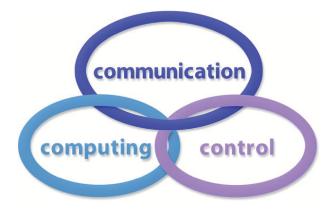
INTERNATIONAL JOURNAL

of

COMPUTERS, COMMUNICATIONS & CONTROL

With Emphasis on the Integration of Three Technologies

IJCCC



Year: 2011 Volume: 6 Number: 2 (June)

Agora University Editing House



www.journal.univagora.ro

International Journal of Computers, Communications & Control



EDITOR IN CHIEF: Florin-Gheorghe Filip

Member of the Romanian Academy Romanian Academy, 125, Calea Victoriei 010071 Bucharest-1, Romania, ffilip@acad.ro

ASSOCIATE EDITOR IN CHIEF:

Ioan Dzitac Aurel Vlaicu University of Arad, Romania Elena Dragoi, 2, Room 81, 310330 Arad, Romania ioan.dzitac@uav.ro

MANAGING EDITOR:

Mişu-Jan Manolescu Agora University, Romania Piata Tineretului, 8, 410526 Oradea, Romania rectorat@univagora.ro

EXECUTIVE EDITOR:

Răzvan Andonie

Central Washington University, USA 400 East University Way, Ellensburg, WA 98926, USA andonie@cwu.edu

TECHNICAL SECRETARY:

Cristian Dziţac R & D Agora, Romania rd.agora@univagora.ro Emma Margareta Văleanu R & D Agora, Romania evaleanu@univagora.ro

EDITORIAL ADDRESS:

R&D Agora Ltd. / S.C. Cercetare Dezvoltare Agora S.R.L. Piaţa Tineretului 8, Oradea, jud. Bihor, Romania, Zip Code 410526 Tel./ Fax: +40 359101032 E-mail: ijccc@univagora.ro, rd.agora@univagora.ro, ccc.journal@gmail.com Journal website: www.journal.univagora.ro

DATA FOR SUBSCRIBERS

Supplier: Cercetare Dezvoltare Agora Srl (Research & Development Agora Ltd.) Fiscal code: RO24747462 Headquarter: Oradea, Piata Tineretului Nr.8, Bihor, Romania, Zip code 410526 Bank: MILLENNIUM BANK, Bank address: Piata Unirii, str. Primariei, 2, Oradea, Romania IBAN Account for EURO: RO73MILB000000000932235 SWIFT CODE (eq.BIC): MILBROBU

International Journal of Computers, Communications & Control



EDITORIAL BOARD

Boldur E. Bărbat

Lucian Blaga University of Sibiu Faculty of Engineering, Department of Research 5-7 Ion Rațiu St., 550012, Sibiu, Romania bbarbat@gmail.com

Pierre Borne

Ecole Centrale de Lille Cité Scientifique-BP 48 Villeneuve d'Ascq Cedex, F 59651, France p.borne@ec-lille.fr

Ioan Buciu

University of Oradea Universitatii, 1, Oradea, Romania ibuciu@uoradea.ro

Hariton-Nicolae Costin

Faculty of Medical Bioengineering Univ. of Medicine and Pharmacy, Iaşi St. Universitatii No.16, 6600 Iaşi, Romania hcostin@iit.tuiasi.ro

Petre Dini

Cisco 170 West Tasman Drive San Jose, CA 95134, USA pdini@cisco.com

Antonio Di Nola

Dept. of Mathematics and Information Sciences Università degli Studi di Salerno Salerno, Via Ponte Don Melillo 84084 Fisciano, Italy dinola@cds.unina.it

Ömer Egecioglu

Department of Computer Science University of California Santa Barbara, CA 93106-5110, U.S.A omer@cs.ucsb.edu

Constantin Gaindric

Institute of Mathematics of Moldavian Academy of Sciences Kishinev, 277028, Academiei 5, Moldova gaindric@math.md

Xiao-Shan Gao

Academy of Mathematics and System Sciences Academia Sinica Beijing 100080, China xgao@mmrc.iss.ac.cn

Kaoru Hirota

Hirota Lab. Dept. C.I. & S.S. Tokyo Institute of Technology G3-49,4259 Nagatsuta,Midori-ku,226-8502,Japan hirota@hrt.dis.titech.ac.jp

George Metakides

University of Patras University Campus Patras 26 504, Greece george@metakides.net

Ştefan I. Nitchi

Department of Economic Informatics Babes Bolyai University, Cluj-Napoca, Romania St. T. Mihali, Nr. 58-60, 400591, Cluj-Napoca nitchi@econ.ubbcluj.ro

Shimon Y. Nof

School of Industrial Engineering Purdue University Grissom Hall, West Lafayette, IN 47907, U.S.A. nof@purdue.edu

Stephan Olariu

Department of Computer Science Old Dominion University Norfolk, VA 23529-0162, U.S.A. olariu@cs.odu.edu

Horea Oros

Dept. of Mathematics and Computer Science University of Oradea, Romania St. Universitatii 1, 410087, Oradea, Romania horos@uoradea.ro

Gheorghe Păun

Institute of Mathematics of the Romanian Academy Bucharest, PO Box 1-764, 70700, Romania gpaun@us.es

Mario de J. Pérez Jiménez

Dept. of CS and Artificial Intelligence University of Seville Sevilla, Avda. Reina Mercedes s/n, 41012, Spain marper@us.es

Dana Petcu

Computer Science Department Western University of Timisoara V.Parvan 4, 300223 Timisoara, Romania petcu@info.uvt.ro

Radu Popescu-Zeletin

Fraunhofer Institute for Open Communication Systems Technical University Berlin, Germany rpz@cs.tu-berlin.de

Imre J. Rudas

Institute of Intelligent Engineering Systems Budapest Tech Budapest, Bécsi út 96/B, H-1034, Hungary rudas@bmf.hu

Athanasios D. Styliadis

Alexander Institute of Technology Agiou Panteleimona 24, 551 33 Thessaloniki, Greece styl@it.teithe.gr

Gheorghe Tecuci

Learning Agents Center George Mason University, USA University Drive 4440, Fairfax VA 22030-4444 tecuci@gmu.edu

Horia-Nicolai Teodorescu

Faculty of Electronics and Telecommunications Technical University "Gh. Asachi" Iasi Iasi, Bd. Carol I 11, 700506, Romania hteodor@etc.tuiasi.ro

Dan Tufiş

Research Institute for Artificial Intelligence of the Romanian Academy Bucharest, "13 Septembrie" 13, 050711, Romania tufis@racai.ro

Lotfi A. Zadeh Department of Computer Science and Engineering University of California Berkeley, CA 94720-1776, U.S.A. zadeh@cs.berkeley.edu



International Journal of Computers, Communications & Control



Short Description of IJCCC

Title of journal: International Journal of Computers, Communications & Control **Acronym:** IJCCC

International Standard Serial Number: ISSN 1841-9836, E-ISSN 1841-9844 Publisher: CCC Publications - Agora University

Starting year of IJCCC: 2006

Founders of IJCCC: Ioan Dzitac, Florin Gheorghe Filip and Mişu-Jan Manolescu Logo:



Number of issues/year: IJCCC has 4 issues/odd year (March, June, September, December) and 5 issues/even year (March, September, June, November, December). Every even year IJCCC will publish a supplementary issue with selected papers from the International Conference on Computers, Communications and Control. Coverage:

- Beginning with Vol. 1 (2006), Supplementary issue: S, IJCCC is covered by Thomson Reuters SCI Expanded and is indexed in ISI Web of Science.
 - Journal Citation Reports/Science Edition 2009:
 - Impact factor = 0.373
 - Immediacy index = 0.205
 - Beginning with Vol. 2 (2007), No.1, IJCCC is covered in EBSCO.
 - Beginning with Vol. 3 (2008), No.1, IJCCC, is covered in SCOPUS.

Scope: IJCCC is directed to the international communities of scientific researchers in universities, research units and industry. IJCCC publishes original and recent scientific contributions in the following fields: Computing & Computational Mathematics; Information Technology & Communications; Computer-based Control.

Unique features distinguishing IJCCC: To differentiate from other similar journals, the editorial policy of IJCCC encourages especially the publishing of scientific papers that focus on the convergence of the 3 "C" (Computing, Communication, Control).

Policy: The articles submitted to IJCCC must be original and previously unpublished in other journals. The submissions will be revised independently by at least two reviewers and will be published only after completion of the editorial workflow.

Copyright © 2006-2011 by CCC Publications

Contents

An Improved Computational Model for Adaptive Communication Channel Estimation	i-
S.A. Akinboro, G.A. Aderounmu, E.A. Olajubu A.O. Ajayi, I.K. Ogundoyin, O.M Olaniyan	204
Enhanced Security Protocol in Wireless Sensor Networks	
T.C. Aseri, N. Singla	214
Perturbation in Population of Pulse-Coupled Oscillators Leads to Emergence of Structure	
A. Bartha, D. Dumitrescu	222
A Time-Bound Ticket-Based Mutual Authentication Scheme for Cloud Computing	
Z. Hao, S. Zhong, N. Yu	227
A Network Coding based DTN Convergence Layer Reliable Transport Mechanism over InterPlaNetary Networks	
S. Haoliang, L. Lixiang, H. Xiaohui	236
Multi-Objective Optimization for the m-PDPTW: Aggregation Method With Use of Genetic Algorithm and Lower Bounds	h
I. Harbaoui Dridi, R. Kammarti, M. Ksouri, P. Borne	246
An Authenticated Key Agreement Protocol Using Isogenies Between Ellipti Curves	ic
D. He, J. Chen, J. Hu	258
On the Power of Small Size Insertion P Systems	
A. Krassovitskiy	266
Modelling, Implementation and Application of a Flexible Manufacturing Cell F. Leighton, R. Osorio, G. Lefranc	278
A Hybrid Artificial Bee Colony Algorithm for Flexible Job Shop Scheduling Problems	
J. Li, Q. Pan, S. Xie, S. Wang	286

The Communication in Distributed Client - Server Systems Used for Management of Flexible Manufacturing Systems	
V. Lupu, D.E. Tiliute	297
Fuzzy Based Packet Dropping Scheme in Wireless Cellular Networks	
J.D. Mallapur, S.S. Manvi, D.H. Rao	305
Classification Performance Using Principal Component Analysis and Differen Value of the Ratio R	ıt
J. Novakovic, S. Rankov	317
Modeling Uncertainty in a Decision Problem by Externalizing Information	
I. Parpucea, B. Pârv, T. Socaciu	328
Impact of Poor Requirement Engineering in Software Outsourcing: A Study of Software Developers' Experience	n
I. Perera	337
Accessing Information Sources using Ontologies	
D. Sun, H. Jung, C. Hwang, H. Kim, S. Park	349
Digital Control of a Waste Water Treatment Plant	
R. Vilanova, J.D. Rojas, V.M. Alfaro	367
Erratum	375
Author index	376

An Improved Computational Model for Adaptive Communication Channel Estimation

S.A. Akinboro, G.A. Aderounmu, E.A. Olajubu A.O. Ajayi, I.K. Ogundoyin, O.M Olaniyan

S.A. Akinboro, O.M Olaniyan

Computer Science and Technology Bells University of Technology, Ota, Nigeria E-mail: akinboro2002@yahoo.com

G.A. Aderounmu, E.A. Olajubu, A.O. Ajayi, I.K. Ogundoyin

Computer Science and Engineering Obafemi Awolowo University, Ile-Ife, Nigeria

> **Abstract:** Channel estimation is an important and necessary function performed by modern wireless receivers. The goal of channel estimation is to measure the effects of the channel on known or partially known transmission. The usual practice in acquiring knowledge about a channel is to model the channel and then acquire the parameters involved in the model. This paper proposes a variable partial update model for adaptive communication channel estimation with a view to improving signal error at the receiver station. The proposed model is composed of finite impulse response transversal adaptive filter and least mean square adaptation algorithm. The performance of the proposed model was compared with the full update model. The evaluation results indicated that the proposed model performed better than the full update model in terms of computational complexity, memory load, and convergence rate.

