differs throughout the season and in a specific round another source which is not the most accurate can have the highest success rate in percentage. However, what should be done when the source with the success rate of 70% suggests the team A to win, but three other teams with the success rate of 60% forecast the team B to become the winner? Furthermore, the forecast itself from one source does not always have 100% probability. Especially the programmes and systems for forecasting often generate even the predicted score of the match (or the point difference between teams) from which it is possible to estimate the probability of the specific team's victory. When a very trustworthy forecast is required, it is necessary to choose reliable sources that give each of them specific weight and realise (calculate) the probability of the team's win.

During the current season with 16 games each week, the manual application of the mentioned process is quite complicated and mainly lengthy. Last but not least, the regular monitoring of various sources is timely and inconvenient. Mobile or desktop applications which could solve the problems and would provide reliable forecasts of the NFL games (according to the possibilities) can bring comfort and save time for the user [9-12].

3 Existing Applications

Mobile applications which are connected to NFL or generally American football are relative common. Few of them serve as an overview of news and played or future games. The rest of the applications look more or less as games and quizzes. Despite of the expectations, there are currently a lot of mobile applications (not considering the platform) which would be aimed at forecasting of NFL matches and at least partially fulfil the mentioned requirements. Relatively popular are the applications in which the user suggests his/hers favoured choices for the winner and by comparison of the successes of his/hers forecasts competes with other friends. Even with a high number of players, it is possible to estimate certain results of matches, but in most cases these fans' forecasts are not very reliable. If the given application is intended for the serious use in estimating the results of games, it often uses only one source of forecasts. Either it is its own system for games forecasting or an application which uses the preferred choices from different experts. Then some of these applications estimate results on the basis of comparing the basic statistics of both teams. As an example, this paper introduces three applications which can be obtained for free.

The first of these is the IronRank: NFL Predictions which is intended for the Android platform. The application for the calculation of the forecasts compares the basic statistics of both teams and according to the given key estimates the probability of the winner [Fig. 1a].

The second application for Android is called Winning Sports Picks, [Fig. 1b]. In this case the basis for this games forecasting is a mathematical model. The application generates the suggestions for the winner even for other sports and also offers a section for the communication with other users.

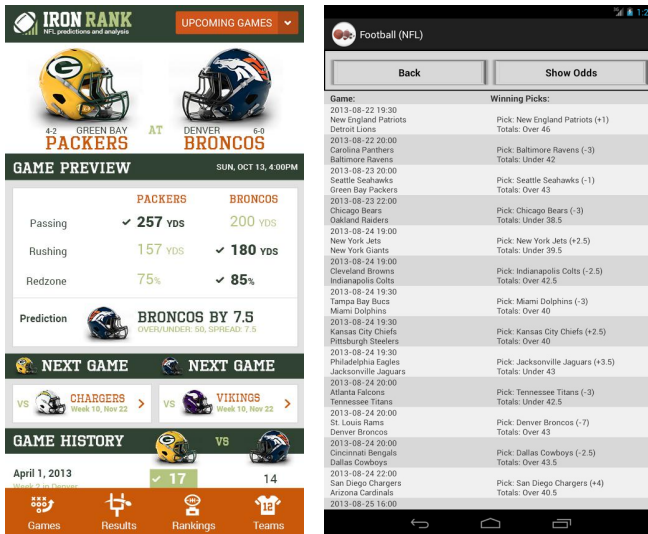


Fig. 1. a, b An example of applications IronRank: NFL Predictions and Winning Sports Picks.

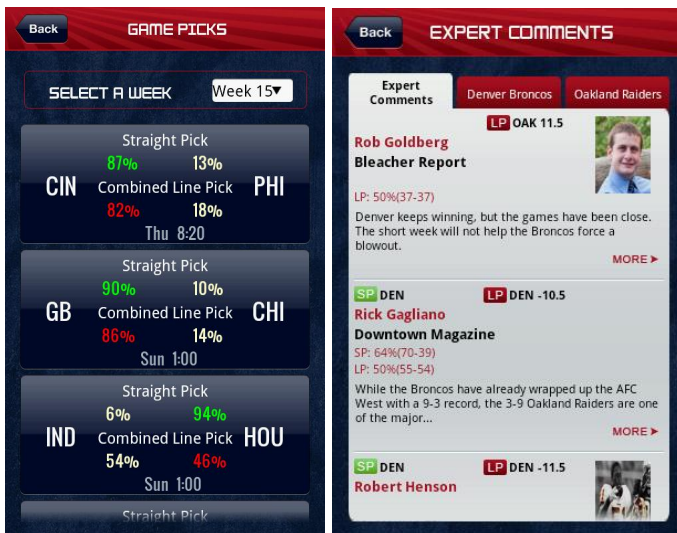


Fig. 2. a, b Example of the application Pick Factor Pro Football

The last application is Pick Factor Pro Football, [Fig. 2a, b]. This application offers forecasts for more sports. Other than NFL games it is for example basketball, hockey, European football. The forecasts are based on the analysis of experts. Apart from the estimations there are also news and commentary of the chosen experts. The application is intended for the Android and iOS platforms.

All of the mentioned applications use only one chosen source. Apart from the rare case where the chosen source is especially reliable, it is useful to compare forecasts from more different sources. Also the smart approaches cannot be found in the mentioned applications. The applications themselves do not take into account for example the success of chose experts.

4 Suggestion for the Mobile Application

4.1 Data Sources

The suggested application can obtain data from various sources, but as the main source the experts' forecasts were chosen. The application can further work with this data. Many web servers offer the forecasts from their own experts. These are mainly former players or current TV commentators. Therefore, it is possible to load data from particular servers. However, the use of servers which collect and freely provide favoured choices of a larger number of experts can be more effective. For example the

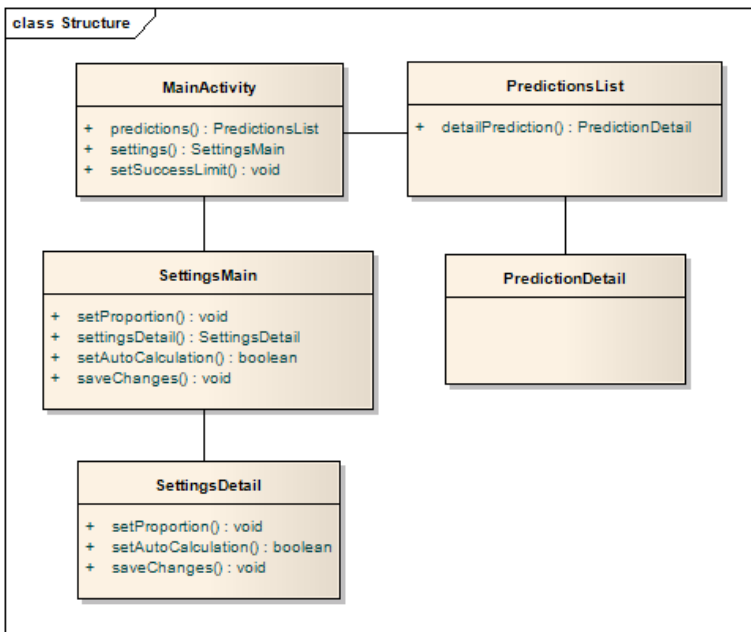


Fig. 3. The structure of the application

NFL Pickwatch server can provide forecasts for the current round of a match from almost 100 NFL experts. The individual successes of given experts are also available. The application would obtain data from the mentioned server which simplifies the implementation.

4.2 The Structure of the Application

The mobile application is developed for the Android platform. The basic element of the application is a list of forecasts for all games of the given NFL round (week). For each game, the application could suggest the winner and the probability of the forecast. The estimations of the experts (sources) who are put into groups (for example according to the TV station) can be shown in detail for each match. The mobile application would analyse individual sources during the season and favour (adjust weights of the forecasts' calculations) specific sources according to their success. The user could also manually adjust the weights. The structure of the application is shown on the [Fig. 3].

4.3 Application's Aim

The main purpose of the application is to save the time of the user, who would otherwise be forced to analyse sources, and to offer clear forecasts of the NFL games according to the success of the individual sources. The advantage of the suggested application is the ability to automatically react to the current success of individual sources and therefore offer the most precise forecasts.

4.4 Possible Problems

Possible problems during the implementation of the application include the building of the algorithm for the calculation of the weights for individual sources according to their success. Furthermore, the calculation of the probability of a win for a specific team which is necessary to be build according to the ratios of particular sources. The aim is to reach the state in which the application prefers sources with the highest success rate in a suitable ratio and, without the interference of the user, reacts (adjusts the ratios) to the change of success rate of individual sources during the season.

Another possible problem is the parsing of data from the chosen web servers. The websites, usually dynamic, often do not provide a well-arranged and valid code. This problem is partially solved by a Jsoup library which is able to work with HTML codes of a worse quality.

5 Implementation of the Mobile Application

5.1 Working with Data

The mobile application was implemented according to the previous draft. The application loads data from web servers. Namely, it is the CBSSports.com from which

all the games for the given week are loaded and NFL Pickwatch which serves as the source of the experts' forecasts. Data parsing is solved using the Jsoup library. Regarding the further work the data are saved into internal SQLite database selected based on [6-8]. The objective-relational mapping and the ORMLite framework is used in the application. The structure of the database is shown in the [Fig. 4].

After the first loading or saving of the changes in the sources' settings, these calculations follow:

- Calculation of the ratio (weight) of each expert in terms of the expert group (according to the success rate)
- Calculation of the success rate for each group(according to the successes and ratio of experts in the group)
- Calculation of the ratio for each group (according to the calculated success rate from the previous result)
- Calculation of the probability and bids for the winner for each group (depending on the bids and ratios of experts)
- Calculation of the probability and bids for winners for the each game (according to the bids and ratios of expert groups)

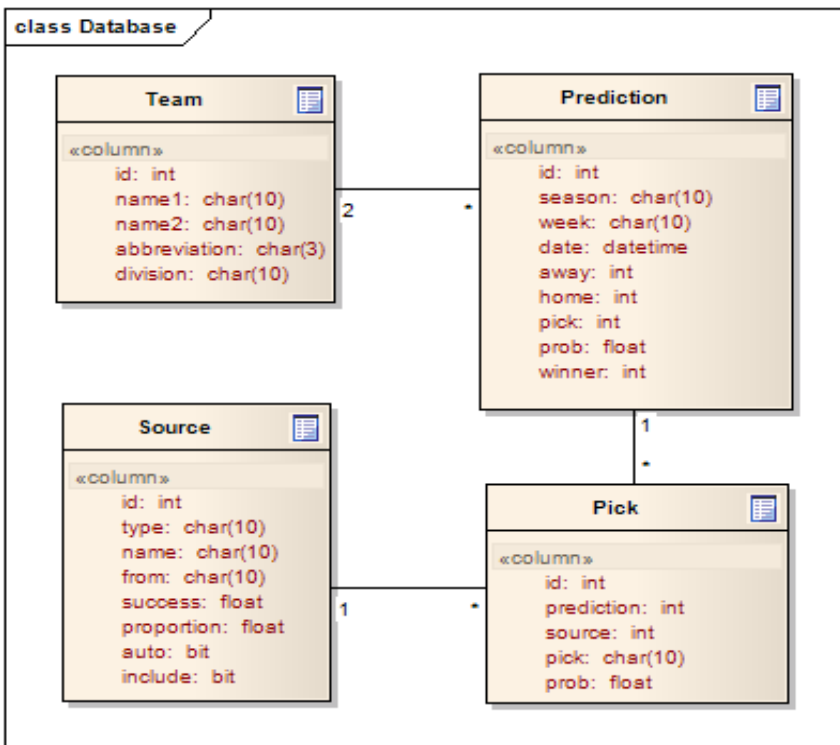


Fig. 4. Structure of the database

5.2 Description of the Application

After the start of the application, the main menu is displayed [Fig. 5a]. The user can see the forecasts for the current week [Fig. 5b]. The forecasts contain the names and logos of competing teams and the probabilities for the win in % (the addition of the probabilities in each game is then always 100%). Regarding the fact that the current season has finished, only the forecasts from the 14th week are shown.

After the selection of the chosen forecast, the detail is displayed in which all the sources can be seen (expert groups) with related information (name, success rate, bids for the winner in the form of the logo of the chosen team and the probability of the forecast). It is possible to display the expert's details for each group (name, success rate, bids for the winner in the form of the logo of the chosen team) [Fig. 6a, b].