> **Keywords:** Adaptation algorithm, Computational complexity, Memory load, convergence rate, Partial update

1 Introduction

Channel refers to one way telecommunication link or transmission medium through which information or signal is transmitted from a transmitter to a receiver. It may either be physical or logical depending on the application, for example cables, and radio frequency are physical channels while control and traffic channels within the ratio frequency channel are logical channels. In estimation theory, it is assumed that the desired information is embedded in a noisy signal. Noise adds uncertainty and if there is no uncertainty then there would be no need for estimation. The channel deforms the transmitted signals often in unpredictable ways. To retrieve the information that is transmitted, the received signal has to be processed. The retrieval of information about the channel either from the received signal or from the signal sent is known as channel estimation. The major sources of impairment in wireless channels include channel time-variation, Inter Symbol Interference, and Co Channel Interference [6]. In order to deal with these problems, the transmitted signal needs to be processed at the receiver station. A core feature of many modern communication systems is their ability to adapt to the working environment. The technology at the heart of these flexible systems is an adaptive digital filter whose coefficients change in response to external conditions [3]. An adaptive filter has coefficients that are updated by some types of adaptive algorithms to improve or optimize its response to a desired performance. In general, adaptive filters consist of two basic parts: the filter which applies the required processing on the incoming signal to be filtered and an adaptive algorithm which adjusts the coefficient of the filter to improve its performance [1]. Many computationally efficient algorithms have been developed for adaptive filtering [2]. They are based on either a statistical approach such as least-mean square (LMS) algorithm or a deterministic approach such as recursive least-square (RLS) algorithm. The major advantage of the LMS algorithm is its computational simplicity. The RLS algorithm conversely, offers faster convergence but with a higher degree of computational complexity. Increased popularity of mobile phones and other wireless communication products provide the demand for developing appropriate techniques to improve the performance of existing system for reliable transmission of data over wireless communicational system. Many applications in wireless communication like channel estimation and echo cancellation require the adaptive filter to have a very large number of co-efficient, and updating the entire coefficient is costly in terms of power, memory and computation, and sometimes impractical for mobile units [2]. In this paper, an improved computational model for adaptive channel estimation is proposed. The adaptive channel estimator is modelled using finite impulse response traversal adaptive filter. The adaptation process of the filtering is performed using Variable Partial Update LMS algorithm. The coefficient of the adaptive filter is partially updated to reduce computational cost and memory load, while at the same time updating the step size parameter to enhance the speed and the accuracy of convergence. The outline of this paper is as follows. In section 2, the review of related works is carried out, and in section 3, the proposed model is described together with the adaptation algorithm used. Simulation and results are presented in section 4, and section 5 concludes the study.

2 Related work

Several research works have been done in developing computational models and adaptation algorithms for adaptive communication channel estimation but all the models failed to address the issues of convergence rate, memory load and computational complexity efficiently. In particular, no research work has considered the comparison of the full update and the partial update of the adaptive digital filter coefficients using the parameters memory load, convergence rate and computational complexity. In an attempt to reduce computational complexity and improve asymptotic performance, [4] Proposed active tap detection LMS algorithm. Though the research work reduces computational burdens as well as unsatisfactory poor convergence rate asymptotic performance of the adaptive tap in a long channel but storage location is provided for the entire adaptive tap which is quite expensive. Also, [2] proposed partial updating of the LMS adaptive filter to reduce the cost of power, memory load and computation. Sequential partial update LMS is employed in their work. They analyzed the alternating odd/even partial update LMS algorithm and derived stability bound on step size parameter for wide sense stationary and cyclo-stationary signals based on external properties of the matrix 2-norm but comparing with the proposed model, the memory load and computational complexity is still large. The behavior of three variants of variable step size LMS algorithm for training based multi-user detection in a CDMA system was studied by [10]. Two of the algorithm have smaller computational complexity and memory load but still suffers from the fact that their steady state error and speed of convergence depend on the same parameters (the step size), therefore complementary pair variable step size LMS was introduced. Although the proposed algorithm has an increased computational complexity and memory load, but it has better speed performance and more simple parameters setup which are very important in practical applications. A number of previous LMS algorithms were analyzed by [5]. They pointed out their weaknesses and proposed a variant of modified variable step size LMS algorithm which was tested and can ensure convergence in any cases and can provide a higher speed of convergence and a better level of tracking ability.

The normalized LMS algorithm is adapted to get higher speed of convergence by adjusting the step size through the power of input signals. However, we can hardly make an accurate estimate of the auto correlation matrix and mean square error practically. They also said that variants of variable step size LMS algorithms have been proposed, but in all of them, the step size equation can be written as $\mu(n + 1) = \alpha \mu(n) + \gamma p$ where p is a function and depends on the different VSS-LMS algorithm, α and γ are constant, μ is the step size and n is the time index. Previous research works focused on improving the effect of channel on transmitted signal through proposition of variants of adaptation algorithm and design of improved adaptive digital filters, this work however, compare the performance of full-update and variable partial update variants of the adaptation algorithm with design of adaptive digital filter using index factors. The three performance metrics considered are convergence rate, memory load and computational complexity.

3 Architecture of the proposed model

The block diagram for Variable Partial Update LMS Model is applied in this work, Figure 1 depict the model, while the adaptation module of the proposed model is an improvement on [2] which is depicted in Figure 2. The difference between the model developed in [2] and this work is the updates of the coefficients. This model uses index factors of three and five for updates while [2] uses alternate even and odd updates of the coefficients. The model consists of four major modules vis-f-vis; unknown channel, FIR filter, summer (Σ), and adaptation algorithm. The unknown channel is positioned parallel to the FIR filter so that the same input signal can be transmitted simultaneously. The estimation is adaptive and parallel, this is because in wireless situation the paths that a signal takes between the transmitter and the receiver may keep changing. The signal transmitted through the FIR filter output is the estimated training signal y(k). The two signals are transmitted through the summer to give the estimated error signal e(k). The estimated error signal is then used to update the coefficient of the FIR filter using the adaptation algorithm called Variable step size Partial Update Least Mean Square algorithm (VPU-LMS) [6].

It is assumed that the variable partial update LMS filter in fig 2 is a standard transversal FIR filter of length $L \ge 5$. Let $(x_{i,n})$ be the input sequence and let $(w_{i,n})$ denote the coefficients of the adaptive filter.

 $Wn = [w1, n, w2, nw3, n...wL, n]^T$

Xn = [x1, n, x2, nx3, n...xL, n]

Where the terms define above are for the instant n and T denotes the transpose operator. Also, let d(n) denotes the desired response. It represents a known training signal which is transmitted over a noisy channel with unknown FIR transfer function. It was assumed that d(n) obeys a FIR model given by

$$d(n) = x_n + k(n) \tag{1}$$

Where k(n) is the autoregressive noise signal that is independent of the input sequence X_n . We also assumed that the filter coefficient is a mutually exclusive set. The elements of the set are coefficients with index factors of three and five i.e. filter length that is divisible by three and five. Therefore the set S is defined as $S = \{w_3, w_5, w_6, w_9, w_{10}...\}$

The proposed algorithm called variable partial update LMS algorithm. The algorithm updates only the coefficient of the adaptive filter with index factors of three and five. It also makes sure that only the active coefficients (i.e. value not equal to zero) are used for the update process. The step size parameter is updated to ensure convergence of the algorithm. The algorithm is described in following steps:

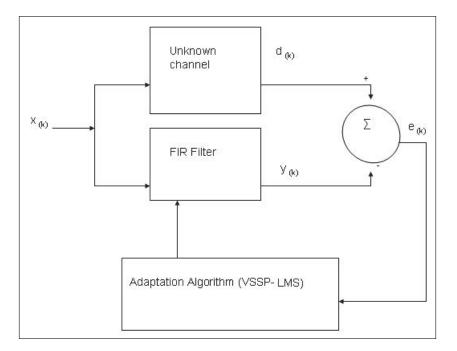


Figure 1: The block diagram for Variable Partial Update LMS Model

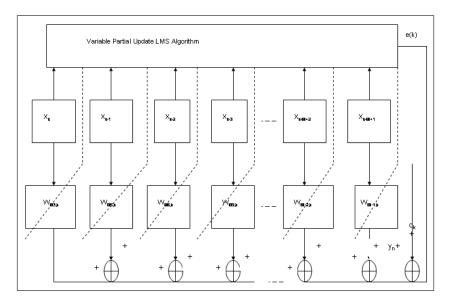


Figure 2: Transversal FIR Filter Structure for the Variable Partial Update-LMS Model

• Compute the output of the adaptive filter

$$y(n) = w_{nj}^T x_{nj} \tag{2}$$

where $j \in S$ and T is the Transpose operator.

• Compute the output error

$$e(n) = d(n) - y(n) \tag{3}$$

where $d(n) = x_n + K(n)$ is the desired output of the transfer filter

• Update the coefficient of the adaptive filter using Equation(1).

$$W_{n+1} = \begin{cases} W(n-1), j + \mu e(n) * X(n), j & \text{if } j \in S_i \text{ and } W \neg 0, \\ W_{(n-1),j} & \text{otherwise.} \end{cases}$$
(4)

• Update the step size of the adaptive filter

$$\mu(n+1) = \alpha \mu(n-1) + \mu(n)p(n)$$
(5)

$$\mu(n+1) = \begin{cases} \mu_{max} & \text{if } j\mu(n+1) > \mu_{max}, \\ \mu_{min} & \text{if } j\mu(n+1) < \mu_{min}, \\ \mu(n+1) & \text{otherwise.} \end{cases}$$
(6)

$$\gamma(n) = \begin{cases} \beta \gamma(n-1) & \text{if } j\mu(n+1) > \mu_{max}, \\ \gamma(n-1) & \text{otherwise.} \end{cases}$$
(7)

where γ and β are constant values, $0 < \gamma$, $\beta > 1$. The algorithm will adjust the parameter γ with the constant β . To ensure convergence the parameter β must satisfy that $0 < \beta < 1$

4 Simulation and Results

To test the performance of the proposed variable partial update LMS algorithm, we simulated the discrete signal sequence generated in Mat lab environment using pseudo random number generator with zero mean and variance of one, and a noise signal sequence which is obtained by introducing 0.8 noise level to the discrete signal. These signals form the desired signals that were input to the finite impulse response filter. The outputs from the filter form the actual signal which was subtracted from the desired signal to obtain the mean square error. To establish the superiority of the proposed partial update model over the full update model, training was performed using fifty different sets of input data with different value of the step size to obtain the average result of the mean square error (MSE) and the efficiency of the two models.

Table 1 shows the values of the step size and other simulation parameters. The full update LMS algorithm and our partial update LMS algorithm were simulated using various values of step size. Figure 3 shows the comparison of the full update and the partial update model

The proposed algorithm exhibit variable update of the step size. In order to test this algorithm for speed and accuracy of convergence, we used the parameters specified in Table 2 for the case when step size is less than the minimum step size (0.0036). Figure 4 shows a specimen comparison of the proposed partial update model with the equivalent full update model.

We illustrate the case when the step size is set to a large value or possibly larger than the maximum step size. Table 3 shows the parameters used for the simulation. In a normal situation, when the mean square error is magnified, stability of the filter is affected. This is

Step	μ_{max}	μ_{min}	γ	α	β	Length
Size Value						
0.001	1.9	0.001	0.0007	0.2	0.01	50
0.0011						
0.0013						
0.0017						
0.0020						
0.0030						

Table 1: Simulation Parameters for Fixed Step Size (μ)

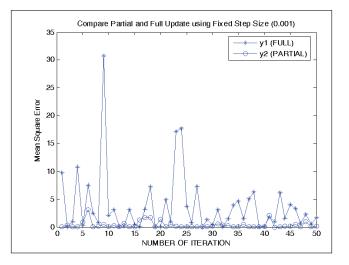


Figure 3: Mean Square Error versus Number of Iteration for Step Size of 0.001

Step Size Value	μ_{max}	μ_{min}	γ	α	β	Length
0.0036	0.09	0.01	0.0007	0.2	0.01	50
$0.0049 \\ 0.0052$						

Table 2: Simulation Parameter for Variable Step Size, for ($\mu < \mu_{min}$)

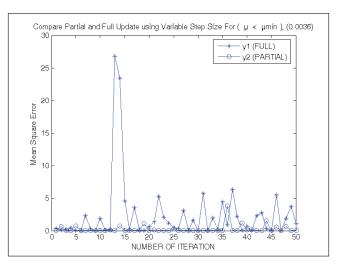


Figure 4: Mean Square Error versus Number of Iteration for $(\mu < \mu_{min})$

shown in Figure 5 where the mean square error was magnified and the accuracy of convergence was adversely affected for the full update model. However, the reverse is the case for the partial update model, with the speed and accuracy of convergence still enhanced.