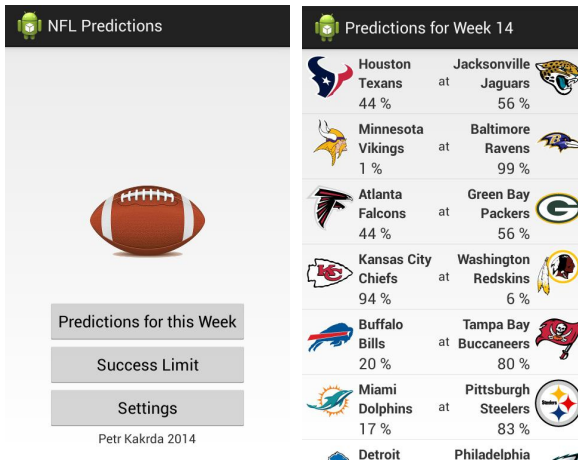


Fig. 5. a, b Example of the new application

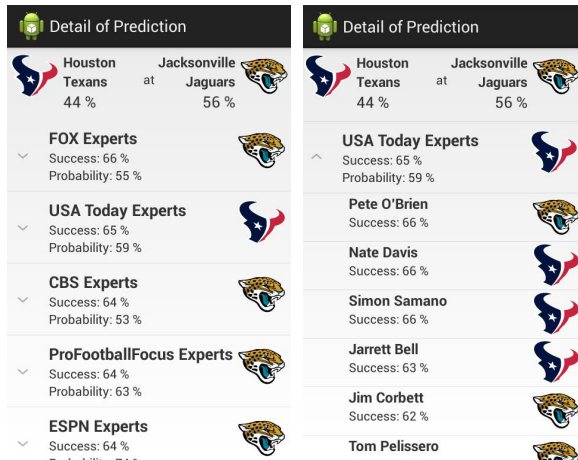


Fig. 6. a, b Example of the new application

Furthermore, the user can set up the minimum border of success rate in % for all the sources (button “Success Limit”). If any of the sources have smaller success rate than the given border, the result is not included in the forecasts.

6 Conclusion

Forecasting of the NFL games with sufficient reliability is not a simple matter and the problems which are raised should be solved. In this paper, it is suggested to implement a mobile application which offers complete forecasts of the NFL games’ results. The advantage of the suggested application is the ability to automatically react to the current success rate of individual sources and therefore, the attempt to provide the most reliable forecasts. The aim of this application is to save user’s time that would otherwise have to be spent on an analysis of the mentioned data sources.

Acknowledgment. This work and the contribution were supported by project “SP/2014 - Smart Solutions for Ubiquitous Computing Environments” Faculty of Informatics and Management, University of Hradec Kralove, Czech Republic.

References

1. Behan, M., Krejcar, O.: Modern Smart Device-Based Concept of Sensoric Networks. *EURASIP Journal on Wireless Communications and Networking* 155(1) (2013)
2. Gantulga, E., Krejcar, O.: Smart Access to Big Data Storage – Android Multi-language Offline Dictionary Application. In: Nguyen, N.-T., Hoang, K., Jędrzejowicz, P. (eds.) *ICCCI 2012, Part I. LNCS*, vol. 7653, pp. 375–384. Springer, Heidelberg (2012)
3. Machacek, Z., Slaby, R., Hercik, R., Koziorek, J.: Advanced system for consumption meters with recognition of video camera signal. *Elektronika ir Elektrotechnika* 18(10), 57–60 (2012)
4. Vanus, J., Novak, T., Koziorek, J., Konecny, J., Hrbac, R.: The proposal model of energy savings of lighting systems in the smart home care. *IFAC Proceedings* 12(pt.1), 411–415 (2013)
5. Vanus, J., Koziorek, J., Hercik, R.: Design of a smart building control with view to the senior citizens’ needs. *IFAC Proceedings Volumes (IFAC-Papers Online)* 12(pt.1), 422–427 (2013)
6. Machaj, J., Brida, P.: Performance comparison of similarity measurements for database correlation localization method. In: Nguyen, N.T., Kim, C.-G., Janiak, A. (eds.) *ACIIDS 2011, Part II. LNCS*, vol. 6592, pp. 452–461. Springer, Heidelberg (2011)
7. Michal, M., Peter, B., Machaj, J.: Modular localization system for intelligent transport. In: Badica, A., Trawinski, B., Nguyen, N.T. (eds.) *Recent Developments in Computational Collective Intelligence. SCI*, vol. 513, pp. 115–124. Springer, Heidelberg (2014)
8. Penhaker, M., Krejcar, O., Kasik, V., Snášel, V.: Cloud Computing Environments for Biomedical Data Services. In: Yin, H., Costa, J.A.F., Barreto, G. (eds.) *IDEAL 2012. LNCS*, vol. 7435, pp. 336–343. Springer, Heidelberg (2012)
9. Behan, M., Krejcar, O.: Smart Communication Adviser for Remote Users. In: Zgrzywa, A., Choroś, K., Siemiński, A. (eds.) *Multimedia and Internet Systems: Theory and Practice. AISC*, vol. 183, pp. 169–178. Springer, Heidelberg (2012)

10. Benikovsky, J., Brida, P., Machaj, J.: Proposal of User Adaptive Modular Localization System for Ubiquitous Positioning. In: Pan, J.-S., Chen, S.-M., Nguyen, N.T. (eds.) ACIIDS 2012, Part II. LNCS, vol. 7197, pp. 391–400. Springer, Heidelberg (2012)
11. Krejcar, O.: Threading Possibilities of Smart Devices Platforms for Future User Adaptive Systems. In: Pan, J.-S., Chen, S.-M., Nguyen, N.T. (eds.) ACIIDS 2012, Part II. LNCS, vol. 7197, pp. 458–467. Springer, Heidelberg (2012)
12. Behan, M., Krejcar, O.: Adaptive Graphical User Interface Solution for Modern User Devices. In: Pan, J.-S., Chen, S.-M., Nguyen, N.T. (eds.) ACIIDS 2012, Part II. LNCS, vol. 7197, pp. 411–420. Springer, Heidelberg (2012)

Security Issues of Mobile Application Using Cloud Computing

Richard Cimler, Jan Matyska, Ladislav Balík, Josef Horalek, and Vladimír Sobeslav

University of Hradec Kralove, Faculty of Informatics and Management,
Department of Information Technologies,
Rokitanskeho 62, Hradec Kralove, 500 03, Czech Republic
{richard.cimler, jan.matyska, ladislav.balik, josef.horalek,
vladimir.sobeslav}@uhk.cz

Abstract. Security issues of the mobile application using cloud computation services are discussed in this paper. Communication between smart phone as a client on the one side and cloud server on the other is described. Security analyse for proposed solution of the health monitoring application is introduced as well. Data about health of the person are one of the most confidential thus need to be secured against different types of threats. Proposed solution is based on the smartphone as a client gathering data and the cloud servers as a computational platform for data storage and analysing. The sensors embedded in the smart phone measure data about monitored person then are partly processed and sent to the server by the internet connection for the deeper analysis.

Keywords: Cloud computing, security, mobile application, health care.

1 Introduction

Current science is characterized by an enormous growth of the ability to measure and gather data in previously unthinkable enormous amounts. Nevertheless, it leads to new requirements for development and application of new methods and ICT capable of storing, transmitting and further analysing data, or potentially capable of using the data as input data for mathematical modelling [1]. Businesses, government agencies, organizations, and individual consumers are rapidly adopting cloud technologies [2]. Businesses are becoming increasingly receptive to the potential for broadening their development, services, and marketing through information technologies [3] [4]. Security of the data is very important part of whole cloud computing concept. Complex view on the issues of cloud computing security can be found at [5] [6]. Both papers are unique by its complexity and described analysis. Several areas of research and dynamic development of the Mobile Cloud computing security are considered. Encryption and key management algorithms, named as Ad hoc Clods in [5] are discussed as well.

As reported by article [5] in its conclusion cloud computing bring various benefits to organizations and users. There are many challenges related to security and privacy

in the Cloud environment. It opens up space for research new techniques for security and privacy in mobile Cloud and ad hoc Cloud. This includes a need for a dynamic security model and better crypto (and key management) algorithms that targets different levels of security and privacy for Cloud computing. With the increasing usage of the Cloud services it is possible to collect sufficient evidence from the cloud providers on the level of trust on each of their services. This can help the service providers, infrastructure providers, and the end-users to better choose the right services from the ever growing Cloud vendors.

There are a lot of benefits of Cloud computing on the one side but a lot of security risk on the other. Many technologies are linked between in the Cloud Computing solutions. Together with its capabilities Cloud "

Computing inherits capabilities of these technologies but its vulnerabilities as well. It is necessary to understand these vulnerabilities to be able to use cloud computing safely.

In by article [7] security issues for cloud models: IaaS, PaaS, and SaaS are presented. Issues vary depending on the model and described storage. Similarly presents solutions for Cloud deployment model and comprehensive paper [8] and [9].

In papers [10] [11] we can find, that Cloud Service as a kind of Web Services is based on Internet service, it faces all kinds of security problems because Internet has many inherent safety defects and also exists in other attacks and threats. Therefore the development of Cloud Service depends on its security deeply, and it is a major significance to consensus on the Cloud Service security.

Security issues of the cloud based application for mobile devices and usage of different frameworks is discussed in many papers, articles and analyses such as [12] [13]. In the paper [13] is described usage of virtualization as a tool for solving security issues of cloud computing using M2M (machine-to-machine) communications. In the M2M computing technologies personal computers, Internet, wireless sensors, and mobile devices are working together. There are many security threats for mobile devices which are the same as for the desktop and server ones. Virtualization technique in M2M communication is described as a way for increasing protection against mobile threats and increase of the performance efficiency.

Modern solution for solving security issues is usage of the frameworks. It also enables to ensure the integrity, secure the data and improves user's identification. Proxy-based multicloud computing framework is introduced at [14]. Several features such as dynamic, on-the-fly collaborations, addressing trust, resource sharing among cloud-based services, privacy issues without pre-established collaboration agreements or standardized interfaces are described in this paper.

Another Secured Mobile-Cloud framework is proposed in [15]. Framework is focused on the security of data transmitted between the components of a mobile cloud application. Two aspects are taken into the account: energy consumptions and users options regarding the security level required for private data. Several distributed components deployed in the cloud or on the mobile device are included in the framework. There is proposed a proof of concept of Android prototype as well.

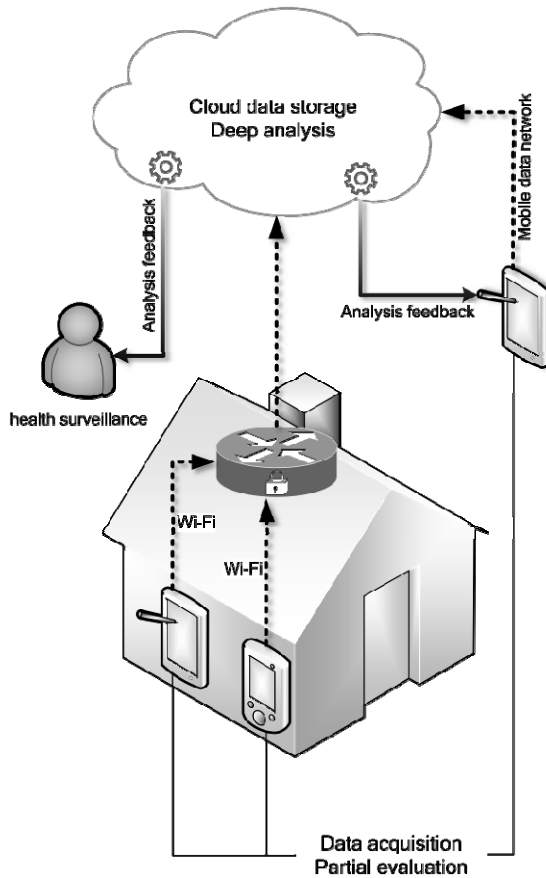


Fig. 1. System scheme

2 Security of Phone and Cloud Communication

In the following part of the paper, security of the communication of specific cloud-based application is described. Application Watch Dog is health care application for monitoring current state of the user by the sensors embedded in the mobile phone. Idea of the whole project is using devices which are people accustomed to carry nearly all day thus they are not constrained to carry some unusual device. Modern smartphone are equipped with several sensors capable of monitoring different physical variables such as gyroscope, accelerometer and thermometer.

Processors and memory of these devices enables to run algorithms for basic evaluation of measured data in order to detect critical situations. Critical situation can be for example fall of the user which can be recognized by the sudden change of the position of the user recorded by the gyroscope. Data from the device are transmitted to the cloud server where deep analysis of the data is executed, see Fig. 1. Results of the

analysis can be sent to the appropriate location. Receivers of the results can be user himself, medical facility or for example personal doctor of the monitored person. There are several communication possibilities of smart phones which enable to monitor person not only indoors but to monitor and transmit data outside as well. Data amount is optimized and usage of mobile communication services such as 3G, LTE or GPRS is taken into the account. Amount of the data should not exceed about 1Gb per month. That enables to use this proposed service even if there is not possibility of Wi-Fi connection.

Security of the communication is essential for whole project. Data about health of a person are one of the most personal and confident. There are a lot of transitions between client and server mostly by the internet connection. Security of the transmitted data in the proposed system is described in the next chapter.

2.1 Communication between Client and Server

Remote data processing requires model specification that will be used for the server to communicate with a client. The basic parameters are data format, frequency of sending and receiving data, volume and a mechanism of confirmation and cataloguing.

Data Transmission Frequency

Given the data character and their use, a frequency of one minute was chosen for the transmission. During this time interval, the client collects data and every minute prepares a data package that is sent to the server. Each one minute interval allows prompt evaluation and sending notifications even for some life-threatening situations.

The system allows sending the data back in case that there was some problem and it was not possible to send the data. The reason can be for example, the lack of data connection. In this case the client creates a data package containing the data from the last confirmed synchronization, but not older than from the last 24 hours.

Format and Volume of the Sent Data

The client sends raw data measured on individual sensors. For this type of data, there is a simple convenient and easily processable structure of ASCII comma-separated values with 300 lines (local data collection in the interval of 200ms). The volume of the data is up to 0.5kB.

Confirmation of Reception and Cataloguing

The client's database stores only a limited amount of data because the old data is always replaced with new data on the basis of round-robin model. This way, the client's database should not exceed the size of 36MB and should allow using the application even on less equipped devices. This is at the same time the maximum possible volume of the data sent and that makes the application easy to use not only with the constant Wi-Fi connection but also with a connection provided by a mobile network operator (GPRS, CDMA, 3G, LTE).

The individual lines of sent data are equipped with a time stamp that serves as a unique key in the server database. The long-term data are stored in raw and also in a converted form.

A check sum of the data package is verified during data transmission so that it would be possible to detect any damage or any possible change in the process. During the communication, the server confirms reception of the individual packages and it can also request completion of some lines from the last 24 hours.

2.2 Security Model of Communication between Client and Server

There is a duplex communication between server and client that always runs on the lines that can be tapped in some way, whether it is using the Wi-Fi network or the Internet. Therefore, it is necessary to encrypt the data and that way ensure their safety. For that reason, the application counts on using [16] AAA security model, Authentication, Authorization and Accounting. Furthermore, the application relies on the security on several layers of ISO/OSI model.

Authentication, Authorization and Accounting

The model AAA provides all of the basic needs of users' administration and because it is a case of a very widespread standard, it allows connecting the application with other already existing services. Authentication, in this case, provides user authentication using a username and a password. With these, the user can access both, the web frontend and mobile applications (where the password is saved).

Authorization covers the access to one's own data and also the data of other users. Authorization includes used mobile devices for the collection of user data. Accounting serves then for billing purposes and for access logging to identify safety incidents.

Verification of the Used Mobile Device

A part of authorization is verification of a mobile device that is used by the user. To secure a stronger security, it is not possible to use for the communication with the server just any mobile device with the installed application using a username and a password. In order to communicate with the application, the server authenticates also IMEI devices that are registered to the user's account. This security model is similar to the one that some banks use for their smartphone banking.

Transmission of Recorded and Processed Data

HTTP protocol is used for the transmission of recorded data [16]. It features a complete model of the security based on certificate use and it includes mechanisms for safe key exchange, symmetric data encrypting and simplex hashing. The server, as well as the application, will support encrypting using SSL 3.0, or alternatively TLS.

Encryption of Transmitted Data

Since the application wants to offer a truly high degree of data security, the data are encrypted even during the transmission on the level of the application itself. It is considered that the most convenient is the use of symmetric encryption with the help of AES protocol with the key length of 256-bit. The encryption using this protocol is directly implemented in the programming language Java and therefore its use is simple and uncomplicated for the tools of mobile devices.

3 Advantages and Disadvantages of Public Cloud Server

Modern mobile applications that require efficient platform for recorded data processing very often use remote data processing by means of servers situated in a local network or in the Internet. These servers then either offer user interface with higher functionality or they only process the data and send them back.

If the concern is an application that needs a remote data processing in order to work and it can make use of the web access as a platform for visualization of those processed data, it can be thought about a backend platform in form of Cloud solution. With Cloud solution there is a choice of its form. From the point of view of offered services, there can be distinguished three basic distribution models. See Fig.2.

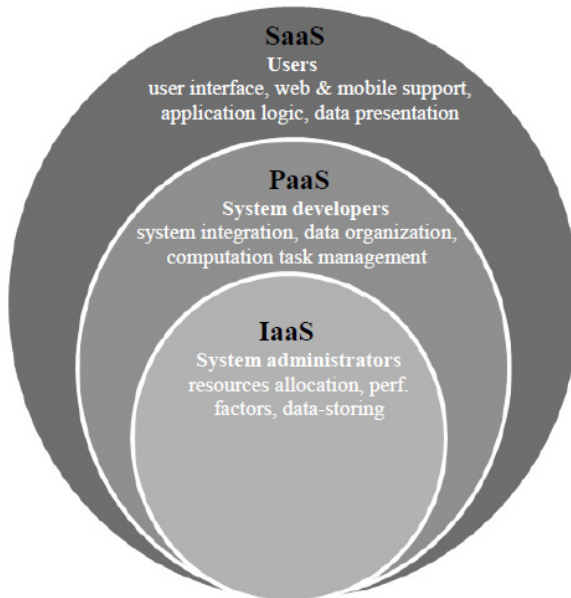


Fig. 2. Cloud computing distribution models

3.1 Watchdog Cloud Computing Solutions

Watchdog application is able to run in three majorly different solutions. Two of them are purely cloud based, and the third is a small environment oriented.

Public Cloud Run as Service

With Watchdog application, Cloud solution for processing and presenting data (in case of public Cloud) has a form of a software run as a service. User has to log into that application and the application will offer three basic services. First, there is a

performance for data processing and processing itself. Second, there is a presentation of processed data. Finally, there is a storage where the processed data are saved.

The advantage of this solution is the already mentioned simplicity of use. The user does not have to care about the infrastructure nor the application. The user only uses the services in accordance with the SaaS model. In contrast, there are all disadvantages of public Clouds. The most important is undoubtedly the safety. Others are limited data control, impossibility to intervene in the form of application, etc.

Private Cloud for Gathered Data Processing

In case of running a Cloud application, we can think about the usual architecture Software as a Service as mentioned above, or about the possibilities of a private Cloud.

Considering the magnitude of the whole solution, it can be spoken about private Cloud only in connection with the organization that operates this service for their clients. If it was a case of one user of this application who runs server application at home for own purposes, we could not be talking about a Cloud solution. In contrast to that, we can imagine a situation where a nursing home with 500 clients runs this service as a private Cloud. Here, a necessary separation of the hardware infrastructure from the very service is put to work. IT staff takes care of the Cloud and the service is consequently used by care assistants and clients.

In case of the private Cloud the risks of security stay but the operator has an absolute control over the data and the form of the application (Open-source solution is expected) with all its positive and negative results (care of the Cloud, data backup, the risk of data mishandling).

Server Run for Independent User

In the case of the situation described above, where one user wants to operate the service only for own use, we cannot talk about Cloud solution but this form of operation is naturally also possible. However, there is a range of limitations and complications that make this solution, in our opinion, the least recommendable. Mobile application itself can evaluate some basic occurrences but primarily it is designed for constant communication with the server. Therefore, it is necessary to arrange either a constant access to the server through Wi-Fi or make the server accessible from the Internet. For domestic use, there are difficulties arising with public addresses, possible use of dynamic DNS records, and also the server security when we cannot assume that an ordinary user will be at the same time a specialist in the server security area.

4 Safety in Service Run in Cloud Environment

We can classify user data concerning the state of user's health among one of the most personal data that the user has. For that reason, we have to take into consideration the safety of saved data and pay attention to the risks of individual solutions. [18] Safety risk of public Clouds might be their most mentioned weakness.

[19] mentions Cloud computing risk assessment where it refers to problems using usecase. Those problems are in most of the cases not connected to a technical solution. A range of safety risks connected with public Clouds run is not solely technical. Problems connected with cooperation with another subject play a big part here. For example, it is possible to use the situation when the client uses services of supplying company for SaaS. Even if the client verifies the company to find out if it meets the technical requirements, it is reliable, etc., unexpected complications can occur. The company can go bankrupt, it can be merged with another company or bought by another company. Thereafter, our data go to someone else and we cannot influence that. These non-technically orientated issues are the subject of risk management and even though we have to include their risks, they have no direct connection with the operation solution from the technical perspective.

The main technical risk related to public Cloud is a loss of isolation. The isolation is for the run of public Cloud solution absolutely crucial. If the clients run their own service in public Cloud, their operator is obliged to separate their data and processors from other clients even though they share the same hardware. That way the physical disks, processors, RAM memories and network connections are shared and their separation happens in logical layers, in software. The virtualization safety of data storages is very closely connected to that. The user, however, does not have an influence on this operator's environment. The effective protection of application run in Cloud environment resides for the user or operator in a careful selection of the provider and in case that the provider cannot be trusted, it is good to opt for running an own private Cloud.

We have to consider that the run of a private Cloud or a home server does not automatically guarantee a higher safety. This presumption would be possible only in case where we would consider that the operator of the private cloud or home server has unlimited tools or skills for their protection. We would recommend the client to use the services of public Cloud and private Cloud in case that it is convenient to invest resources to building an administration of such solution.

4.1 Gartner's Seven Cloud-Computing Security Risks

When a customer chooses a cloud provider, there are seven general security issues concerning Cloud computing described by Gartner [20]. They include issues like privileged user access which addresses risk of confidentiality disruption. This issue is connected either with data transfers and Cloud service provider. Both are addressed by an application security model and a careful selection of a provider. Another issue concerning Watchdog application is a regulatory compliance. It discusses and addresses the data responsibilities. Customer should choose only those providers, who are responsible enough to allow 3rd party organizations to verify their security level. By the possibility of choice, in the end, the customer is responsible for his own data security. There are several law related issues described by Gartner, which are data location and investigative support. When a customer utilizes a Cloud services, he often doesn't know, where his data are located or even in what country. This information is in fact of a vital importance, because of a different law environment.

Investigative support is an issue related to possibilities to investigate security concerns. These possibilities can be completely restricted by contract, or by law. As described above, data segregation on a common hardware is a key principle of Cloud computing. Effective utilization of encryption should be implemented in the Cloud architecture. Furthermore provider should have a reliable encryption model, because encryption model failure could cause an irreversible corruption of data. Another issue connected to a corruption of data is a recovery protocol implementation. Every Cloud provider should have a robust recovery protocol and be able to provide complete restoration in case of a disaster. Especially single-site infrastructure is vulnerable to various kinds of disasters. Last issue discussed by Gartner is one mentioned earlier in this article. You can never be sure, that your provider won't go bankrupt or be swallowed by another company. Service contract should remember these situations and specify how and under what terms your data will be available in these situations.

These risks address most common issues related with cloud-computing oriented solutions. It can be used as a reference for Watchdog implementation same as for general use with any cloud-computing oriented project.

5 Conclusions

This paper presents the security issues of a mobile application using cloud computing. Data about health status of any person are one of the most personal data and has to be secured thoroughly. The application which is being developed collects data by using the internal sensors of a smart phone. Data are partly evaluated in the device in order to evaluate critical situation and sent to the server for the deeper analysis. There are several layers of security. The application is using AAA security model, the server authenticates IMEI of the device and encrypting using SSL 3.0 as well. Data are also encrypted on the application level by the symmetric encryption by means of the AES protocol with the key length of 256-bit.

Three possible cloud computing solution for the proposed application has been discussed: Public Cloud run as service, Private Cloud for gathered data processing and Server run for independent user. Only first two of them can be called pure cloud solutions but all of them can be used as a secured server for the application.

Acknowledgment. This work was supported by the project of specific research-Project.No. 3/2014/2101. Faculty of Informatics and Management, University of Hradec Kralove.