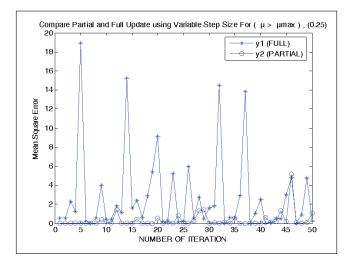


Figure 5: Mean Square Error versus Number of Iteration for $(\mu > \mu_{max})$

Step Size Value	μ_{max}	μ_{min}	γ	α	β	Length
0.25	1.9	0.001	0.0007	0.2	0.01	50
0.30						
0.40						

Table 3: Simulation Parameter for Variable Step Size for $(\mu > \mu_{max})$

The filter length was varied with a constant step size to evaluate the system performance. The filter length is equivalent to the number of coefficients required for the update of the filter. The parameter considered here is the memory load. Each of the coefficient requires a unit storage, therefore the lesser the number of coefficient the lesser the memory requirement. For the full update model, all the coefficients are used for the update process therefore the memory load (M) is equivalent to the filter length (L). In the case of partial update model, only the coefficients with index factors of three and five are used for the update. Therefore memory load $M = \frac{L}{3} + \frac{L}{5}$. Where $\frac{L}{3}$ and $\frac{L}{5}$ are the sets of filter length divisible by factors of three and five respectively. Simulation was carried out for fifty iterations using the parameters in Table 4. The result shown in Figure 6 revealed that the performance of the system is below 50% for full update model and up to 86% for the partial update model.

4.1 Computational Complexity

Computational complexity refers to the number of hardware resources required to implement the system. The complexity of an algorithm determines the hardware requirement and computational cost. The hardware required to implement the full update and the partial update finite impulse response transversal adaptive filter are the multiplier, summer and the memory. The multiplier is use to multiply the input with the corresponding weight, memory to store the weights, and summer to perform the addition. The computational complexity of our model was estimated by counting the number of hardware resources as described in [12] such as multipliers,

Filter	μ_{max}	μ_{min}	γ	α	β	Step size	Memory Load	
Length							Full Update	Partial Update
20	0.02	0.01	0.0007	0.2	0.01	0.003	20	9
50							50	23
70							70	33
90							90	42
100							100	47
130							130	61
150							150	70
170							170	79
200							200	93
250							250	126

 Table 4: Simulation Parameter for System Performance Evaluation

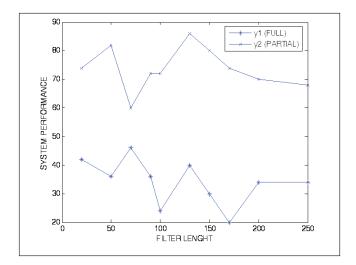


Figure 6: System Performance versus Filter Length

Filter	Memories	Summers	Multipliers	Memories	Summers	Multipliers
length	(FU)	(FU)	(FU)	(PU)	(PU)	(PU)
20	20	20	20	9	9	9
50	50	50	50	23	23	23
70	70	70	70	33	33	33
90	90	90	90	42	42	42
100	100	100	100	47	47	47
130	130	130	130	61	61	61
150	150	150	150	70	70	70
170	170	170	170	79	79	79
200	200	200	200	93	93	93
250	250	250	250	126	126	126

summers and memories required for a single iteration for each model as shown in Table 5. The evaluation result shown in Figure 7 revealed that the computational complexity of the proposed partial update model is considerably lower when compared with the full update model.

Table 5: Evaluation of Computational Complexity

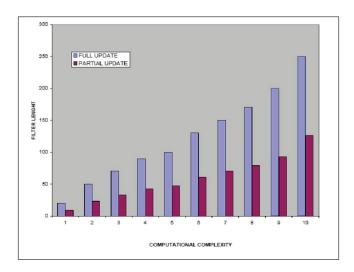


Figure 7: Filter Lengths versus Computational Complexity

5 Conclusions

To achieve the continuous update of the coefficient using adaptive algorithm, an improved computational model was proposed in this research, the novelty of which is the adoption of finite impulse response transversal adaptive filter to filter the noise signal from the transmitted signal. The adaptation process employed the concept of variable partial update least mean square algorithm. In the update process, only the coefficients with the factors of three and five are used. The performance of the proposed model was compared with the full update model using the following parameters convergence rate, memory load, and computational complexity in Mat lab environment. The simulation results revealed a better performance of the proposed model over the full update model. The proposed framework will particularly be suitable for wireless communication environment where the characteristics of the channel changes with time. The results obtained in the study will go a long way to reducing the effect of channel time variation, inter-symbol interference and co channel interference on the transmitted signal.

Bibliography

- C. S. Douglas, and W. Pan, Exact Expectation Analysis of the LMS Adaptive Filter, *IEEE Transaction on Signal Processing*, 43(12),2863-2871, 1995.
- [2] M. Godavarti and A. O. Hero, Analysis of the Sequential Partial Update LMS Algorithm, Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing, pp. 3857-3860, 2001.
- [3] S. Mullins and C. Heneghan, Alternative Least Mean Square Adaptive Filter Architectures for Implementation on Field Programmable Gate Arrays, Department of Electronic and Electrical Engineering, University College Dublin available at http://195.134.67.70/eurasip/proceedings/eusipco/2002/articles/paper389.pdf, 2002.
- [4] S. Singh, R.K. Bansal, and S. Bansal, Improved Channel Estimation with Auto-Regressive Prewhitening Techniques for Color Inputs, International Conference on Next generation Communication System: ICONGENCOM-06, pp. 9-14, 2006.
- [5] Y. Li and W. Xinan, A Modified VS LMS Algorithm, *IEEE Transaction on Signal Processing*, Vol. 2, No. 5, 615-618, 2005.
- [6] S. Diggavi, B. Chong, and A. Paulraj, An Interference Suppression Scheme with Joint Channel Data Estimation, *IEEE Journal on Communication*, Vol. 17, No.11,465-469, 1998.
- J. Sanubari, Fast Convergence LMS Adaptive Filters Employing FuzzyPartial Updates, *IEEE Transaction on Signal Processing*, Vol. 4, pp. 1334-1337, 2003.
- [8] S. Jo, J. Choi and Y. Lee, Modified Leaky LMS Algorithm for Channel Estimation in DS-CDMA System, *IEEE Communications Letters*, Vol. 6, No 5, 202-204 2002.
- [9] M. Gadhiok, Channel Estimation for Fast Fading Environment, *IEEE Transaction on Signal Processing*, Vol. 19, pp. 14-22, 2004.
- [10] K. Egiazarian, P. Kuosmanen, and R. C. Bilcu, Variable Step Size LMS Adaptive Filters for CDMA Multiuser Detection, Vol.17 pp. 259-264, 2004.
- [11] R. Bilcu, P. Kuosmanen and C. Rusu, A Novel Complementary Variable Step LMS Algorithm, *IEEE Transaction on Signal Processing*, Vol. 23 pp. 13-19, 2000.
- [12] V. Debrunner and D. Zhou, Hybrid Filtered Error LMS Algorithm: Another Alternative to Filtered-x LMS, *IEEE Transactions on Circuit and System*, Vol. 53, No. 3, pp.653-661, 2006.

Enhanced Security Protocol in Wireless Sensor Networks

T.C. Aseri, N. Singla

Trilok C. Aseri, Neha Singla

PEC University of Technology
Computer Science & Engineering Department
House No. 808, PEC Campus, Sector-12, Chandigarh-160012 (India)
E-mail: a_trilok_chand@yahoo.com, nehasingla1409@gmail.com

Abstract: The need for security in communications is in fact not new. This need has existed in military communications for thousands of years. In this paper, we focus on network protocols that provide security services. Wireless sensor network is an emerging technology that shows applications both for public as well as military purposes. Monitoring is one of the main applications. A large amount of redundant data is generated by sensor nodes. This paper compares all the protocols which are designed for security of wireless sensor network on the basis of security services and propose an improved protocol that reduces communication overhead by removing data redundancy from the network. By using the message sequence number we can check whether it is old message or new message. If the message is old then no need to send that message thereby reducing overhead. It also integrates security by data freshness in the protocol.

Keywords: data freshness, protocol, security, wireless sensor network.

1 Introduction

Sensor networks are typically data driven, i.e., the whole network cooperates in communicating data from sensors (information sources) to information sinks. Low-cost, low power, multifunctional sensor nodes that are small in size and communicate unterhered in short distances have been developed due to the recent advances in wireless communication. These tiny sensors have the ability of sensing, data processing, and communicating with each other. Wireless Sensor Networks (WSN) which rely on collaborative work of large number of sensors are realized. A WSN is a wireless network consisting of spatially distributed autonomous devices using sensors to cooperatively monitor physical or environmental conditions, such as temperature, sound, vibration, pressure, motion or pollutants, at different locations. In addition to one or more sensors, each node in a sensor network is typically equipped with a radio transceiver or other wireless communications device, a small microcontroller, and an energy source, usually a battery. A sensor network normally constitutes a wireless ad-hoc network, meaning that each sensor supports a multi-hop routing algorithm. Wireless sensor network is one of the most exciting and challenging research areas.

Nodes in sensor networks have restricted storage, computational and energy resources; these restrictions place a limit on the types of deployable routing mechanisms. Additionally, ad hoc routing protocols for conventional wireless networks support IP style addressing of sources and destinations. They also use intermediate nodes to support end-to-end communication between arbitrary nodes in the network. It is possible for any-to-any communication to be relevant in a sensor network; however this approach may be unsuitable as it could generate unwanted traffic in the network, thus, results the extra usage of already limited node resources. Many-to-one communication paradigm is widely used in regards to sensor networks since sensor nodes send their data to a common sink node for processing. This many-to-one paradigm also results in non-uniform energy drainage in the network.

The applications for WSNs are many and varied, but typically involve some kind of monitoring, tracking, and controlling, intelligent buildings, transportation, space exploration, disaster detection. In order to operate these applications successfully, it is necessary to maintain privacy and security of the transmitted data.

The rest of the paper is organized as follows: section 2 explains security requirements, section 3 presents a review of the relevant work, section 4 presents the proposed protocol, section 5 shows the results and discussion, and section 6 concludes the paper.

2 Security Requirements

2.1 Confidentiality

Confidentiality means keeping information secret from unauthorized parties. A sensor network should not leak sensor readings to neighboring networks. The confidentiality objective is required in sensors environment to protect information traveling between the sensor nodes of the network or between the sensors and the base station from disclosure, since an adversary having the appropriate equipment may eavesdrop on the communication. By eavesdropping, the adversary could overhear critical information such as sensing data and routing information.

2.2 Authentication

In a sensor network, an adversary can easily inject messages, so the receiver needs to make sure that the data used in any decision-making process originates from the correct source. As in conventional systems, authentication techniques verify the identity of the participants in a communication, distinguishing in this way legitimate users from intruders. In the case of sensor networks, it is essential for each sensor node and base station to have the ability to verify that the data received was really sent by a trusted sender and not by an adversary that tricked legitimate nodes into accepting false data. If such a case happens and false data are supplied into the network, then its behavior could not be predicted, and most of the times the mission of WSN will not be accomplished as expected.

2.3 Integrity

Data integrity ensures the receiver that the received data is not altered in transit by an adversary. Lack of integrity could result in many problems since the consequences of using inaccurate information could be disastrous, for example, for the healthcare sector where lives are endangered. Integrity controls must be implemented to ensure that information is not altered in any unexpected way.

2.4 Freshness

One of the many attacks launched against sensor networks is the message replay attack where an adversary may capture messages exchanged between nodes and replay them later to cause confusion to the network. Data freshness implies that the data is recent, and it ensures that an adversary has not replayed old messages. To achieve freshness, network protocols must be designed in a way to identify duplicate packets and discard them preventing potential mix-up.

2.5 Availability

Availability ensures that services and information can be accessed at the time they are required. In sensor networks there are many risks that could result in loss of availability such as sensor node capturing and denial of service attacks. The availability of a sensor and sensor network may decrease for the following reasons [1]:

- Additional computation consumes additional energy. If no more energy exists, the data will no longer be available.
- Additional communication also consumes more energy. Besides, as communication power increases so does the chance of a communication conflict or interference. A single point failure exists if we use the central point scheme such as a single sink or gateway. This greatly threatens the availability of the network.