References

- [1] Hřebíček, J., et al.: Scientific computing in mathematical biology, MU (2012), <http://www.iba.muni.cz/res/file/ucebnice/hrebicek-vedecke-vypocty.pdf>

- [2] Bureš, V., Otčenášková, T., Čech, P., Antoš, K.: A Proposal for a Computer-Based Framework of Support for Public Health in the Management of Biological Incidents: the Czech Republic Experience. *Perspectives in Public Health* 132(6), 292–298 (2012), doi:10.1177/1757913912444260, ISSN 1757-9139
- [3] Allan, R.: Cloud and Web 2.0 resources for supporting research (2012), <http://tyne.dl.ac.uk/NWGrid/Clouds/>
- [4] Bureš, V., Brunet-Thornton, R.: Knowledge Management: The Czech Situation, Possible Solutions and the Necessity for Further Research. In: *Proceedings of the 6th International Conference on Intellectual Capital and Knowledge Management*, McGill University, Montréal, Canada, pp. 95–102 (2009) ISBN 978-1-906638-45-0
- [5] Chirag, M., Dhiren, P., Bhavesh, B., Avi, P., Muttukrishnan, R.: A survey on security issues and solutions at different layers of Cloud computing. *The Journal of Supercomputing* 63(2), 561–592 (2013), doi:10.1007/s11227-012-0831-5, ISSN 0920-8542
- [6] Subashini, S., Kavitha, V.: A survey on security issues in service delivery models of cloud computing. *Journal of Network and Computer Applications* 34(1), 1–11 (2011), <http://dx.doi.org/10.1016/j.jnca.2010.07.006>, ISSN 1084-8045
- [7] Keiko, H., Rosado, D.G., Fernández-Medina, E., Fernandez, E.B.: An analysis of security issues for cloud computing. *Journal of Internet Services and Applications* 4(5) (2013), doi:10.1186/1869-0238-4-5, ISSN 1867-4828
- [8] Fernandes, D.A.B., Soares, L.F.B., Gomes, J.V., Freire, M.M., Inácio, P.R.M.: Security issues in cloud environments: a survey. *International Journal of Information Security* (2013), doi:10.1007/s10207-013-0208-7, ISSN 1615-5262
- [9] Lee, H., Kim, J., Lee, Y., Won, D.: Security Issues and Threats According to the Attribute of Cloud Computing. In: Kim, T.-H., Stoica, A., Fang, W.-C., Vasilakos, T., Villalba, J.G., Arnett, K.P., Khan, M.K., Kang, B.-H. (eds.) *SecTech, CA, CES3 2012*. CCIS, vol. 339, pp. 101–108. Springer, Heidelberg (2012)
- [10] Jia, W., Sun, S.: Research on the Security Issues of Cloud Computing. In: Du, Z. (ed.) *Intelligence Computation and Evolutionary Computation*. AISC, vol. 180, pp. 845–848. Springer, Heidelberg (2013)
- [11] Mouratidis, H., Islam, S., Kalloniatis, C., Gritzalis, S.: A framework to support selection of cloud providers based on security and privacy requirements. *Journal of Systems and Software* 86(9), 2276–2293 (2013), <http://dx.doi.org/10.1016/j.jss.2013.03.011>, ISSN 0164-1212
- [12] Sujithra, M., Padmavathi, G.: Mobile device security: A survey on mobile device threats, vulnerabilities and their defensive mechanism. *International Journal of Computer Applications* 56(14) (2012), doi:<http://dx.doi.org/10.5120/8960-3163>, ISSN: 09758887
- [13] Cagalaban, G., Kim, S., Kim, M.: A mobile device-based virtualization technique for M2M communication in cloud computing security. In: Kim, T.-H., Stoica, A., Fang, W.-C., Vasilakos, T., Villalba, J.G., Arnett, K.P., Khan, M.K., Kang, B.-H. (eds.) *SecTech, CA, CES3 2012*. CCIS, vol. 339, pp. 160–167. Springer, Heidelberg (2012)
- [14] Singhal, M., Chandrasekhar, S., Ge, S.R., Krishnan, R., Ahn, G.J., Bertino, E.: Collaboration in Multicloud Computing Environments: Framework and Security Issues. *Computer* 46(2), 76–84 (2013) ISSN: 0018-9162, WOS:000314943300019
- [15] Popa, D., Boudaoud, K., Borda, M.: Secure mobile-cloud framework - implementation on the mobile device. *Acta Technica Napocensis* 54(4), 7–12 (2013) ISSN: 12216542
- [16] Wood, J., Aboba, B.: RFC 3539 - Authentication, Authorization and Accounting (AAA) Transport Profile (2003), <http://tools.ietf.org/html/rfc3539>

- [17] Singhal, M., Chandrasekhar, S., Ge, S.R., Krishnan, R., Ahn, G.J., Bertino, E.: Collaboration in Multicloud Computing Environments: Framework and Security Issues. *Computer* 46(2), 76–84 (2013) ISSN: 0018-9162, WOS:000314943300019
- [18] Gejibo, S., Mancini, F., Mughal, K.A., Valvik, R., Klungsøyr, J.: Challenges in Implementing an End-to-End Secure Protocol for Java ME-Based Mobile Data Collection in Low-Budget Settings. In: Barthe, G., Livshits, B., Scandariato, R., et al. (eds.) *ESSoS 2012*. LNCS, vol. 7159, pp. 38–45. Springer, Heidelberg (2012)
- [19] European Union Agency for Network and Information Security, *Cloud Computing Risk Assessment — ENISA* (2009), <http://www.enisa.europa.eu/activities/risk-management/files/deliverables/cloud-computing-risk-assessment>
- [20] Brodtkin, J.: Gartner: Seven cloud-computing security risks. *Infoworld*, 1–3 (2008)

SIFT-Based Arabic Sign Language Recognition System

Alaa Tharwat^{1,4}, Tarek Gaber^{2,4}, Aboul Ella Hassanien^{3,4},
M.K. Shahin¹, and Basma Refaat¹

¹ Faculty of Eng. Suez Canal University, Ismailia, Egypt

² Faculty of Computers & Informatics , Suez Canal University, Ismailia, Egypt

³ Faculty of Computers and Information, Cairo University, Egypt

⁴ Scientific Research Group in Egypt (SRGE)

egyptscience.net

Abstract. The literature contains many proposed solutions for automatic sign language recognition. However, the ArSL (Arabic Sign Language), unlike ASL (American Sign Language), did not take much attention from the research community. In this paper, we propose a new system which does not require a deaf wear inconvenient devices like gloves to simplify the process of hand recognition. The system is based on gesture extracted from 2D images. The Scale Invariant Features Transform (SIFT) technique is used to achieve this task as it extracts invariant features which are robust to rotation and occlusion. Also, the Linear Discriminant Analysis (LDA) technique is used to solve dimensionality problem of the extracted feature vectors and to increase the separability between classes, thus increasing the accuracy of the introduced system. The classifiers, Support Vector Machine (SVM), k-Nearest Neighbor (k-NN), and minimum distance will be used to identify the Arabic sign characters. Experiments are conducted to check the performance of the proposed system and it showed that the accuracy of the obtained results is around 99%. Also, the experiments proved that the proposed system is robust against any rotation and they achieved an identification rate near to 99%. Moreover, the evaluation shown that the system is comparable to the related work.

Keywords: SIFT, Arabic Sign Language Recognition, SVM, k-NN, LDA, classification.

1 Introduction

For thousands of years, sign languages are the default communication languages between deaf people. These languages have been **used** to successfully teach generations of deaf children. However, a communication gap between deaf and non-deaf people is obvious. This is because the normal people find it difficult to learn and comprehend sign languages [1–3].

A sign language is a collection of gestures, movements, postures, and facial expressions corresponding to letters and words in natural languages. So, there

should be a way for the non-deaf people to recognize the deaf language (i.e. sign language). Such process is known as a *sign language recognition*. The aim of the sign language recognition is to provide an accurate and convenient mechanism to transcribe sign gestures into meaningful text or speech so that communication between deaf and hearing society can easily be made. To achieve this aim, many proposal attempts are designed to make fully automated systems or Human Computer Interaction (HCI) to facilitate interaction between deaf and non-deaf individuals [1, 4].

The sign language recognition is mainly based on gesture recognition. There are two main categories for gesture recognition glove-based systems and vision-based systems.

- **Glove-based systems:** In these systems, electromechanical devices are used to collect data about deaf's gestures. With this systems, the deaf person should wear a wired glove connected to a number of sensors to collect the gestures of the person's hand. So, such gestures can be recognized through a computer interface. This way gives a good result but it is inconvenient because the user must always carry wired sensors (gloves) and this is not natural way to communicate between deaf and non-deaf people [4].
- **Vision-based systems:** These systems make image processing and machine learning techniques to identify, recognize and interpret hand gestures. Such system can overcome the inconvenience problem of gloved-based systems as there is no need for the deaf users to wear any electromechanical devices. In other words, vision-based systems are more flexible to use [4].

As sign language varies from country to country (in some case, it even varies according to the regions), our focus in this paper will be on the Arabic Sign Language (ArSL). Many efforts [1, 4–8] have been made to establish the sign language used in Arab countries by trying to standardize the language and spread it among members of the deaf community and those concerned [1].

K. Assaleh et al. [1] proposed automatic ArSL recognition system which uses the polynomial networks as a classification engine. To evaluate the performance of their system, they have used real ArSL data collected from deaf people. Their system was based on glove based system. They used Adaptive Neuro Fuzzy Inference System (ANFIS) in recognition and they achieved accuracy near to 93.41%.

Al-Jarrah et al. [6], have presented two neural network systems for the recognition of the sign language alphabets. They used vision based system; hence they deal with signs as images. They extracted the orientation and location of the hand within the image; then they used two different classifiers, namely Feed Forward Neural Network (FFNN) and Probabilistic Neural Networks (PNN). Their system achieved 94.4% when they used FFNN classifier, while PNN achieved 91.3%. In another experiment [9], the images of bar hands are processed and robust features against to rotation, translation and scaling were extracted and they achieved 97.5%.

Nashwa El-Bendary et al. [4], proposed a system to translate Arabic sign language into text based on vision based system. They used three features to

represent the position of the hand. Then, they used the direction of the hand wrist to categorize the Arabic Alphabets into three classes. To extract features, they first detect the orientation point, and then about 50 points equally spaced by a certain angle are used as features. Minimum distance and Multi-layer Perceptron (MLP) classifiers were used and achieved 91.37% and 83.7%, respectively.

Aliaa et al. [10], designed as ArSL based on Hidden Markov Model (HMM). They collected a large dataset to recognize 20 isolated words from real videos taken for deaf people in different clothes and skin colors and they achieved recognition rate near to 82.22%.

Nadia et al. in [7] have designed an ArSL to identify the sign character from high resolution video. Their system consists of two stages. In the first stage, their system detects a hand movement in each frame to decide the final region of processing (i.e. Region of Interest). Then, a motion detection technique was applied to detect the time or frame of recognition. In the second stage, a sign language recognition was implemented based on Fourier descriptors feature extraction method and k-NN classifier. Their system has achieved an accuracy rate of 90.55%.

In [11, 12], Mohandes used Hu's moments as features and SVM classifier and achieved 87%.

In this paper, the proposed system based on extracting features that robust against to rotation, scaling, and shifting using Scale Invariant Feature Transform (SIFT). In SIFT algorithm, the dimension of feature vector is not fixed, but it depends on size and contents of the image. Thus, many images have feature vectors with high dimensions. This leads to high dimensionality problem. LDA is used to solve this problem and improve the accuracy of the system. Finally, three different classifiers are used namely, SVM, k-NN, and Nearest Neighbor. The proposed system is a part of our framework in [13].

The remainder of this paper is organized as follows. Section (2) gives an overview of the SIFT, LDA, and classifiers SVM, k-NN, and nearest neighbor used in our proposed system which is presented in Section (3). Section 4 introduces the implementation phase, the experiment scenarios, and the discussion. Conclusions and future works are reported in Section (5).

2 Preliminaries

2.1 Scale Invariant Features Transform (SIFT)

The SIFT algorithm takes an image and transforms it into a collection of local feature vectors. Each of these feature vectors is supposed to be distinctive and invariant to any scaling, rotation or translation of the image. SIFT algorithm consists of the following steps:

- **Creating the Difference of Gaussian Pyramid (Scale-Space Peak Selection):** The first stage is to construct a Gaussian, "scale space", function from the input image. This is formed by convolution (filtering) of the original image with Gaussian functions of varying scales. The difference of

Gaussian (DoG), $D(x, y, \sigma) = L(x, y, k\sigma) - L(x, y, \sigma)$, where $L(x, y, \sigma)$ and $L(x, y, k\sigma)$, are two images that produced from the convolution of Gaussian functions with an input image $I(x, y)$ with σ and $k\sigma$ respectively, and $G(x, y, \sigma) = \frac{1}{2\pi\sigma^2} \exp[-\frac{x^2+y^2}{\sigma^2}]$ represents Gaussian function [14].

- **Extrema Detection:** In this step, keypoints are detected. To do so, the local maximum and minimum of $D(x, y, \sigma)$ is computed by compared the pixel with the pixels of all its 26 neighbors.
- **Unreliable Keypoints Elimination:** This stage attempts to eliminate some points from the candidate list of keypoints by finding those that have low contrast (sensitive to noise) or are poorly localised on an edge.
- **Orientation Assignment:** This step aims to assign one or more orientation to the keypoints based on local image properties. The histogram is formed from $m(x, y) = \sqrt{(L(x+1, y) - L(x-1, y))^2 + (L(x, y+1) - L(x, y-1))^2}$ and $\theta(x, y) = \arctan(L(x, y+1) - L(x, y-1)) / (L(x+1, y) - L(x-1, y))$ which represent gradient and orientation of sample points within a region around the keypoint [14–16].

Each sample is weighted by its gradient magnitude and by a Gaussian-weighted circular window with a σ that is 1.5 times that of the scale of the keypoint.

We locate the highest peak in the histogram and use this peak and any other local peak to create a keypoint with that orientation. Some points will be assigned multiple orientations if there are multiple peaks of similar magnitude. The assigned orientation, location and scale for each keypoint enables SIFT features to be robust to rotation, scale and translation

- **Descriptor Computation:** In this stage, the goal is to create descriptive for the patch that is compact, highly distinctive and to be robust to changes in illumination and camera viewpoint. The image gradient magnitudes and orientations are sampled around the keypoint location. These values are illustrated with small arrows at each sample location as in Fig (1). In order to achieve orientation invariance, the coordinates of the descriptor and the gradient orientations are rotated relative to the keypoint orientation. In our implementation, a 16×16 sample array is computed and a histogram with 8 bins is used

The Number of features depends on image content, size, and choice of various parameters such as patch size, number of angles and bins, and peak threshold. These parameters will be briefly described below.

Peak Threshold (*PeakThr*) parameter is used to determine the dimension of feature vectors because *PeakThr* represents the amount of contrast to extract a keypoint. The optimum value of *PeakThr* is 0.0 because when the value of *PeakThr* parameter increased, the number of features decreased and more keypoints are eliminated; thus, the robustness of feature matching will decreased [4].

The patch size (*Psize*) parameter is used to extract different fine grained of features. Increasing the size of patches will decrease the dimension of feature vector, but at the same time it increases CPU time. The feature vector at each keypoint is *bins* \times *bins* \times *angles*. The problem is to experiment many patch sizes

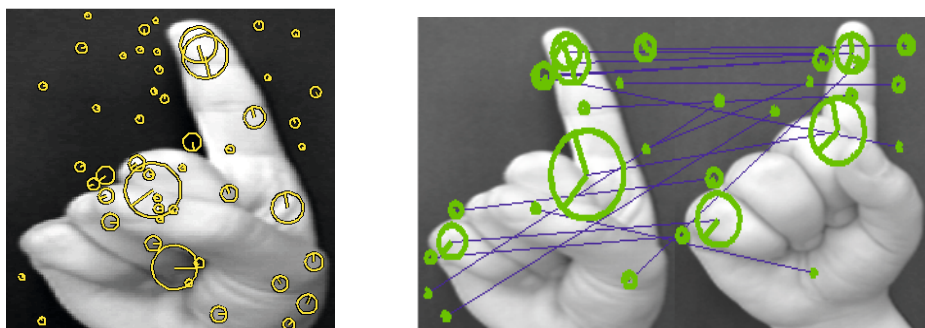


Fig. 1. Keypoints and matching between two different characters of Arabic sign language images, (a) Keypoints or Extrema of SIFT, (b) Matching between two images based on SIFT features

to reach to the $Psize$ that gives better accuracy. As reported in [4], increasing $Psize$ of SIFT will produce global features. Also, decreasing $Psize$ will not extract enough features for identification.

Using different number of angles ($Nangels$) and number of bins ($Nbins$) will collect features in different orientations. Increasing the $Nangels$ increases the number of features hence improve the accuracy of the system. on the other hand decreasing $Nangels$ and $Nbins$ leads to small number of features and the identification rate will decreased specially when the images are rotated because the features became variant against any rotation.

2.2 Linear Discriminant Analysis (LDA)

LDA is one of the most famous dimensionality reduction method used in machine learning. LDA attempts to find a linear combination of features which separate two or more classes [17].

The goal of LDA is to find a matrix $W = \max \frac{S_b}{S_w}$ that maximizing Fisher's formula. $S_w = \sum_{j=1}^c \sum_{i=1}^{N_j} (x_i^j - \mu_j)(x_i^j - \mu_j)^T$ represents a within-class scatter matrix, where x_i^j is the i^{th} sample of class j , μ_j is the mean of class j , c is the number of classes, and N_j is the number of samples in class j . $S_b = \sum_{j=1}^c (\mu_j - \mu)(\mu_j - \mu)^T$ is a between-class scatter matrix, where μ represents the mean of all classes. The solution of Fisher's formula is a set of eigen vectors (V) and eigen values (λ) of the fisher's formula.

2.3 Classifiers

In this paper, we have applied three classifiers to assess their performance with our approach. An overview of these classifiers is given below.

Support Vector Machine (SVM). SVM is one of the classifiers which deals with a problem of high dimensional datasets and gives very good results. SVM tries to find out an optimal hyperplane separating 2-classes basing on training cases [18].

Given a training dataset, $\{x_i, y_i\}$, where $i = 1, 2, 3, \dots, N$, where N is the number of training samples, x_i is a features vector, and $y_i \in \{-1, +1\}$ is the target label, $y = +1$, for samples belong to class C_1 and $y = -1$ denotes to samples belong to class C_2 . Classes C_1 and C_2 are linearly separable classes. Geometrically, the SVM modeling algorithm finds an optimal hyperplane with the maximal margin to separate two classes, which requires solving the optimization problem in equation 1.

$$\begin{aligned} \text{maximize } & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j (x_i, x_j) \\ \text{subject to: } & \sum_{i=1}^n \alpha_i y_i, 0 \leq \alpha_i \leq C \end{aligned} \quad (1)$$

where, α_i is the weight assigned to the training sample x_i (if $\alpha_i > 0$, then x_i is called a support vector); C is a regulation parameter used to find a trade-off between the training accuracy and the model complexity so that a superior generalization capability can be achieved; and K is a kernel function, which is used to measure the similarity between two samples.

K-Nearest Neighbor (k-NN). In k-NN classifier, unknown patterns are distinguished based on the similarity to known samples by computing the distances from unknown patterns to every sample and select the K-nearest samples as the base for classification. The unknown pattern is assigned to the class containing the most samples among the K-nearest samples [19].

Nearest-Neighbor Classifier. The nearest neighbor or minimum distance classifier is one of the oldest known classifiers. Its idea is extremely simple as it does not require learning. Despite its simplicity, nearest neighbors has been successful in a large number of classification and regression problems. To classify an object I , first one needs to find its closest neighbor X_i among all the training objects X and then assigns to unknown object the label Y_i of X_i . Nearest neighbor classifier, works very well in low dimensions. The distance can, in general, be any metric measure: standard Euclidean distance is the most common choice.

3 The Proposed ArSL Recognition System

As shown in Fig (2), all training and testing images are collected and features are then extracted by applying SIFT algorithm which is described above. The output of the SIFT algorithm is feature vectors with high dimension for each image. Processing these high dimension vectors takes more CPU time. So, we need

to reduce the dimensions of the features by applying dimensionality reduction method, such as LDA which is one of the most suitable methods to reduce the dimensions and increase the separation between different classes. LDA also decreases the distance between the objects belong to the same class. The feature extraction and reduction are performed in both training and testing phases. The matching or the classification of the new images are only done in the testing phase.

The proposed system, as illustrated in Fig (2), is mainly consists of two phase: Training and Testing.

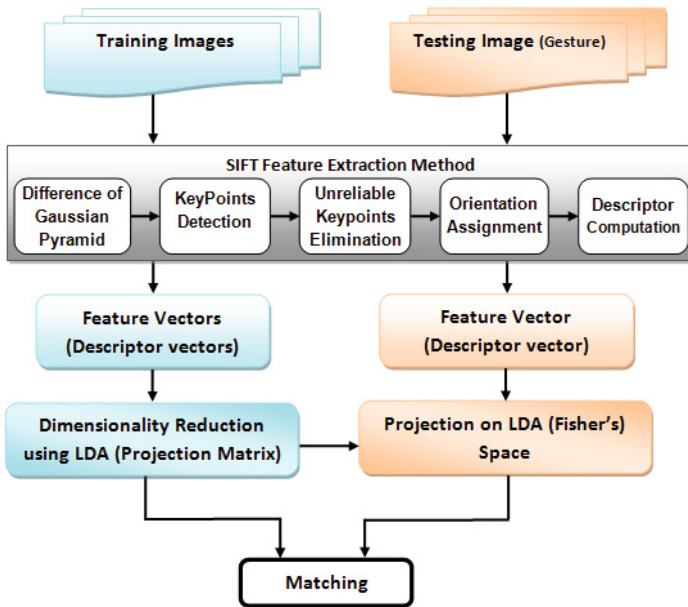


Fig. 2. Block diagram of the proposed system

- **Training Phase:** The detailed steps of training phase of proposed algorithm can be summarized as follows.
 1. Collecting all training images (i.e. gestures of Arabic Sign Language).
 2. Extracting features using SIFT from each images.
 3. Representing each image by one feature vector. To reduce the number features; we used LDA as a dimensionality reduction method.
- **Testing Phase:** The detailed steps of testing phase of proposed algorithm can be summarized as follows.
 1. Collecting the testing image,
 2. Extract the features of testing images and convert them to feature vectors.

3. Feature vectors are projected on LDA space to reduce their dimensionality.
4. Matching or classifying the test feature vector to identify final decision (i.e. identify the character corresponding to this image (sign)).

4 Experimental Results

To evaluate our proposed system, we have used Matlab platform under windows 32-bit operating system to implement it and run some **experiments**. The experiments have been conducted using a PC with the following specifications: Intel(R) Core(TM) i5-2400 CPU @ 3.10 GHz, and 4.00 GB RAM.

The dataset used in our experiments is an ArSL database images which were collected in Suez Canal University. The database consisting of 210 gray level ArSL images with size 200×200 . These images represent 30 Arabic characters (7 images for each character). The images are collected in different illumination, rotation, quality levels, and image partiality. All images are centered and cropped to 200×200 . Examples of these images are shown in Fig (3) in which two different samples of the first two characters (e.g. Alef and Baa) are shown. These samples (gestures) are collected from various individuals who have different sizes of hand. However, this does not matter as our method depends on SIFT features.

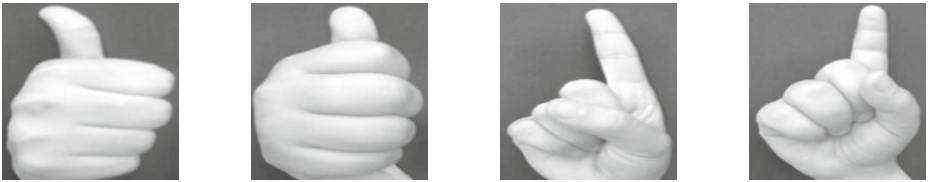


Fig. 3. A sample of collected ArSL gestures representing two characters (a and b represent (Alef) while c and d represent (Ba)

4.1 Experiment Scenarios and Discussions

To test our proposed system, four scenarios are designed. We conducted the first scenario to understand the effect of changing the number of training data and to evaluate the performance stability over the standardize data (without occlusion nor rotation). In this scenario, SIFT technique is used to extract features form testing and training images. Then, SVM using Gaussian kernel, k-NN ($k=5$), and Nearest Neighbor classifiers are applied to testing images to recognize the ArSL character. In this experiment scenario, the number of training images ranged from one to six images and the rest images are used as testing images. The parameters of this experiment are 16×16 *Psize*, *Nbins*=4, *Nangels*=8, and *PeakThr*=0. Table (1) summarize the results obtained from this scenario.