3 Related Work

The various protocols which have been proposed for security in wireless sensor network by various authors are SPIN, LEAP, TINYSEC, ZIGBEE, SM. In SPIN (Sensor Protocols for Information via Negotiation), nodes use three types of messages ADV, REQ and DATA to communicate. ADV is used to advertise new data, REQ to request for data and DATA is the actual message itself. The protocol starts when a SPIN node obtains new data that it is willing to share. It does so by broadcasting an ADV message containing meta-data. If a neighbor is interested in the data, it sends an REQ message for the DATA and the DATA is sent to this neighbor node. The neighbor sensor node then repeats this process to its neighbors as a result of which the entire sensor area will get a copy. It consists of two secure building blocks SNEP (Sensor Network Encryption Protocol) and ÎźTESLA (Timed Efficient Stream Loss-tolerant Authentication). In addition to integrity, SNEP is used to provide confidentiality through encryption and authentication using a message authentication code (MAC). It lowers communication overhead adding only 8 bytes per message [2]. TESLA authenticates the initial packet using the digital signature. For an authenticated packet to be sent, the base station computes a MAC on the packet with the key that is secret at that point in time. When a node gets a packet, it can confirm that the base station did not yet disclose the corresponding MAC key [3].

The goal of LEAP (Localized Encryption and Authentication Protocol) is to satisfy the security properties of authentication and confidentiality in a wireless environment where the intruder may eavesdrop, inject packets, and replay messages [4]. LEAP, as a key management protocol for sensor networks, is designed to support in-network processing, while restricting the impact of a compromised node to the network. In order to support the in-network processing necessary for most applications of these networks while at the same time providing security properties, such as security and authentication, similar to those of pairwise symmetric keys, LEAP specifies four types of keys: individual keys, pairwise shared keys, cluster keys and group keys. Individual keys are symmetric keys shared between the base station and each of the nodes. For example, a node might use the individual key to notify the base station of a suspicious neighbor. Pairwise shared keys are symmetric keys shared between a node and each of its neighbors. While pairwise shared keys are used to establish cluster keys, they prevent passive participation which is desirable for in-network processing. Cluster keys are symmetric keys shared between a node and all of its neighbors. These cluster keys can be used for locally broadcast messages such as a routing protocol might use and are also used for updating the group key. The group key, a symmetric key shared between the base station and all of the nodes, allows encrypted and authenticated messages to broadcast through the whole network.

In the next protocol, TINY SEC, the dominant traffic pattern in sensor networks is many-toone, with many sensor nodes communicating sensor readings or network events over a multihop topology to a central base station. However, neighboring nodes in sensor networks often witness the same or correlated environmental events, and if each node sends a packet to the base station in response, precious energy and bandwidth are wasted. To prune these redundant messages to reduce traffic and save energy, sensor networks use in-network processing such as aggregation and duplicate elimination [5, 6]. Since in-network processing requires intermediate nodes to access, modify, and suppress the contents of messages, it is unlikely we can use end-to-end security mechanisms between each sensor node and the base station to guarantee the authenticity, integrity, and confidentiality of these messages. With authenticated encryption, TinySec encrypts the data payload and authenticates the packet with a MAC [7]. Single shared global cryptographic key, link layer encryption and integrity protection cryptography is based on a block cipher. TinySec is a research platform that is easily extensible and has been incorporated into higher level protocols.

In ZIGBEE, the concept of a Trust Center is introduced in the specification. Generally the ZigBee coordinator performs this duty. This trust center allows other devices to join the network and also distributes the keys. There are three roles played:

- trust manager, whereby authentication of devices requesting to join the network is done,
- network manager, maintaining and distributing network keys, and
- configuration manager, enabling end-to-end security between devices [8].

It operates in both Residential Mode and Commercial Mode. The Trust Center running Residential Mode is used for low security residential applications. Commercial Mode is designed for high-security commercial applications. In Residential Mode, the Trust Center will allow devices to join the network, but does not establish keys with the network devices. It therefore cannot periodically update keys and allows for the memory cost to be minimal, as it cannot scale with size of the network. In commercial mode, it establishes and maintains keys and freshness counters with every device in the network, allowing centralized control and update of keys. This results in a memory cost that could scale with the size of the network. There are three types of keys employed, the Master Key, the Link Key and the Network Key. Master keys are installed first, either in the factory or out of band. They are sent from the Trust Center and are the basis for long-term security between two devices. The Link key is a basis of security between two devices and the Network keys are the basis of security across the entire network. Link and Network keys, which are either installed in the factory or out of band, employ symmetrical keykey exchange (SKKE) handshake between devices. The key is transported from the Trust Center for both types of keys. This operation occurs in commercial mode, as residential mode does not allow for authentication.

In the latest protocol, SM (Security Manager), a new method of key agreement has been proposed in [9], whereby, when a new device joins a network, the Security Manager (SM) gives static domain parameters at the base station such as the order of the curve and the elliptic curve coefficients. After calculating a public key using the base point and a private key, the device sends a public key to the SM. Therefore the SM would have the public key list for all the devices in the network. Authentication is achieved by using either Diffie-Hellman or Elliptic Curve Equation. Confidentiality is achieved by using message authentication protocol. This shows that SM protocol offers more services than the other existing protocols.

4 Proposed Protocol

A security protocol refers to a set of rules governing the interaction between peer processes to provide a certain type of security service. We propose a new security protocol. Security Manager (SM) [9] does not guarantee data freshness; and so, we suggest a protocol to make up for the weakness of SM. This provides a solution to maintain data freshness by checking the message sequence number. When the message is sent, it is checked by the message sequence number whether it is already sent or not. If the message is old then no need to send that message. By this way, we can reduce overhead. By reducing the overhead we can make the protocol more efficient. Authentication and confidentiality are also provided. Authentication is defined to provide assurance about the originator of a message. This prevents an attacker from mimicking the operation of another device in any attempt to compromise the network.

Confidentiality means keeping information secret from unauthorized parties. The standard solution to keep sensitive data secret is to encrypt the data with a secret key that only the intended receivers possess, hence achieving confidentiality.

Additionally, this protocol provides freshness through the use of freshness checks. These checks prevent replay attacks, as devices maintain incoming and outgoing messages. Whenever a node wants to send message, message sequence number will be checked.

One of the many attacks launched against sensor networks is the message replay attack where an adversary may capture messages exchanged between nodes and replay them later to cause confusion to the network. Data freshness implies that the data is recent, and it ensures that an adversary has not replayed old messages. To achieve freshness, network protocols must be designed in a way to identify duplicate packets and discard them preventing potential mix-up. This extra feature shows that proposed protocol offers more security services than the existing one.

In the proposed protocol, time interval Δt accounts for request time, response time, delay, as shown in (1) and (2). Freshness is computed as the difference between the time a data item is generated and the time it is received at the sink.

With data freshness,

$$\Delta t = t_{Rq} + t_{Rs} + \Delta d + t \tag{1}$$

Without data freshness,

$$\Delta t = t_{Rq} + t_{Rs} + \Delta d \tag{2}$$

Where t_{Rq} = request transmission, t_{Rs} = response transmission, Δd = sum of delays (delays in transmission and propagation), t = time when message sequence number is checked. In case of without data freshness, no message is checked.

Accuracy is measured as the ratio of total number of messages received at the sink to the total number of messages generated. In case of without data freshness, no message is discarded but in case of with data freshness the message, which is already sent then no need to send again. Therefore, the ratio of without data freshness is always one but it is less than one in case of with data freshness. We can calculate the packet ratio by (3) and (4) as follows:

Without data freshness,

$$PacketRatio = send/receive = 1 \tag{3}$$

With data freshness,

$$PacketRatio = send/receive < 1 \tag{4}$$

5 Results and Discussion

Sensor network is a promising and upcoming technology with usage in important applications. The resource constraint hardware, specialized software, low energy devices and hostile environment makes the security in wireless sensor networks a challenging task as and when compared to the traditional computer networks.

Energy efficiency can be achieved by reducing the number of packets transmitted. Without data freshness, each node will send a packet that will be forwarded to the sink whereas with data freshness no need to send all packets, this reduced number of packets transmitted improve the efficiency. Figure 1 shows the average latency. In case of with data freshness, average latency is constant as the number of nodes is increased while it is increased in case of without data freshness.

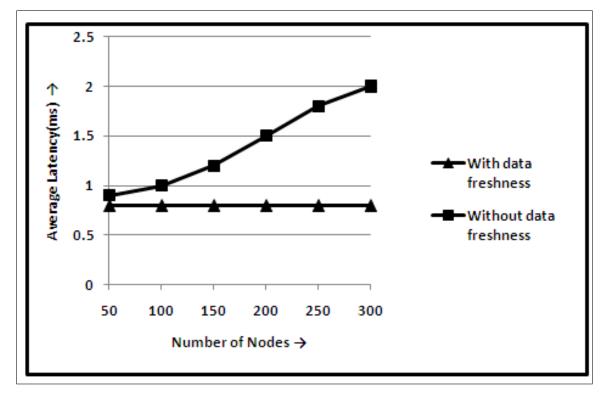


Figure 1: Number of nodes vs. Average Latency

Figure 2 shows the average packet delivery ratio. In case of with data freshness, average packet ratio is 100% because how many packets sent will definitely be received but in case of without data freshness some may already sent. This reduced number of packet transmitted also improves the efficiency.

The discussion of the security protocol and authentication mechanism allow for the construction of comparison table as in given Table 1, where they can be compared under similar headings. It can be seen from the table that new protocol is better than the existing protocol and offers more security services than the earliest one.

6 Conclusion

In this paper, firstly we propose a new security protocol for wireless sensor network. Secondly, we compared the performances of all the existing protocol with proposed protocol. SPIN

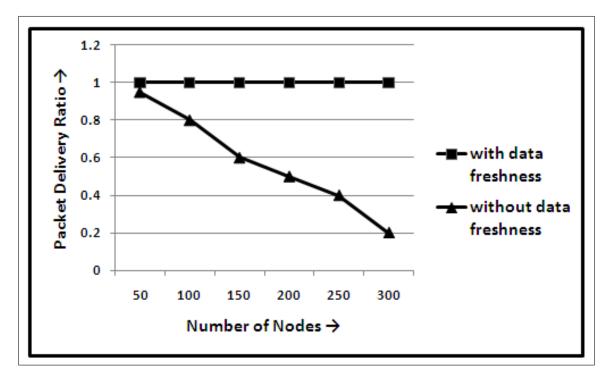


Figure 2: Number of nodes vs. Packet Delivery Ratio

Protocol / Service	С	F	Ι	Ava	IA	А
SPIN	YES	YES	YES	NO	YES	NO
LEAP	YES	NO	NO	NO	YES	NO
TINYSEC	YES	NO	NO	-	YES	YES
ZIGBEE	YES	YES	YES	NO	YES	YES
SM	YES	NO	NO	-	YES	YES
OUR PROTOCOL	YES	YES	-	-	YES	YES

Table 1: Security Architecture Comparison

C=Confidentiality, F=Freshness, I=Integrity, Ava=Availability, IA=Implicit Authentication, A=Authentication of user

was found to perform better in smaller size networks because of its efficiency and high latency properties. The use of SPIN in large scale networks could potentially exhaust system resources in a much faster pace. Our protocol has one extra feature i.e. freshness. Freshness reduces the overhead. This extra feature shows that this is superior to the existing protocols. This also improves the efficiency.

Bibliography

- [1] J.P. Walters, Zh. Liang, W. Shi, V. Chaudhary, Security in Distributed, Grid, and Pervasive Computing, Chapter 17, CRC Press, 2006.
- [2] A. Perrig, R. Szewczk, J.D. Tygar, V. Wen, D.E. Culler, SPINS: Security Protocols for Sensor Networks, *Wireless Networking*, Vol. 8, No. 5, pp. 521-534, Sept 2002.

- [3] A. Perrig, R. Canneti, J. D. Tygar, D. Song, The TESLA Broadcast Authentication Protocol, *CryptoBytes*, Vol. 5, No. 2, pp. 2-13, 2002.
- [4] D. Boyle, T. Newe, Security Protocols for use with Wireless Sensor Networks: A Survey of Security Architectures, Proceedings of the 3rd International Conference on Wireless and Mobile Communications, Guadeloupe, French Caribbean, pp. 54, 04-09 March 2007.
- [5] S. Madden, M.J. Franklin, J.M. Hellerstein, W. Hong, TAG: a Tiny Aggregation Service for Ad-Hoc Sensor Networks, Proceedings of the 5th Symposium on Operating System Design and Implementation (OSDI), Boston, Massachusetts, USA, pp. 131-146, 09-11 December 2002.
- [6] S. Madden, R. Szewczyk, M.J. Franklin, D. Culler, Supporting Aggregate Queries Over Ad-Hoc Wireless Sensor Networks, *Proceedings of the 4th IEEE Workshop on Mobile Computing* and Systems Applications (WMCSA), Callicoon, NY, USA, pp. 49-58, 20-21 June 2002.
- [7] C. Karlof, N. Sastry, D. Wagner, TinySec: A Link Layer Security Architecture for Wireless Sensor Networks, Proceedings of the 2nd ACM International Conference on Embedded Networked Sensor Systems, Baltimore, MD, USA, pp. 162-175, 03-05 November 2004.
- [8] ZigBee Alliance, ZigBee Security Specification Overview, [Online] Available: http://www.zigbee.org/en/events/documents/december2005_open_house_presentations/ zigbee_security_layer_technical_overview.pdf.
- [9] J. Heo, C.S. Hong, Efficient and Authenticated Key Agreement Mechanism in Low-Rate WPAN Environment, Proceedings of the 1st IEEE International Symposium on Wireless Pervasive Computing, Physics, Thailand, pp. 1-5, 16-18 January 2006.

Perturbation in Population of Pulse-Coupled Oscillators Leads to Emergence of Structure

A. Bartha, D. Dumitrescu

Attila Bartha

Babes-Bolyai University of Cluj-Napoca Romania, 400084 Cluj-Napoca, 1 Mihail Cogalniceanu E-mail: abartha@cs.ubbcluj.ro

D. Dumitrescu

Babes-Bolyai University of Cluj-Napoca Romania, 400084 Cluj-Napoca, 1 Mihail Cogalniceanu E-mail: ddumitr@cs.ubbcluj.ro

> **Abstract:** A new synchronization model based on pulse-coupled oscillators is proposed. A population of coupled oscillators is represented as a cellular automaton. Each cell periodically enters a firing state. Firing of a cell is sensed by other cells in a neighborhood of radius R. As a result the sensing cell may change its firing rate. The interaction strength between a firing and a sensing cell decreases with the squared distance between the two cells. For most starting conditions waves of synchronized firing cells emerge. Simulations indicate that for certain parameter values the emergence of synchronization waves occurs only if there is dispersion in the intrinsic firing frequencies of the cells. Emergence of synchronization waves is an important feature of the model.

> **Keywords:** Cellular Automaton, Synchronization, Emergence, Pulse-Coupled Oscillators

1 Introduction

Many physical and biological systems can be described by mathematical models of coupled oscillators. Examples include earthquake formation [1], synchronously firing neurons [2], synchronous flashing of fireflies [3] and heart pacemaker cells [4].

Pulse-coupled oscillators are a special case when the interaction between oscillators is pulsatile. Each oscillator can enter a firing state when emits a pulse of physical signal. This signal is received by other oscillators in the population and as a result they may change their own oscillation frequency.

A model of pulse-coupled oscillators in the area of cardiology has been first described by Peskin [4]. The model has proved itself to be applicable in many other areas and was extensively studied [5,6]. Many variations of the model have been considered, including different coupling mechanisms and topolgies.

A population of pulse-coupled oscillators on a 2D grid is considered. Changes in each oscillator occur at discrete time moments. Oscillators (grid cells) interact only if the distance between them is less than a given radius. The interaction strength decreases with the squared distance between the interacting cells.

In many models of pulse-coupled oscillators a total synchronization of the population is eventually observed. However, for the proposed model, extensive simulations indicate that a total synchronization of the population does not occur. Instead, waves of synchronously firing oscillators emerge (see section 3). This is an interesting feature of the model that indicate a good potential for applications.

2 Proposed Oscillator Population Model

2.1 Grid model

A population of pulse-coupled oscillators arranged on a 2-dimensional grid is considered. The oscillator's states are synchronously updated at regular intervals, called generations. The structure of each oscillator consist of a classical integrate-and-fire mechanism combined with a separate mechanism for sensing and integrating the signal emitted by neighbouring oscillators.

A classical integrate-and fire oscillator consists of a unit which integrates an input signal until the output value of the integrator reaches a certain threshold. Then the oscillator enters firing state, the integrator discharges and the process repeats by itself. The time elapsed between two consecutive firing states is called the firing period of the oscillator. The discharge time is usually considered negligable compared to the firing period.

The firing period of the oscillator can be changed in at least two ways: by changing the slope of the integrator or by changing the firing threshold.

In the proposed model the second approach is considered. The integrator has a constant slope (which is a parameter of the oscillator). It generates a signal which increases proportionally with time (the number of generations in our discrete model).

A separate mechanism in each oscillator is responsible for sensing the signals emitted by other oscillators. The energy of the detected signals are integrated by a leaky integrator, i.e. the accumulated energy of the integrator decays exponentially in time.

An oscillator meats the firing condition if the output value of the integrate-and-fire mechanism summed with the weighted value of the sensing integrator is greater than one. This is equivalent to a classic integrate-and-fire oscillator with a non-constant firing threshold: at each moment the firing threshold is one minus the (weighted) value of the sensing integrator.

A particular feature of the model is that the intensity of the signal emitted by a firing oscillator decreases with the squared distance. This models the behaviour of a physical signal (like light or sound) propagating in the three-dimensional space.

The proposed model limits the interaction distance between two cells to a certain value called the interaction radius. This corresponds to a sensitivity threshold of a cell.

2.2 Oscillator state and update rules

The internal state of an oscillator is represented by a set of three variables (P, S, F) with the following meaning :

- P, called potential, represents the current value of the constant-slope integrator;

- S, called signal, represents the current value of the sensing integrator;

- F, called firing state, is a boolean value which is true when the oscillator is in firing state in the given generation.

At each generation the state variables of each oscillator are updated according to the following rules:

The potential P is incremented with a constant value E. The value of E is an oscillator-specific parameter and determines the slope of the integrator.

$$P <= P + E \tag{1}$$

The signal S is calculated by summing the intensities of the signals emitted by all oscillators in the interaction radius which are in firing state, then decrementing S by an amount proportional to the value of S in the previous generation:

$$S <= \sum_{k} \frac{1}{d_k^2} \tag{2}$$

where d_k is the distance between the oscillator and a neighbouring oscillator k in the interaction domain $d_k \leq R$ which is in firing state and D is an oscillator-specific decay parameter.

The firing state F of an oscillator is determined by comparing the weighted summ of the potential P and signal S with a constant threshold:

 $F \le true$ if $P + KS \ge 1$ $F \le false$ otherwise

where K is an oscillator-specific coupling parameter.

A cell is in firing state for only one generation. In the next generation returns to normal state and the potential P is reset to zero:

 $P \ll 0$ if F = true in the previous generation

The model implies that if the oscillator has no firing neighbors inside the interaction radius R than it will fire periodically with a firing rate determined by the parameter E.

The potential is incremented in each generation with E and when it becomes grater than 1, the oscillator enters firing state for one generation. On the next generation the potential is reset to zero and the process repeats by itself.

However, if there are firing oscillators in the neigbourhood of the oscillator then the 'radiated energy' of the firing increases the signal value of the oscillator. This value, multiplied by the coupling parameter K is added to the potential and together determine the firing condition of the oscillator. Thus, if many oscillators fire in the neigbourhood of the oscillator, the oscillator will fire more frequently.

3 Experimental Results

The proposed automaton model is simulated for different values of the parameters E, D and K. Different interaction radii R are also considered.

The potential P of each oscillator is initialized to a random value uniformly distributed in the interval [0, 1). The starting value of the signal S is set to zero for each oscillator.

In each experiment the decay and the coupling parameters D and K are set to the same value for each cell.

In some experiments the parameter E which determines the intrinsic firing frequency of the cell is chosen to have the same value. In other experiments the E value is peturbed by a small random noise, so that each cell is initialized with different values of E.

In experiments where the parameter is the same for all cells, the evolution of the cell array displays a random pattern of firing cells. This pattern is maintained along more than 100,000 generations.

However, if the value of the parameter E is perturbed with a small noise, in a few hundred generations clusters of firing cells emerge. Waves of synchronously firing cells emerge tipically after 600-1000 generations.

Figure 1 depicts a typical evolution of the automaton on a grid of 150 x 150 cells.

Parameter values are E = 0.05, D = 0.7, K = 0.1 and R = 30. The initial value of parameter E for each cell is perturbed by a random value in the range [0, 0.005] around the central value of 0.05. In the diagrams of Figure 1 the potential values of the cells are displayed in red color with intensity proportional to the potential. Cells in the firing state are indicated in a bright color.

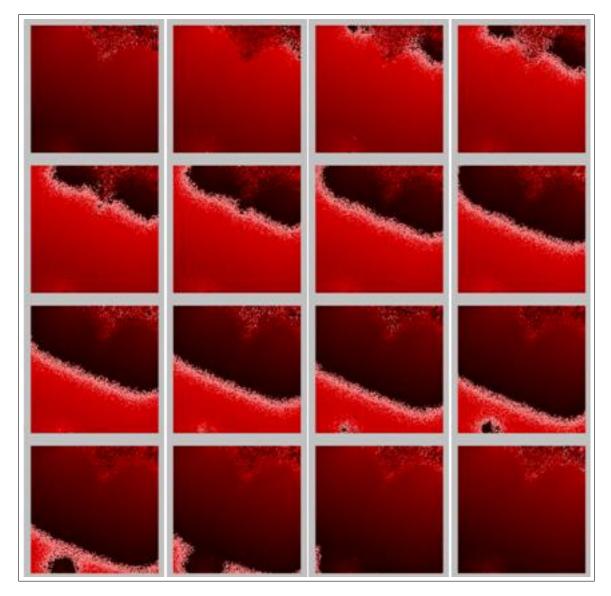


Figure 1: A typical evolution of the oscillator population. Parameter setting: E = 0.05, D = 0.7, K = 0.2, R = 30. The variance of the E parameter was set to 0.005. The sequence of images represents 16 consecutive states of the oscillator array for generations 2400 to 2415. The emergence of a synchronized front of firing oscillators can be observed in this sequence. This front travels through the oscillator array and disappears at the array boundary. After 16 generations the oscillator array returns to a similar state and the process repeats itself.

4 Conclusions

A new model of pulse-coupled oscillators is proposed and investigated. A population of pulsecoupled oscillators on a 2D grid is considered. Changes in each oscillator occur at discrete time intervals. The interaction strength decreases with the squared distance between interacting cells. The interaction distance is limited by a maximum radius. A total synchronization of the population is eventually observed in many models of pulse-coupled oscillators. Simulations indicate that total synchronization of the population does not occur in the propsed model. Instead some synchronization patterns emerge. These emergent patterns are observed as waves of synchronously firing oscillators. The condition of emergence is the existence of a small perturbation in the intrinsic frequencies of individual oscillators. If all oscillators have the same intrinsic frequency then the synchronization patterns do not emerge. The proposed model describes a new auto-organization paradigm conditioned by the presence of a noise mandatory in inducing an emergent synchronization process. Therefore the model provides a way of understanding self-organization in some natural or artificial systems.