Table 1. Accuracy (in %) of ArSL gestures when applying SIFT feature extraction method using different training images

Classifier	No. of Training Images		
	5	3	1
Min. Dist.	100	99.2	98.9
k-NN (k=5)	100	98.9	98.9
SVM	100	99.2	98.9

In the second scenario, we run an experiment to investigate the effect of changing parameter values of SIFT algorithm. We used four training images and the rest images (i.e. 3 images) are used as testing images. The accuracy of identification is computed when the *Peakthr* parameter ranged from 0 to 0.2 and different *Nangels* (2, 4 and 8). Also, we investigated the accuracy when the *Psize* increased by 4 from 4×4 to reach to 32×32 (i.e we run this experiment 4 times with different parameters). The results of this experiment is shown in Table (2).

Table 2. Identification rate (in %) of ArSL based on SIFT feature extraction method using four training images and different values of *Peakthr*, *Psize* and *Nangels*

Classifiers	PeakThr			Psize				Nangels		
	0	0.1	0.2	4x4	8x8	16x16	32x32	2	4	8
NN	100	97.7	94.2	94.2	99.2	100	93.2	94.2	98.9	100
k-NN	100	98.9	96.3	96.3	99.2	100	93.6	96.3	98.9	100
SVM	100	99.2	98.9	97.7	100	100	94.2	96.3	98.9	100

In the third scenario, we run an experiment to prove that our proposed system can overcome the problems of image rotation in different angles. In this scenario, we used SIFT Feature Extraction Methods(F.E.M.); and four training images. In the testing phase, the testing images are rotated, then it used to identify the ArSL character. Different orientations are used in our experiment, i.e., the images are rotated in the following angles: (0° , 45° , 90° , 135° , 180° , 225° , 270° 315°). The results of this experiment is shown in Table (3).

Table 3. Accuracy (in (%)) of our system when using rotated images in different angles

F.E.M.	Matching	Angles of rotation ($^\circ$)							
		0	45	90	135	180	225	270	315
SIFT	Min Dist.	100	98.9	97.8	96.7	100	97.8	100	98.9
	k-NN 5	100	100	100	96.7	100	98.9	100	100
	SVM	100	100	98.9	98.9	100	98.9	100	100

In fourth scenario, we tested our system when the gestures are shifted or occluded horizontally and vertically in different percentage of its sizes. In other word, we try to investigate whether our proposed **system** is robust against occluded images while identifying the character. In this scenario, we used four training images and the testing images are occluded to identify the ArSL character. The results of this experiment **are** shown in Table (4).

Table 4. Accuracy (in %) of our system in case of image occlusion

F.E.M.	Matching	Percentage of Occlusion					
		Horizontal			Vertical		
		20	40	60	20	40	60
SIFT	Nearest Neighbor	98.9	93.3	34.4	98.9	95.6	32.2
	k-NN 5	97.8	95.6	38.9	97.8	96.7	53.3
	SVM	98.9	95.6	52.2	98.9	96.7	45.6

In order to evaluate the performance of our proposed system, we have considered the percentage of the total number of Arabic sign characters that were correct as the main factor. The results, summarized in Table (1, 2, 3 and 4), will be discussed to confirm this factor.

From Table (1), we can notice that the accuracy of identifying ArSL based on gestures achieved very high results. It can be also seen that the accuracy slightly decreased when the number of training images are decreased. Moreover, it can be remarked that the accuracy achieved by applying SVM classifier is better than the one accomplished by applying the minimum distance and k-NN classifiers. Also, we note that, SIFT achieved excellent results. Also we note from the experiment that, the length of each feature vector reached to 28800, which needs more CPU time. But, after applying LDA the feature length becomes only 210 features and then the classification process becomes more faster. Table (2) shows that the best accuracy achieved when the *PeakThr*=0.0, *Psize*= 16×16 and *Nangels* is 8. When the *PeakThr* parameter increases, the number of features decreases and then the accuracy decreases. The highest accuracy achieved when *Psize*= 16×16 . When the *Psize* increases, then SIFT will consider as global features and decreases *Psize* will not extract more features which are necessary in identification process. Increasing the *Nangels* and *Nbins* will extract and collect features in different orientations and solve rotation problem.

Table (3) shows that our proposed system is robust against to rotation in different angels (i.e. any rotation of Arabic sign images has no effect on the accuracy of our proposed system). In addition, it is clear that the use of SVM classifier to recognize the sign characters achieved best results among the other classifier.

Table (4) proves that our system is not affected when the testing images are horizontally and vertically occluded in different percentage of the image's size. Thanks to SIFT technique which extracts features robust to image occlusion. From this table, it can be noticed that the SVM classifier, among other classifiers,

achieved best recognition rate of Arabic sign characters when the images are occluded. Also we note that, SIFT are robust against occlusion and achieved good results. In comparison to related works, i.e., Al-Jarrah's systems which achieved (94.4% and 96.5%) and Assaleh's system which accomplished (93.5%), our proposed system, which achieved around (99.5%) is better than these two systems as shown in Table 5. In addition, our system has accomplished the early the same percentage (99%) in case of rotated and occluded images.

Table 5. A comparison between proposed system and previous systems

Author	Accuracy in (%)
K. Assaleh et al. [1]	93.5%
Al-Jarrah et al. [6]	94.4%
Al-Jarrah et al. [9]	97.5%
Mohandes et al. [12]	87%
Our proposed	99%

5 Conclusions

In this paper, we have proposed a system for ArSL recognition based on gesture extracted from Arabic sign images. We have used SIFT technique to extract these features. The SIFT is used as it extracts invariant features which are robust to rotation and occlusion. Then, LDA technique is used to solve dimensionality problem of the extracted feature vectors and to increase the separability between classes, thus increasing the accuracy for our system. In our proposed system, we have used three classifiers, SVM, k-NN, and minimum distance. The experimental results showed that our system has achieved an excellent accuracy around 98.9%. Also, the results proved that our approach is robust against any rotation and they achieved an identification rate of near to 99%. In case of image occlusion (about 60% of its size), our approach has accomplished an accuracy (approximately 50%). In our future work, we are going to find a way to improve the results of our system in case of image occlusion and also increase the size of the dataset to check its scalability. Also, we will try to identify characters from video frames and then try to implement real time ArSL system.

References

1. Assaleh, K., Al-Rousan, M.: Recognition of arabic sign language alphabet using polynomial classifiers. EURASIP Journal on Applied Signal Processing 2005, 2136–2145 (2005)
2. Samir, A., Aboul-Ela, M.: Error detection and correction approach for arabic sign language recognition. In: 2012 Seventh International Conference on Computer Engineering & Systems (ICCES), pp. 117–123. IEEE (2012)
3. Tolba, M., Elons, A.: Recent developments in sign language recognition systems. In: 2013 8th International Conference on Computer Engineering & Systems (ICCES), pp. xxxvi–xlii. IEEE (2013)

4. El-Bendary, N., Zawbaa, H.M., Daoud, M.S., Nakamatsu, K., et al.: Arslat: Arabic sign language alphabets translator. In: 2010 International Conference on Computer Information Systems and Industrial Management Applications (CISIM), pp. 590–595. IEEE (2010)
5. Tolba, M., Samir, A., Aboul-Ela, M.: Arabic sign language continuous sentences recognition using pcnn and graph matching. *Neural Computing and Applications* 23(3-4), 999–1010 (2013)
6. Al-Jarrah, O., Halawani, A.: Recognition of gestures in arabic sign language using neuro-fuzzy systems. *Artificial Intelligence* 133(1), 117–138 (2001)
7. Albelwi, N.R., Alginahi, Y.M.: Real-time arabic sign language (arsl) recognition. In: International Conference on Communications and Information Technology (IC-CIT 2012), Tunisia, pp. 497–501 (2012)
8. Tolba, M., Samir, A., Abul-Ela, M.: A proposed graph matching technique for arabic sign language continuous sentences recognition. In: 8th IEEE International Conference on Informatics and Systems (INFOS), pp. 14–20 (2012)
9. Al-Jarrah, O., Al-Omari, F.A.: Improving gesture recognition in the arabic sign language using texture analysis. *Journal of Applied Artificial Intelligence* 21(1), 11–33 (2007)
10. Youssif, A.A., Aboutabl, A.E., Ali, H.H.: Arabic sign language (arsl) recognition system using hmm. *International Journal of Advanced Computer Science and Applications* (IJACSA) 2(11) (2011)
11. Mohandes, M., Deriche, M., Liu, J.: Image-based and sensor-based approaches to arabic sign language recognition. *IEEE Transactions on Human-Machine Systems* 44(4), 551–557 (2014)
12. Mohandes, M., Deriche, M.: Image based arabic sign language recognition. In: *Proceedings IEEE of the Eighth International Symposium on Signal Processing and Its Applications*, vol. 1, pp. 86–89. IEEE (2005)
13. El-Gayyar, M., Yamany, H.F.E., Gaber, T., Hassanien, A.E.: Social network framework for deaf and blind people based on cloud computing. In: Ganzha, M., Maciaszek, L.A., Paprzycki, M. (eds.) *FedCSIS*, pp. 1301–1307 (2013)
14. Meng, Y., Tiddeman, B., et al.: Implementing the scale invariant feature transform (sift) method. *Citeseer* (2008), <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.102.180>
15. Lowe, D.G.: Object recognition from local scale-invariant features. In: *The Proceedings of the Seventh IEEE International Conference on Computer Vision*, vol. 2, pp. 1150–1157. IEEE (1999)
16. Cheung, W., Hamarneh, G.: n-sift: n-dimensional scale invariant feature transform. *IEEE Transactions on Image Processing* 18(9), 2012–2021 (2009)
17. Scholkopf, B., Mullert, K.R.: Fisher discriminant analysis with kernels. In: *Proceedings of the IEEE Signal Processing Society Workshops, Neural Networks for Signal Processing IX*, August 23-25, pp. 41–48 (1999)
18. Elhariri, E., El-Bendary, N., Fouad, M.M.M., Platoš, J., Hassanien, A.E., Hussein, A.M.M.: Multi-class SVM based classification approach for tomato ripeness. In: Abraham, A., Krömer, P., Snášel, V. (eds.) *Innovations in Bio-inspired Computing and Applications*. AISC, vol. 237, pp. 175–186. Springer, Heidelberg (2014)
19. Lee, Y.: Handwritten digit recognition using k nearest-neighbor, radial-basis function, and back propagation neural networks. *Neural Computation* 3(3), 440–449 (1991)

Stock Market Forecasting Using LASSO Linear Regression Model

Sanjiban Sekhar Roy¹, Dishant Mittal¹, Avik Basu¹, and Ajith Abraham^{2,3}

¹ School of Computing Science and Engineering¹, VIT University
Vellore, Tamilnadu, India

²IT4Innovations, VSB - Technical University of Ostrava, Czech Republic

³Machine Intelligence Research Labs (MIR Labs), Washington 98071, USA
{s.roy,dishant.mittal2011,avik.basu2011}@vit.ac.in,
ajith.abraham@ieee.org

Abstract. Predicting stock exchange rates is receiving increasing attention and is a vital financial problem as it contributes to the development of effective strategies for stock exchange transactions. The forecasting of stock price movement in general is considered to be a thought-provoking and essential task for financial time series' exploration. In this paper, a Least Absolute Shrinkage and Selection Operator (LASSO) method based on a linear regression model is proposed as a novel method to predict financial market behavior. LASSO method is able to produce sparse solutions and performs very well when the numbers of features are less as compared to the number of observations. Experiments were performed with Goldman Sachs Group Inc. stock to determine the efficiency of the model. The results indicate that the proposed model outperforms the ridge linear regression model.

Keywords: Stock price prediction, LASSO regression.

1 Introduction

Prediction of stock price is a crucial factor considering its contribution to the development of effective strategies for stock exchange transactions. The Stock market plays a crucial role in the country's economy. This is due to the fact that stock market helps in flourishing the commerce and industry that ultimately has an effect on the country's economy. Whenever the company requires funds for expanding its business or if it is setting a new venture it has two options. Either a loan can be taken from a financial organization or shares can be issued through the stock market. A company can issue its shares that are in part ownership. For issuing shares for investment in the stocks, a company must get listed in the stock exchange and after this they can accumulate the funds needed for its business. Another important function that the stock market plays is that it provides a generic platform for the sellers and buyers of stocks listed on the stock market. The buyers and sellers are basically retail and institutional investors. These people are the traders who provide funds for the

businesses by investing in stocks. If the stock's future price can be predicted, it can preclude significant losses and can certainly increase the profits.