Acknowledgment

This work was supported by CNCSIS - UEFISCSU, project number PNII - IDEI 508/2007. Investing in people! PhD scholarship, Project co-financed by the European Social Fund, SEC-TORAL OPERATIONAL PROGRAMME HUMAN RESOURCES DEVELOPMENT 2007-2013 Babeş-Bolyai University, Cluj-Napoca, Romania

Bibliography

- Herz, A. V. M., Hopeld, J. J.: Earthquake cycles and neural reverberations: Collective oscillations in systems with pulse-coupled threshold elements. *Phys. Rev. Lett*, 75(6), 1222-1225 (Aug. 1995)
- [2] Izhikevich, E. M.: Weakly pulse-coupled oscillators, FM interactions, synchronization, and oscillatory associative memory. IEEE Trans. on Neural Networks, 10(3), 508-526 (May, 1999)
- [3] Buck, J., Buck, E., Case J., Hanson, F.: Control of repteroptyx cribellata. J. Comp. Physiology A, 144(3), 630-633 (Sep, 1981)
- [4] Peskin C.: Mathematical aspects of heart physiology. Courant Institute of Mathematical Sciences, New York University (1975)
- [5] Mirollo, R.E., Strogatz, S.H.: Synchronization of pulse-coupled biological oscillators, SIAM Journal on Applied Mathematics, 50, 1645-1662 (1990)
- [6] Strogatz, S.H.: From Kuramoto to Crawford: exploring the onset of synchronization in populations of coupled oscillators. Physica D 143, 1-20 (Sep, 2000)

A Time-Bound Ticket-Based Mutual Authentication Scheme for Cloud Computing

Z. Hao, S. Zhong, N. Yu

Zhuo Hao

 University of Science and Technology of China Department of Electronic Engineering and Information Science Hefei, Anhui 230027, P.R.China, and
 State University of New York at Buffalo Department of Computer Science and Engineering 201 Bell Hall, Amherst, NY 14260, USA E-mail: hzhuo@mail.ustc.edu.cn

Sheng Zhong

State University of New York at Buffalo Department of Computer Science and Engineering 201 Bell Hall, Amherst, NY 14260, USA E-mail: szhong@buffalo.edu

Nenghai Yu

University of Science and Technology of China Department of Electronic Engineering and Information Science Hefei, Anhui 230027, P.R.China E-mail: ynh@ustc.edu.cn

> **Abstract:** Cloud computing is becoming popular quickly. In cloud computing, people store their important data in the cloud, which makes it important to ensure the data integrity and availability. Remote data integrity checking enables the client to perform data integrity verification without access to the complete file. This service brings convenience to clients, but degrades the server's performance severely. Proper schemes must be designed to reduce the performance degradation.

> In this paper, a time-bound ticket-based mutual authentication scheme is proposed for solving this problem. The proposed authentication scheme achieves mutual authentication between the server and the client. The use of timebound tickets reduces the server's processing overhead efficiently. The correspondence relationship between the digital ticket and the client's smart card prevents user masquerade attack effectively. By security analysis, we show that the proposed scheme is resistant to masquerade attack, replay attack and password guessing attack. By performance analysis, we show that the proposed scheme has good efficiency. The proposed scheme is very suitable for cloud computing.

> **Keywords:** cloud computing, mutual authentication, digital ticket, masquerade attack.

1 Introduction

Cloud Computing [1,2] has become a very popular technology. By using cloud computing, both the data storage and the computation resources are moved from personal computers into the cloud. Recently many companies provide the cloud storage services, including Amazon Simple

Storage Service (S3) [3], Microsoft SkyDrive [4], Nirvanix CloudNAS [5], etc. As people store their important data in the cloud without keeping a local copy, it is important for the cloud clients to be able to verify the data integrity and availability.

Remote data integrity checking protocols [6–10] are proposed to ensure the data integrity in the cloud. Ateniese et al. [6] propose the technology of provable data possession to enable the client to verify the data integrity without retrieving it from the server. Juels and Kaliski [7] propose the scheme called proofs of retrievability to enable the server to produce a concise proof about the data availability. Other schemes [8–10] have also been proposed with increased scalability and efficiency. However, as the clients of cloud services are numerous and are also increasing very quickly, these technologies bring a lot of extra computation and storage overhead to the server. Mechanisms should be designed to reduce these overhead as much as possible.

One possible solution is to limit the client's data verification frequency. In this method, the server releases a certain number of digital tickets to the client at a constant frequency, e.g., 12 tickets per year. The client uses one ticket for one time of data verification, and after that, the used ticket becomes invalid and cannot be reused. If the client needs more tickets, she must purchase them from the server. On the other hand, the client's identity should be properly authenticated to the server before she uses these services.

Recently Lin and Chang [11] propose a countable and time-bound password-based user authentication scheme, which is a possible method for solving this problem. In the Lin-Chang scheme, the tickets $\{TK_1, TK_2, ..., TK_t\}$ are generated by the server in a way that the authentication value in TK_j is one of the square roots of the value in TK_{j+1} . The Lin-Chang scheme is secure against masquerading attack and replay attack, and also achieves good efficiency at the client side. However, the Lin-Chang scheme has several disadvantages, which make it not suitable in this environment. Firstly, the ticket in Lin-Chang scheme is not associated with the client's identity, which allows anyone who obtains the ticket to be able to use it. Secondly, the client's tickets can only be used sequentially. Whenever TK_j expires, the tickets $TK_1, TK_2, ..., TK_{j-1}$ cannot be used anymore. Thirdly, the Lin-Chang scheme involves high-cost modular exponentiation operations at the server side, which bring large amount of overhead as the quantity of cloud clients increases quickly. Finally, the Lin-Chang scheme does not achieve mutual authentication.

In this paper, we propose a new ticket-based mutual authentication scheme which has improvements in these aspects. Firstly, we use smart card in the mutual authentication scheme. Each client has her own unique smart card. The client's tickets are associated with her smart card, so that even when her tickets are lost, they cannot be used by other clients. Secondly, in the proposed scheme, tickets are relatively independent of one another, so that the tickets need not be used sequentially. When one ticket expires, all other tickets that do not expire can still be used. Thirdly, the proposed scheme involves only lightweight exclusive-or and hash computations, which makes it very efficient at both the client side and the server side. Finally, the proposed scheme achieves mutual authentication, so that both the server and the client authenticates each other when performing the data verification. By security analysis, we show that the proposed scheme is secure against lost smart card attack, lost ticket attack, masquerade attack, replay attack, etc. By performance analysis, we show that the proposed scheme achieves good performance at both the server side and the client side. The proposed scheme is suitable for cloud computing.

The rest of this paper is organized as follows. Section 2 presents the system model and introduces common notations used throughout this paper. In section 3, the proposed timebound ticket-based mutual authentication scheme is presented. In section 4, we show that the proposed scheme is secure against adversary's attacks. In section 5, we show that the proposed scheme is cost-efficient. Finally, we conclude in section 6.

$\mathbf{2}$ System Model and Common Notations

In the proposed scheme, there are two different entities: the cloud server and the client. The cloud server provides data storage services to a lot of clients, and also provides data integrity verification services to the clients. Furthermore, the cloud server is in charge of the client registration, client authentication and digital ticket management. Clients put their data at the server and retrieves the data on demand. Each client has a unique identification and a password with which she can prove her identity to the server. In addition, similar to the Lin-Chang scheme [11], a bulletin board is maintained by the server for preventing repeated usage of one ticket.

For convenience of description, we denote the cloud server by S and the clients by $\{U_i, i =$ $1, 2, 3, \dots$ Denote the identification and password of U_i by ID_i and PW_i . We use a cryptographic hash function h() and a keyed hash function $h_K()$, with K as a cryptographic key. S has two long-term secret keys, denoted by K_1 and K_2 .

3 The Proposed Mutual Authentication Scheme

In this section, we present the proposed time-bound ticket-based mutual authentication scheme. The proposed scheme consists of 4 phases: registration phase, verification request phase, mutual authentication phase and password change phase.

3.1**Registration Phase**

When the client U_i initially registers herself to S, the registration phase is invoked. The registration phase consists of the following steps:

1. U_i selects her own ID_i , her password PW_i and a random number b. Then U_i computes $IPB_i = H(ID_i||H(PW_i \oplus b))$ and sends $\{ID_i, IPB_i, t\}$ to the server, in which t is the number of digital tickets U_i needs. At the same time U_i pays the corresponding ticket fee to S. The ticket fee can be payed in the form of either real or virtual currency, which depends on the policy of S.

2. After S receives the message and ticket fee from U_i , S generates t tickets for U_i . Denote the *j*th ticket of U_i by $T_i^{(j)}$. Denote by $TID_i^{(j)}$ and $VP_i^{(j)}$ the ticket ID and valid period of $T_i^{(j)}$ respectively. Specifically, *S* generates $\{(TID_i^{(j)}, VP_i^{(j)}), j = 1, 2, ..., t\}$ and computes as follows:

$$W_i = IPB_i \oplus H(ID_i, K_1),$$
$$\alpha_i^{(j)} = H_{K_2}(ID_i \parallel TID_i^{(j)} \parallel VP_i^{(j)}),$$
$$\beta_i^{(j)} = \alpha_i^{(j)} \oplus IPB_i.$$

 $T_i^{(j)}$ has two parts:

$$T_i^{(j)} = (T_i^{(j)_1}, T_i^{(j)_2}),$$

in which

$$\begin{split} T_i^{(j)_1} &= (TID_i^{(j)}, \ VP_i^{(j)}), \\ T_i^{(j)_2} &= \beta_i^{(j)}. \end{split}$$

S also computes $Z_i = H_{K_2}(ID_i) \oplus IPB_i$, which is used during the password change phase.

3. S writes ID_i , t, W_i , Z_i and $T_i^{(j)}$, j = 1, 2, ..., t into a smart card and sends it to U_i . 4. When U_i receives her smart card, she writes b into the card.

3.2 Verification Request Phase

As the client receives t tickets, she can use these tickets to perform data verification at most t times. The description below assumes that this is the kth verification request that U_i invokes.

- 1. U_i inserts her smart card into a card reader and enters ID_i and PW_i .
- 2. The smart card generates a nonce r_U according to system time, and computes

$$IPB_{i} = H(ID_{i}||H(PW_{i} \oplus b)), H_{i} = W_{i} \oplus IPB_{i}$$
$$C_{1} = r_{U} \oplus H_{i}, C_{2} = H(r_{U}) \oplus T_{i}^{(k)_{2}} \oplus IPB_{i}.$$

Then U_i 's smart card sends $\{ID_i, T_i^{(k)_1}, C_1, C_2\}$ to S.

3.3 Mutual Authentication Phase

When S receives the verification request message from U_i , it executes the following steps:

1. S checks the validity of ID_i and rejects the service request if ID_i is invalid.

2. S checks whether the ticket ID $TID_i^{(k)}$ is on the bulletin board. If it is on the bulletin board, then S rejects U_i 's service request and terminates the process.

3. S checks whether the current date is in the range of $VP_i^{(k)}$ or not. If not, then S rejects U_i 's service request and terminates the process.

4. S computes $D_0 = H(ID_i, K_1), D_1 = C_1 \oplus D_0$ and $D_2 = H(D_1) \oplus C_2$.

5. S computes $H_{K_2}(ID_i \parallel TID_i^{(k)} \parallel VP_i^{(k)})$ and checks whether it is equal to D_2 . If they are not equal, then S rejects U_i 's service request and terminates the process. Otherwise, S authenticates U_i successfully.

6. S generates a random nonce r_S , computes $C_3 = D_0 \oplus r_S$ and $C_4 = H(r_U, r_S)$ and sends $\{C_3, C_4\}$ to U_i . S also computes $K_S = H(D_0, r_U || r_S)$, which will be used as a subsequent session key.

7. U_i 's smart card computes $D_3 = C_3 \oplus H_i$. Then it compares $H(r_U, D_3)$ with C_4 . If they are equal, U_i authenticates S successfully. Then the smart card computes $K_C = H(H_i, r_U || r_S)$ and uses K to communicate with S in the subsequent data verification process. Note that $K_C = H(H_i, r_U || r_S) = H(H(ID_i, K_1), r_U || r_S) = K_S$. When the kth data verification is finished, U_i deletes the ticket $T_i^{(k)}$ from its smart card, and S publishes $TID_i^{(k)}$ to its bulletin board. Finally the smart card deletes the nonce r_U , and S deletes r_S from its memory, to prevent replay attack.

8. After the mutual authentication, the data verification is performed. The data verification schemes [6–10] are independent of the proposed authentication scheme, so we don't describe the details. The keys K_C and K_S can be used for secret communication.

3.4 Password Change Phase

This phase is invoked when U_i needs to change her password PW_i to a new one. The password change phase consists of the following steps:

1. U_i inserts her smart card into a card reader, and enters ID_i and PW_i .

2. The smart card generates a nonce r_U according to system time, and computes

$$IPB_i = H(ID_i||H(PW_i \oplus b)),$$

$$C_1 = r_U \oplus W_i \oplus IPB_i, \ C_2 = H(r_U) \oplus Z_i \oplus IPB_i.$$

Then U_i 's smart card sends { $update, ID_i, C_1, C_2$ } to S, in which update is the message type indicating this is a password change request message.

3. When S receives this message, it checks the validity of ID_i . If ID_i is invalid, it rejects the service request.

4. S computes $D_1 = C_1 \oplus H(ID_i, K_1)$ and $D_2 = H(D_1) \oplus C_2$. After that S checks whether D_2 is equal to $H_{K_2}(ID_i)$. If they are not equal, S rejects the password change request. Otherwise, S authenticates U_i successfully and accepts the password change request.

5. S generates a random nonce r_S , computes $C_3 = H(ID_i, K_1) \oplus r_S$ and $C_4 = H(r_U, r_S)$ and sends $\{C_3, C_4\}$ to U_i .

6. The smart card computes $D_3 = C_3 \oplus W_i \oplus IPB_i$. Then it compares $H(r_U, D_3)$ with C_4 . If they are equal, the smart card successfully authenticates S and prompts U_i to enter a new password.

7. U_i enters a new password PW_i^{new} . Then the smart card computes $IPB_i^{new} = H(ID_i||H(PW_i^{new} \oplus b))$. After that the smart card computes $W_i^{new} = W_i \oplus IPB_i \oplus IPB_i^{new}$, which yields $H(ID_i, K_1) \oplus IPB_i^{new}$, and computes $Z_i^{new} = Z_i \oplus IPB_i \oplus IPB_i^{new}$, which yields $H_{K_2}(ID_i) \oplus IPB_i^{new}$. The smart card also updates $T_i^{(j)_2}$ to $T_i^{(j)_2} \oplus IPB_i \oplus IPB_i^{new}$ for all remaining tickets, which yields $\alpha_i^{(j)} \oplus IPB_i^{new}$.

4 Security Analysis of the Proposed Scheme

In this section, we present security analysis of the proposed scheme. In section 4.1, we show that the server's secret key cannot be obtained by an adversary who monitors the communication. In sections 4.2-4.5, we show that the proposed scheme is resistant to lost smart card attack, masquerade attack, replay attack and valid period extending attack.

4.1 Secrecy of the server's secret key

In the proposed scheme, communications between U_i and S are through a common channel. So we assume the adversary can eavesdrop this channel and get any messages transmitted between U_i and S.

During the kth verification request and mutual authentication phase, the adversary can get messages $\{ID_i, T_i^{(k)_1}, C_1, C_2, C_3, C_4\}$, in which ID_i and $T_i^{(k)_1}$ have no relationship with K_1 or K_2 , and C_1, C_2, C_3, C_4 are as follows:

$$C_1 = r_U \oplus W_i \oplus IPB_i = r_U \oplus H(ID_i, K_1)$$

$$C_2 = H(r_U) \oplus T_i^{(k)_2} \oplus IPB_i = H(r_U) \oplus H_{K_2}(ID_i \parallel TID_i^{(k)} \parallel VP_i^{(k)})$$

$$C_3 = H(ID_i, K_1) \oplus r_S$$

$$C_4 = H(r_U, r_S)$$

As r_U and r_S are both random nonces generated by U_i and S, the adversary cannot guess their values. In addition, due to the one-wayness of hash function, the adversary cannot get r_U or r_S from $H(r_U, r_S)$. So the adversary cannot get the value of $H(ID_i, K_1)$ or $H_{K_2}(ID_i \parallel TID_i^{(k)} \parallel VP_i^{(k)})$. As a result, the adversary cannot get any information on S's secret keys by eavesdropping channels.

4.2 Resistance to Attacks Based on Lost Smart Card

A user's smart card may get lost due to an accident or the user's carelessness. In this section, we show that when an adversary gets a lost smart card, he cannot carry out attacks to the proposed scheme. We first show that the proposed scheme is resistant to offline password guessing attack. After that, we show that the lost tickets cannot be used by the adversary.

Resistance to Offline Password Guessing Attack

When an adversary gets the lost smart card of U_i , he can extract the stored data from it by monitoring the power consumption [12] or analyzing the leaked information [13]. The values stored in U_i 's smart card are

$$ID_i, t, W_i, Z_i, b, T_i^{(j)}, j = 1, 2, ..., t.$$

The adversary can pick up a password candidate PW'. Then he can computes $IPB'_i = H(ID_i||H(PW'_i \oplus b))$. From $W_i = IPB_i \oplus H(ID_i, K_1)$, $Z_i = H_{K_2}(ID_i) \oplus IPB_i$ and $\beta_i^{(j)} = \alpha_i^{(j)} \oplus IPB_i = H_{K_2}(ID_i \parallel TID_i^{(j)} \parallel VP_i^{(j)}) \oplus IPB_i$ he can further compute $H'(ID_i, K_1) = W_i \oplus IPB'_i$, $H'_{K_2}(ID_i) = Z_i \oplus IPB'_i$ and $H'_{K_2}(ID_i \parallel TID_i^{(j)} \parallel VP_i^{(j)}) = \beta_i^{(j)} \oplus IPB'_i$. However, as K_1 and K_2 are S's secret keys, the adversary cannot get them. So the adversary cannot test whether PW' is equal to PW_i from these values.

From the above analysis, we can see that the proposed scheme is resistant to offline password guessing attack.

Resistance to Attacks Based on Lost Tickets

When an adversary gets the lost smart card of U_i , he can extract the stored tickets $T_i^{(j)}$, j = 1, 2, ..., t from it. However, as the adversary cannot get PW_i from offline guessing attack (refer to 4.2), he cannot compute $IPB_i = H(ID_i||H(PW_i \oplus b))$. As a result, the adversary cannot make a valid verification request message $\{C_1, C_2\}$, because both the computations of C_1 and C_2 require the knowledge of IPB_i .

From the above analysis, we can see that the proposed scheme is secure when an adversary gets a lost ticket.

4.3 Resistance to Masquerade Attack

Suppose the adversary gets a set of history messages during the past channel eavesdropping. In order to masquerade as U_i , the adversary has to forge a valid message $\{ID_i^*, T_i^{(k)}, C_1^*, C_2^*\}$ which can pass S's authentication process.

It's easy to see that the computation of $C_1^* = r_U \oplus H(ID_i, K_1)$ requires the adversary to get knowledge of $H(ID_i, K_1)$. From 4.1 we know that the adversary cannot get $H(ID_i, K_1)$, so he cannot forge a valid verification request message either.

On the other hand, if the adversary wants to masquerade as the server, he must be able to compute a valid message $\{C_3, C_4\}$, in which $C_3 = H(ID_i, K_1) \oplus r_S$ and $C_4 = H(r_U, r_S)$. Because the computation of C_3 requires the knowledge of $H(ID_i, K_1)$, which the adversary does not have, the adversary cannot masquerade as the server.

From the above analysis, we can see that the proposed scheme is resistant to masquerade attack.

4.4 Resistance to Replay Attack

In the proposed scheme, the unique ticket ID and the nonces are used for preventing replay attack.

Assume that the adversary gets a valid message $\{ID_i, T_i^{(k)_1}, C_1, C_2\}$ by eavesdropping communications between U_i and S. Then he tries to resend the message to S after a period of time when the mutual authentication process finishes, with the expectation of obtaining data verification service. However, when S receives this message, it will find that the ticket ID $TID_i^{(k)}$ is already on the bullet in board. So S will reject the adversary's service request and terminate the process.

On the other hand, if the adversary gets a message from S to U_i : $\{C_3 = H(ID_i, K_1) \oplus r_S, C_4 = H(r_U, r_S)\}$. The adversary may try to resend it to U_i after a period of time when the mutual authentication process is finished. However, when U_i 's smart card receives $\{C_3, C_4\}$, the nonce r_U has already been deleted from its storage. So this message will be discarded immediately.

From the above analysis, we can see that the proposed scheme is resistant to replay attack.

4.5 Resistance to Valid Period Extending Attack

Whenever U_i wants to use a ticket whose valid period does not include the current date, the protocol can ensure that the ticket cannot be used even if U_i tries to modify the ticket's begin date or expiration date.

Assume U_i wants to use ticket $T_i^{(k)}$ by changing the ticket's $VP_i^{(k)}$ to $\hat{VP}_i^{(k)}$, so that the current date is included in $\hat{VP}_i^{(k)}$. Denote the modified ticket by $\hat{T}_i^{(k)}$. When U_i sends the message $\{ID_i, \hat{T}_i^{(k)_1}, C_1, C_2\}$ to S, S computes $D_1 = C_1 \oplus H(ID_i, K_1)$ and $D_2 = H(D_1) \oplus C_2$. Then S computes $H_{K_2}(ID_i \parallel TID_i^{(k)} \parallel \hat{VP}_i^{(k)})$ and compares it with D_2 . Since $D_2 = H(D_1) \oplus C_2 = H(C_1 \oplus H(ID_i, K_1)) \oplus C_2 = \alpha_i^{(k)} = H_{K_2}(ID_i \parallel TID_i^{(k)} \parallel VP_i^{(k)}) \neq H_{K_2}(ID_i \parallel TID_i^{(k)} \parallel \hat{VP}_i^{(k)})$, S will reject U_i 's verification request and terminate the process immediately.

From the above analysis, we can see that the server can limit the data verification frequency by specifying a valid period for each ticket. For example, the typical coverage of the valid period can be one week, one month, one year, etc. The proposed scheme can prevent the client from using a ticket whose valid period does not include the current date.

5 Performance Analysis

In this section we present performance analysis of the proposed scheme. The main operations include the hash computation and the exclusive-or operations, which are summarized in Table 1. We use the hash-based message authentication code (HMAC) [14] as the keyed hash, which incurs 2 exclusive-or and 2 common hash operations.

	operation	verification request phase	mutual authentication phase
client	xor	5	1
	hash	3	2
	keyed hash	0	0
	xor	0	3
server	hash	0	4
	keyed hash	0	1

Table 1: Performance Analysis of the Proposed Scheme

Compared with [11], which uses expensive operations like modular exponentiations, our scheme is much more efficient. According to the Crypto++ benchmarks¹ [15], the SHA-256 hash algorithm [16] can achieve a throughput of 111MiB/Sec under Intel Core2 1.83GHz processor. In the proposed scheme, the length of a message to be hashed does not exceed 1KB, so the average time for one hash computation is about 0.009ms. If the Lin-Chang scheme uses

 $^{^1\}mathrm{The}$ Crypto++ 5.6.0 benchmark is evaluated in Intel Core
2 1.83 GHz processor under Windows Vista in 32-bit mode

the primes of length 1024bits, then one time of modular exponentiation will cost about 1.46ms under the same platform [15]. Our scheme is at least 40 times more efficient than the Lin-Chang scheme [11] in case of the server's computation time. This is a tremendous performance improvement to the cloud servers.

Performance Benefit by Limiting Verification Frequency The cloud server can also benefit its performance by limiting the client's data verification frequency. This is achieved by controlling the valid period of the digital ticket. From the performance evaluation of [6] we know that for a 30MB file, the computation time of the data integrity verification is approximately 1 second^2 at the server side. Assume the cloud storage system has 100,000 clients, and each client performs data integrity verification once a week. Then the average computation time of the cloud server during one day is approximately:

 $\frac{100,000}{7} \cdot 1s \approx 3.97 \ hours.$

By using the proposed scheme, the server can limit the data verification frequency to once per month, or even once per quarter, so that the average computation time of the cloud server is reduced to 56 minutes per day and 18.5 minutes per day respectively.

6 Conclusions

In this paper, we propose a time-bound ticket-based mutual authentication scheme. In the proposed scheme, the digital tickets are associated with the client's smart card, which effectively prevents the ticket from being used by other clients. By designing a mutual authentication based on the client's smart card, both the server and the client are assured of each other's identity. The proposed authentication scheme can efficiently decrease the server's processing overhead by limiting the data verification frequency. By security analysis and performance analysis, the proposed scheme is shown to be both secure and efficient. It is very suitable for providing mutual authentication in cloud computing.

Acknowledgment

This work was supported by NSF CNS-0845149, NSF CCF-0915374 and Knowledge Innovation Program of Chinese Academy of Sciences (No. YYYJ-1013).

Bibliography

- [1] B. Hayes, "Cloud computing," Commun. ACM, vol. 51, no. 7, pp. 9–11, 2008.
- [2] C. Cachin, I. Keidar, and A. Shraer, "Trusting the cloud," SIGACT News, vol. 40, no. 2, pp. 81–86, 2009.
- [3] Amazon.com, "Amazon Web Services (AWS)," http://aws.amazon.com/s3/, 2009.
- [4] Microsoft.com, "Microsoft Windows SkyDrive," http://windowslive.com/online/skydrive, 2009.