Recently, the prediction of the stock price has gathered significant interests among investors and in incorporating variable historical series into computer algorithms in order to produce estimations of estimated price fluctuations. Due to the blaring environment the prediction of the stock price becomes very complex. Traders often rely on technical indicators based on stock data which can be collected daily. Though the usage of these indicators provides them some information about the prices, but still it is difficult to have an accurate prediction of daily to weekly trends.

For a person who is not well trained trading of stocks is risky. However a neat pile in quick intraday deals can be made if one has a fixation on spotting the trends in the market. There was a generalized mindset in recent times when depending on the beliefs of people trading was considered as game of buying and selling of stocks. Now, some new tools have been devised by the investors by utilizing a method known as technical analysis for predicting future prices from historical price data. On general technical analysis is based on technical indicators. A technical indicator for the stock price is a function that returns a worth for given stock price in some given span of time in history. Information on whether a trend will continue or whether a stock is oversold and overbought can be got from such technical indicators [1].

Apart from the technical analysis, there is a method known as fundamental analysis, which is concerned with the company that underlies the stock itself. It assesses a company's past performance as well as the trustworthiness of its accounts. Many performance ratios are produced that aid in evaluating the rationality of a stock. With the arrival of the digital computer, stock market prediction has since progressed into the technological world. Artificial neural networks (ANNs) and genetic algorithms are involved in some of the most noticeable techniques. ANNs can be considered as approximations of mathematical functions. The use of ANN mimics how the human brain functions, by serving computers with the immense data to mimic human thinking. The feed forward network using the backward propagation of errors algorithm to update the network weights is the most accepted form of ANN for stock market prediction that is currently in use [10]. These networks are commonly referred to as Back propagation networks. Time delay neural network (TDNN) or the time recurrent neural network (TRN) is another type of ANN that is more convenient in forecasting the stock price.

We propose a system which is based on generalized linear regression model and use it for stock market forecasting. In this paper, we present a model that we implemented for the prediction of stock price based on the LASSO method which outperforms the ridge method and the artificial neural network model in terms of accuracy.

2 Related Works

Deng et al. [1] introduced *a stock price prediction model, which extracts features from time series data and social networks for prediction of stock prices and*

calculates its performance. The stock price movements were modeled as a function of these input features and was solved as a regression problem in a Multiple Kernel Learning regression framework by them. Yoo et al. [2] in their work explored various global events and their concerns on forecasting stock markets. They found that integrating event information with the prediction model plays very significant roles for more exact prediction. Schumann et al. [3] presented two models, namely ARIMA and ANN for the prediction of the stock price. Pakdaman Naeini et al. [4] presented two kinds of neural networks, an Elman recurrent network and a feed forward multilayer Perceptron (MLP) that were in turn utilized to measure a company's stock value based on the record of its stock share value. They demonstrated that the application of MLP neural network is more capable of calculating stock value changes rather than Elman recurrent network and linear regression method. Aseervatham et al. [5] claimed that the ridge logistic regression reaches the same performance as the Support Vector Machine. Ticknor [6] presented a Bayesian regularized artificial neural network as a unique method to estimate the financial market behavior. Data Sets from the corporations like Goldman Sachs Group Inc. and Microsoft Corp. were used to perform experiments. Nair et al. [7] proposed an automated decision tree-adaptive neuro-fuzzy hybrid automated stock market prediction system. Ajith Abraham et al [8] proposed a genetic programming method for forecasting of stock market prices. They experimented on Nasdaq stock market and S and P cnx nifty data. They combined two multiobjective optimization algorithms with techniques like svm and neural networks to get the best results. Yuehui Chen et al. [9] introduced flexible neural tree method for organized representation of stock market data, later they used genetic programming for optimization of this method.

3 Generalized Linear Models

There is a set of methods proposed for regression in which the target value is likely to be a linear combination of the input variables. In mathematical notion, if Y is the predicted value, then its value can be obtained from equation 2 as $h(x)$.

Ordinary least squares method is one of the generalized linear regression model in which linear regression fits a linear model with coefficients $w = (w_1, \dots, w_p)$ to minimize the residual sum of squares between the witnessed responses in the dataset, and the responses forecasted by the linear approximation. Mathematically, it solves a problem of minimizing the expression of type as shown in equation 3.

The linear regression can take in its fit method arrays X , y and will accumulate the coefficients w of the linear model in its `coef_member`. However, the freedom of the model terms decides coefficient assessments for Ordinary Least Squares. This method calculates the least squares solution using a singular value decomposition of X . If X is a matrix of size (n, p) then this method has a cost of $O(np^2)$, if $n \geq p$.

3.1 Mathematical Formulation

Consider the set of training vectors (x_i, y_i) , x_n belongs to R^n , y_n belongs to R ,

$$i = 1, \dots, N \quad (1)$$

The hypothesis or the linear regression output is given by

$$h(x) = \sum_{j=0}^d w_j x_j = w^T x \quad (2)$$

where w is the weight vector and d is the dimensionality of the problem or the number of features.

Also, $x_0 = 1$ has been added to make equation (2) valid.

The cost function or the squared error function is defined as

$$J(w) = \frac{1}{N} \sum_{i=0}^N (h(x_i) - y_i)^2 = \frac{1}{N} \|Xw - y\|^2 \quad (3)$$

where

$$X = \begin{bmatrix} -x_1^T & - \\ -x_2^T & - \\ \vdots & \\ -x_N^T & - \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \quad (4)$$

We need to minimize the cost function to get the optimal value of the weight vector.

Minimizing the cost function,

$$\nabla J(w) = \frac{2}{N} X^T (Xw - y) = 0 \quad (5)$$

which implies

$$X^T Xw = X^T y \quad (6)$$

hence

$$w = X^+y \tag{7}$$

where

$$X^+ = (X^T X)^{-1} X^T \tag{8}$$

Putting the value of w from (7) the optimal hypothesis is obtained.

The traditional least square method can be modified to Ridge regression model. The cost function for Ridge regression can be specified as

$$J(w) = \frac{1}{N} \sum_{i=0}^N (h(x_i) - y_i)^2 + \lambda \sum_{j=1}^d w_j^2 = \frac{1}{N} \|Xw - y\|^2 + \lambda \sum_{j=1}^d w_j^2 \tag{9}$$

where λ is the regularization parameter. After minimizing the cost function, we get the coefficients as

$$w = (X^T X + \lambda I)^{-1} X^T y \tag{10}$$

4 Suggested Linear Model

Another modification of the least square method is the LASSO model which stands for Least Absolute Shrinkage and Selection Operator. The suggested model is used for the estimation of sparse coefficients. It is valuable in some backgrounds due to its affinity to prefer solutions with fewer parameter values, efficiently decreasing the number of variables upon which the given solution is dependent. The appropriate group of weights which are not zero can be recovered under certain conditions. Mathematically, a linear model trained with l_1 prior as regularize is comprised in it.

To fit the coefficients, the algorithm that is used for the implementation in the class LASSO is coordinate descent [11].

The new objective for LASSO can be defined as

$$J(w) = \frac{1}{N} \sum_{i=0}^N (h(x_i) - y_i)^2 + \lambda \sum_{j=1}^d |w_j| = \frac{1}{N} \|Xw - y\|^2 + \lambda \sum_{j=1}^d |w_j| \tag{11}$$

The added term corresponds to l_1 -norm. The lasso estimate thus explains the minimization of the least square penalty with $\lambda \|w\|_1$ added where λ is regularization parameter and $\|w\|_1$ is the l_1 -norm of the parameter vector.

5 Experimentation Results

The research data utilized for predicting stock market prices in this study was gathered for Goldman Sachs Group, Inc. (GS). The total number of instances considered for this study were 3686 trading days, from 4 May 1999 to 3 January 2014. Each of the samples composed of daily information including low price, high price, opening price, close price, and trading volume. The training data set was selected as the first 70% of the samples, while the testing data consisted the remaining 30% of the samples. The LASSO model was used to predict the price of the chosen stock for the future days.

A comparison study was performed to test the efficiency of the model suggested in this study to another linear regression model that is named as Ridge and a Bayesian regularized artificial neural network by Jonathan L. Ticknor. For this method, the dataset was gathered for Goldman Sachs Group, Inc. In total number of instances considered for this study were 734 trading days, from 4 January 2010 to 31 December 2012. The training data set was chosen as the first 80% of the samples and the remaining 20% was used as testing dataset.

5.1 Algorithm

LASSO REGRESSION()

1. data ← read ('data.csv')
2. (train_features, train_stock_price) ← training_function()
3. (test_features, test_stock_price) ← testing_function()
4. Model ← LASSO_train(train_features, train_stock_price, lambda)
5. stock_price_predict ← LASSO_predict(train_features)
6. MAPE ← mean [abs{(test_stock_price - stock_price_predict) / test_stock_price}] * 100
7. RMSE ← sqrt [mean{(test_stock_price - stock_price_predict)²}]

5.2 Results and Discussion

The performance of the LASSO Linear regression method was measured by computing root mean square error (RMSE) and the mean absolute percentage error (MAPE). These performance metrics have been used in a number of studies and ensures an effective means of deciding the robustness of the model for predicting daily. It can be represented as :

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - p_i)^2}{n}} \tag{12}$$

where n is the total number of trading days p_i is the predicted stock price on day i and y_i is the actual stock price on the same day.

The Mean Absolute Percentage error (MAPE) metric is first found by calculating the absolute value of the variation between the actual stock price and the expected stock price. The MAPE value is calculated using the following equation.

MAPE formula:

$$MAPE = \frac{\sum_{i=1}^n \frac{|y_i - p_i|}{y_i}}{n} \times 100\% \tag{13}$$

where n is the total number of trading days p_i is the predicted stock price on day i and y_i is the actual stock price on the same day.

Table 1. RMSE and MAPE of training and test set (Ridge vs LASSO)

Method	Training RMSE	Test RMSE	Training MAPE	Test MAPE
Ridge	1.7648	3.2272	1.3028	1.8065
LASSO	1.1403	2.5401	0.9304	1.4726

Table 2. Forecast accuracy comparison with Jonathan L. Ticknor

Method	Training MAPE (%)	Testing MAPE (%)
Bayesian Regularized ANN	1.5235	1.3291
LASSO	0.1806	0.6869

Table 1 presents the experimental results for the two methods (Ridge and LASSO) chosen for prediction over the Goldman Sachs (GS) dataset .The RMSE and MAPE values were calculated for the training and testing dataset to monitor the effectiveness of the models. It was deduced that the testing set MAPE of LASSO regression is less than the testing set MAPE of the Ridge regression method indicating that LASSO method is better than Ridge. The same results are reflected in the graphs from Fig. 1-6. This can be said according to the proximity of the regression line for each graph to the data points.

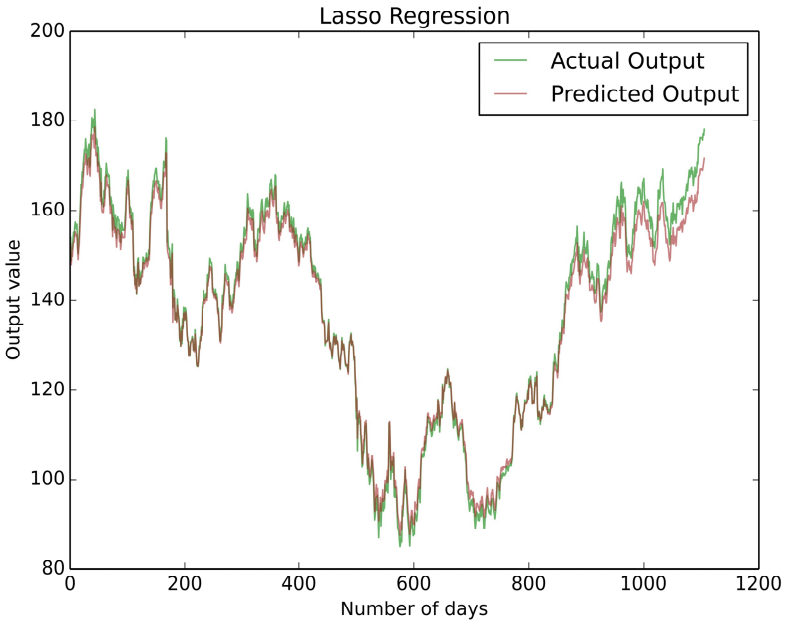


Fig. 1. Future day prediction (LASSO)

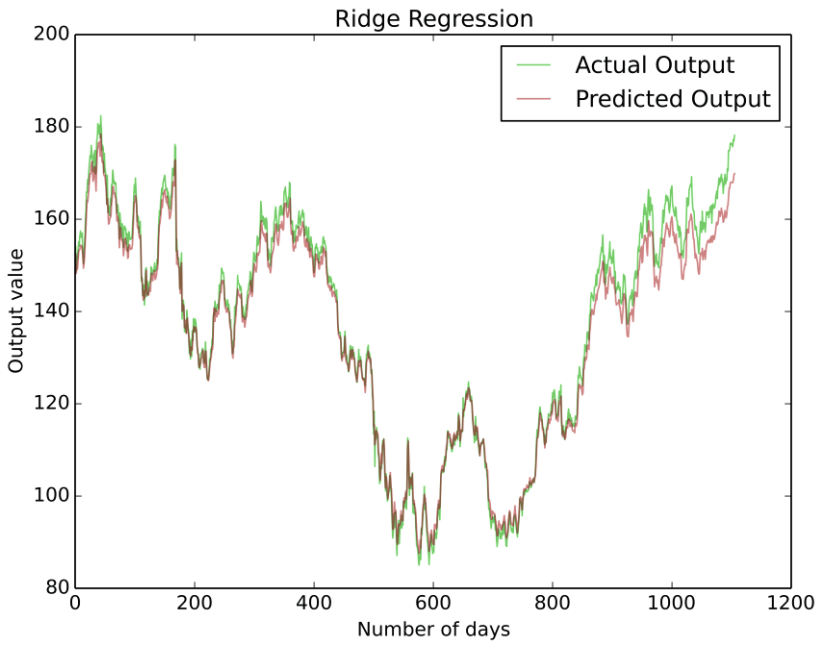


Fig. 2. Future day prediction (Ridge)



Fig. 3. Target vs predicted stock price using training set (LASSO)

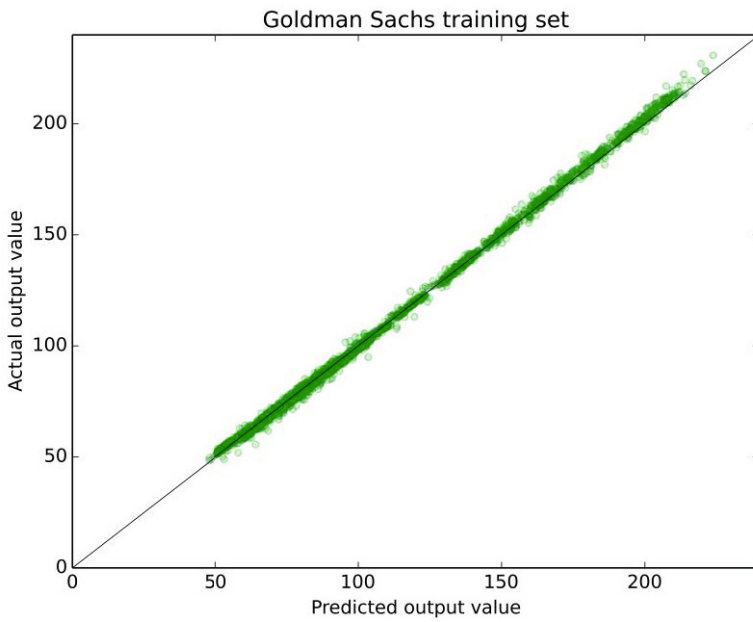


Fig. 4. Target vs predicted stock price using training set (Ridge)

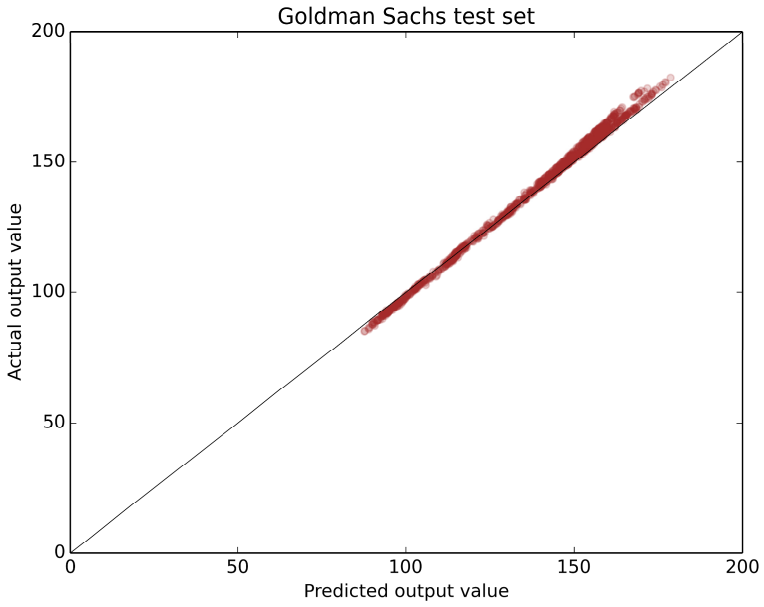


Fig. 5. Target vs predicted stock price using test set (LASSO)

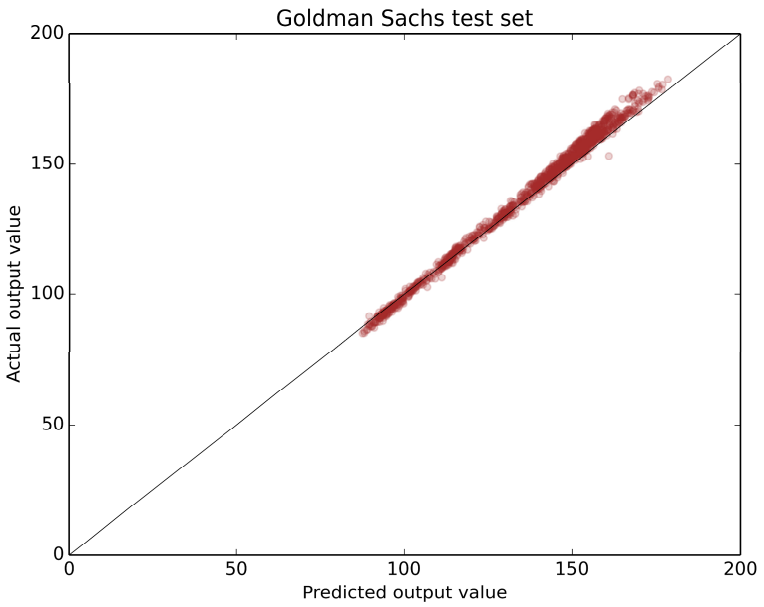


Fig. 6. Target vs predicted stock price using test set (Ridge)

6 Conclusions

Empirical results indicated that the model outperformed the ridge linear regression model and Bayesian regularized artificial model. The model resulted in a MAPE of 0.6869 with respect to 1.3291 that Bayesian artificial neural network method produced for the mentioned dataset. A MAPE value of 1.4726 and RMSE value 2.5401 was produced by using LASSO algorithm, whereas MAPE value 1.8065 and RMSE value 3.2272 was reported by utilizing ridge regression algorithm with respect to the mentioned dataset for 3686 instances. To evaluate the effectiveness of this model, the network was successfully compared with a Bayesian regularized artificial neural network model. Forecast of stock market drifts is very important for the development of effective trading policies.

References

1. Deng, S., Takashi, M., Kei, S., Tatsuro, S.: Akito Sakurai.: Combining technical analysis with sentiment analysis for stock price prediction. In: IEEE Ninth International Conference on Dependable, Autonomic and Secure Computing (DASC), pp. 800–807. IEEE (2011)
2. Yoo, P.D., Kim, M.H., Jan, T.: Machine learning techniques and use of event information for stock market prediction: A survey and evaluation. In: International Conference on Computational Intelligence for Modeling, Control and Automation, 2005 and International Conference on Intelligent Agents, Web Technologies and Internet Commerce, vol. 2, pp. 835–841 (2005)
3. Schumann, M., Lohrbach, T.: Comparing artificial neural networks with statistical methods within the field of stock market prediction. In: Proceeding of the Twenty-Sixth Hawaii International Conference on in System Sciences, vol. 4, pp. 597–606. IEEE (1993)
4. Naeini, M.P., Taremiyan, H., Hashemi, H.B.: Stock market value prediction using neural networks. In: 2010 International Conference on Computer Information Systems and Industrial Management Applications (CISIM), pp. 132–136. IEEE (2010)
5. Aseervatham, S., Antoniadis, A., Gaussier, E., Buret, M., Denneulin, Y.: A sparse version of the ridge logistic regression for large-scale text categorization. *Pattern Recognition Letters* 32(2), 101–106 (2011)
6. Ticknor, J.L.: A Bayesian regularized artificial neural network for stock market forecasting. *Expert Systems with Applications* 40(14), 5501–5506 (2013)
7. Nair, B.B., Minuvarthini, M., Sujithra, B., Mohandas, V.: Stock market prediction using a hybrid neuro-fuzzy system. In: 2010 International Conference on Advances in Recent Technologies in Communication and Computing (ARTCom), pp. 243–247. IEEE (2010)
8. Abraham, A., Grosan, C., Han, S.Y., Gelbukh, A.: Evolutionary multiobjective optimization approach for evolving ensemble of intelligent paradigms for stock market modeling. In: Gelbukh, A., de Albornoz, Á., Terashima-Marín, H. (eds.) MICAI 2005. LNCS (LNAI), vol. 3789, pp. 673–681. Springer, Heidelberg (2005)
9. Chen, Y., Yang, B., Abraham, A.: Flexible neural trees ensemble for stock index modeling. *Neurocomputing* 70(4), 697–703 (2007)
10. Pathak, A.: Predictive time series analysis of stock prices using neural network classifier. *International Journal of Computer Science and Engineering Technology*, 2229–3345 (2014)
11. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-Learn: Machine Learning in Python. *JMLR Journal of Machine Learning Research*, 2825–2830 (2011)

Author Index

- Abraham, Ajith 183, 293, 303, 315, 371
Ahmed, Nada 315
Aldosari, Hamoud M. 303
Alhamedi, Adel H. 303
Ali, Mohamed A. 193
Assefa, Abebayehu 245
Assefa, Dawit 65
- Balík, Ladislav 347
Basu, Avik 371
Berger, Ondrej 335
Berhan, Eshetie 87, 133, 261, 283
Beshah, Birhanu 133, 233
Birhan, Eshetie 271
Buriánek, Tomáš 205
- Cimler, Richard 347
Czopik, Jan 77
- Dagnaw, Girum 87
Đurica, Maroš 99
- Egeonu, Darlington Ihunanyachukwu 109
Emary, E. 1
Enibe, Samuel Ogbonna 109
- Gaber, Tarek 359
Gabralla, Lubna A. 293
Girma, Kidist 133
Grosan, Crina 1
- Habib, Abduletif 261
Hailu, Haftu 271
Hassanien, Aboul Ella 359
Hassenian, Abul Ella 1
- Heckenbergerova, Jana 15
Horalek, Josef 347
- Ježek, David 173
Jilcha, Kassu 261, 271, 283
Jote, Netsanet 233
- Kakrda, Petr 335
Kassaw, Yonas 65
Kitaw, Daniel 233
Kopka, Martin 327
Košinár, Michael Alexander 77
Krejcar, Ondrej 335
Krömer, Pavel 15
Kumaran, Santhi 27
- Lamesgin, Gizeaddis 65
Liebzeit, Radek 173
- Mahersia, Hela 293
Matyska, Jan 347
Messele, Yidnekachew 215, 245
Mittal, Dishant 371
Mulugeta, Mequanent 147
Musilek, Petr 15
- Nasser, Rahma 215
Negash, Andinet 123
Njoku, Howard Okezie 109
- Okolo, Patrick Nwosa 109
- Peterek, Tomáš 205
Platos, Jan 163
Prilepok, Michal 163

- Refaat, Basma 359
Roy, Sanjiban Sekhar 371
Salih, Abdelhamid Salih Mohamed 183
Shahin, M.K. 359
Singh, Nagendra P. 123
Snášel, Václav 163, 205, 303, 327
Sobeslav, Vladimír 347
Solomon, Hermela 283
Soori, Hussein 163
Štolfa, Jakub 77, 327
Štolfa, Svatopluk 77, 327
Tamir, Kassahun 43
Teferi, Amare 87
Tekeba, Menore 43
Tesfaye, Erimas 133
Tharwat, Alaa 359
Yilma, Eyerusalem 215
Zaorálek, Lukáš 205
Zawbaa, Hossam M. 1