 $^{^2{\}rm The}$ experiment environment in [6] is Intel 2.8 GHz Pentium IV system with a 512 KB cache, an 800 MHz EPCI bus, and 1024 MB of RAM.

- [5] Nirvanix.com, "Nirvanix cloudNAS," http://www.nirvanix.com/products-services/, 2009.
- [6] G. Ateniese, R. Burns, R. Curtmola, J. Herring, L. Kissner, Z. Peterson, and D. Song, "Provable data possession at untrusted stores," in CCS '07: Proceedings of the 14th ACM conference on Computer and communications security, (New York, NY, USA), pp. 598–609, ACM, 2007.
- [7] A. Juels and B. S. Kaliski, Jr., "Pors: proofs of retrievability for large files," in CCS '07: Proceedings of the 14th ACM conference on Computer and communications security, (New York, NY, USA), pp. 584–597, ACM, 2007.
- [8] E.-C. Chang and J. Xu, "Remote integrity check with dishonest storage server," in 13th ESORICS, pp. 223–237, Springer Berlin / Heidelberg, 2008.
- [9] A. Heitzmann, B. Palazzi, C. Papamanthou, and R. Tamassia, "Efficient integrity checking of untrusted network storage," in *StorageSS '08*, pp. 43–54, ACM, 2008.
- [10] K. D. Bowers, A. Juels, and A. Oprea, "HAIL: a high-availability and integrity layer for cloud storage," in CCS '09, (New York, NY, USA), pp. 187–198, ACM, 2009.
- [11] I.-C. Lin and C.-C. Chang, "A countable and time-bound password-based user authentication scheme for the applications of electronic commerce," *Information Sciences*, vol. 179, no. 9, pp. 1269 – 1277, 2009.
- [12] P. C. Kocher, J. Jaffe, and B. Jun, "Differential power analysis," in CRYPTO '99: Proceedings of the 19th Annual International Cryptology Conference on Advances in Cryptology, (London, UK), pp. 388–397, Springer-Verlag, 1999.
- [13] T. Messerges, E. Dabbish, and R. Sloan, "Examining smart-card security under the threat of power analysis attacks," *IEEE Transactions on Computers*, vol. 51, no. 5, pp. 541–552, 2002.
- [14] H. Krawczyk, M. Bellare, and R. Canetti, "HMAC: Keyed-Hashing for Message Authentication," *RFC2104*, February 1997.
- [15] "Crypto++ 5.6.0 benchmarks," http://www.cryptopp.com/benchmarks.html.
- [16] "Secure hash standard," Federal Information Processing Standards Publication 180-2, August 2002.

A Network Coding based DTN Convergence Layer Reliable Transport Mechanism over InterPlaNetary Networks

S. Haoliang, L. Lixiang, H. Xiaohui

Sun Haoliang

 Graduate University of the Chinese Academy of Sciences
 Yuquanlu, Beijing, 100049, P. R. China
 National Key Laboratory of Integrated Information System Technology, Institute of Software, Chinese Academy of Sciences
 4# South Fourth Street, Zhong Guan Cun, Beijing 100190 P.R. China E-mail: haoliang08@iscas.ac.cn

Liu Lixiang, Hu Xiaohui

Institute of Software, Chinese Academy of Sciences 4# South Fourth Street, Zhong Guan Cun Street, Beijing, China, 100190

Abstract: The realization of deep space scientific missions are enabled by the developments in the space technologies. TCP can not provide effective communication service in deep space links because of the long propagation delay and high BERs characteristics. Scientific community propose Delay Tolerant Network (DTN) for resolving the communication problem between earth and other planet. DTN introduces a new layer in the protocol stack called Bundle Layer that is placed between application layer and transport layer.DTN calls the transport protocols that it uses to move data across different networks convergence layers. This study present a novel Network Coding based convergence layer Reliable TransPort mechanism(NC-RTP) in order to provide effective communication service between DTN peers. This mechanism transmit a coded bundle every M original bundles (M is related to the packet error rate of the communication link). The coded bundle is a random linear combination of previous M original bundles. Using coded bundle and (M-1)original bundles, the receiver could decode and generate a single lost bundle of previous M original bundles. In this way, NC-RTP could compensate any single lost bundle in (M+1) transmitted bundles (including M original bundles and 1 coded bundle). Our theoretical and simulation performance evaluation results reveal that NC-RTP can enhance transmit reliability and make file transfered faster.

Keywords: DTN, Convergence Layer, Reliable Transport, Network Coding.

1 Introduction

The realization of deep space scientific missions are enabled by the developments in the space technologies. The future space exploration includes missions to deep space that require communication among planets, moons, satellites, asteroids, robotic spacecrafts, and crewed vehicles. Significant amount of scientific data to be delivered to the Earth are produced by these missions. The next step in the design and development of deep space networks is expected to be the Internet of the deep space planetary networks ,called InterPlaNetary (IPN) Internet [1].IPN is an architecture envisioned for interconnecting earth with other planets. The purpose of IPN is to build a communication infrastructure between planets and satellites .

InterPlaNetary (IPN) Internet, in contrast to conventional Internet links, are characterized by:

1. High Link Bit Error Rates. Deep space links have extremely high bit error rates, which may be up to 10 [2], [3].

2. Very High Propagation Delays. The propagation delay increases with the distance. The propagation delay between Earth and Mars varies between 8 and 40 minutes depending on the orbital location of the planets [2] [3].and the propagation delays to outer space planets become even higher.

3. Bandwidth Asymmetry. Asymmetry between the forward and the return path bandwidth may be up to 1000:1 [3] [4].

4. Intermittent Connectivity. Planetary bodies, asteroids or spacecraft may periodically interrupt the communication link between the path endpoints [3] [5].

In this network, classical transport layer mechanisms are not suitable. Scientific community propose Delay Tolerant Network (DTN) [6] for resolving the communication problem between earth and other planet such as Mars. DTN is an architecture particularly useful in scenarios with very long transmission delay or intermittent connectivity, like IPN. DTN introduces a new layer in the protocol stack called Bundle Layer that is placed between application layer and transport layer [7].

The bundle layer resolves the high error rates, long delay, asymmetric data rates and intermittent connectivity by using a store and forward mechanism. It sends a bundle of message fragments to the next-hop node with per-hop error control, which increases the probability of data transmission.

DTN calls the transport protocols that it uses to move data across different networks convergence layers. A convergence layer should make best use of the intermittent and temporary links in DTN. It transfers as much data as possible while the link between peering DTN nodes is up and available. It is widely accepted that regular TCP cannot operate efficiently as a convergence layer transport protocol for InterPlaNetary networks based on DTN architecture [8]. TCP needs 240 minutes to reach Slow Start Threshold equal to 20 packets [5] over a 40 minute Round Trip Time (RTT) path, and longer RTT paths degrade TCP's performance further. Since TCP was designed to operate over wired links, where the link error rate is insignificant, the protocol cannot cope with high link bit error rates [9]. Finally, TCP's transmission rate depends largely on the receiver's feedback.TCP sends one acknowledgment (ACK) for each successfully received data packet. On the presence of bandwidth asymmetry, the large number of ACKs will cause congestion on the reverse path, reducing TCP's transmission rate. As a result, in InterPlaNetary networks based on DTN architecture, a convergence layer transport protocol other than TCP is needed.

The remainder of the paper is organized as follows. Section 2 introduces current convergence layer protocols in DTN architecture and discusses the drawbacks of them. Section 3 proposes our new Network coding based convergence layer reliable transport mechanism(NC-RTP). The theoretical evaluation of NC-RTP is presented in Section 4. And the experimental evaluation is in Section 5. We conclude the paper in Section 6.

2 Related Works

There exist already a number of convergence layer transport protocols for InterPlaNetary networks based on DTN architecture . In this section, we briefly review these proposals.

Saratoga [10] is a reliable rate-based UDP/IP file transfer protocol. It is capable of transferring efficiently both small and very large files. It has been developed by Surrey Satellite Technology Ltd (SSTL) used for mission imaging data. Saratoga was designed for dedicated point-to-point links between DTN Peers, which focuses on transferring data efficiently to the next hop when link connectivity is available. Saratoga achieves efficient transmission by sending out data packets

at the line rate. It also uses a negative acknowledgment strategy in order to deal with channel bandwidth asymmetries. Saratoga is used as a convergence layer to exchange Delay Tolerant Networking bundles [6], [7] between peer nodes [10].

Similarly to Saratoga, the Licklider Transmission Protocol (LTP) [11] is a point-to-point protocol applied as a DTN convergence layer. LTP transfer unnamed blocks of data and introduces the concept of partial reliability by dividing each block of data into two parts: the reliable red part and the unreliable green part. Moreover, LTP send laconic acknowledgments only upon encountering explicit solicitations for reception reports (checkpoints) in the sequence of incoming data segments of the red part of the block. When the communication link is not available, Deferred Transmission is possible as well.

Deep Space Transport Protocol (DS-TP) [12]also can be adapted to serve as an efficient convergencelayer for DTN, by transferring DTN bundles as well as files. It implements a retransmission technique, called Double Automatic Retransmission (DAR), which allows for fast and efficient holefilling at the receiver's buffer. DAR sends each bundle twice, importing some delay between the original transmission and the retransmission. More precisely, one redundant bundle is transmitted every 1/e-1 original bundles(e presents the packet error rate PER). Therefore, in the presence of link errors, corrupted bundles will eventually be replaced by the same correct bundles that arrive later.

Saratoga and LTP achieve efficient transmission by sending out data bundles at the line rate. This ensures that as much data as possible is transferred to the peering node. DS-TP implements Double Automatic Retransmission(DAR) to fast and efficiently compensate lost bundles.

But ,according to the algorithm described in [12],DAR only retransmit finite bundles ((N/(1/e-1) bundles ,N presents the number of bundles of the file). DAR can not provide effective retransmission for all bundles, a more effective retransmission technique is needed to enhance transmit reliability.

3 Network coding based convergence layer reliable transport mechanism

Network coding is an important approach to enhance transmit reliability in lossy link. The central idea of network coding is to transmit a combination of multi-packet to replace a single packet. The receiver could decode and generate needed packets using received packets. It has a good ability to deal with transmit error [13]. The simplest coding approach is random linear coding. We treat packets as vectors over a finite field. To get a coded packet, random linear coding generates a group of random coefficients vector and perform matrix multiplication with the packets vector. The transmission of coded packets are independent of each other. Generally, to simplify computer programming and memory accessing, the size of finite field is set to 8, which is the size of a byte. For example, we have two 16 bits data packets to encode, 3A47H and 5F33H(hex), each packet is considered to be a vector with size 2. Choose two decimal coefficients randomly,4, 30(the coefficients range from 0 to 255). The coded packet is 43A47H+305F33H=E91CH+27FAH=1116H.

Using network coding in DTN convergence layer can enhance transmit reliability and transmit as many bundles as possible before get any ACK information from the receiver. In our study, we transmit a coded bundle every 1/e - 1 original bundles(e presents the packet error rate PER). The coded bundle is a random linear combination of previous 1/e - 1 original bundles. Considering original bundles are liner independent with each other, the coefficients can be set to 1 to all original bundles in the coded bundle. Theoretically, every 1/e bundles will lost one bundle under PER e. Using 1/e - 1 original bundles and the coded bundle of 1/e - 1 original bundles, we can decode

and generate the lost original bundle with Gaussian elimination. Thus we can say ,in this way ,the receiver could get all 1/e-1 original bundles using received 1/e-1 bundles(including original and coded bundles). This approach is called Network Coding based convergence layer Reliable TransPort mechanism(NC-RTP). This mechanism could compensate lost bundle forwardly.

We use an example to illustrate NC-RTP. Let e = 0.2, so 1/e - 1 = 4. A coded bundle is generated and send every 4 bundles.Let b_n presents the n_{th} bundle.The transmission sequence is $b_1, b_2, b_3, b_4, b_{1,2,3,4}(b_{1,2,3,4}\text{presents}$ the coded bundle of b_1, b_2, b_3, b_4 and $b_{1,2,3,4} = b_1 + b_2 + b_3 + b_4$. Theoretically, 1 bundle will be lost during the transmission. Suppose b_3 is lost, and the reciver get b_1, b_2, b_4 , and $b_{1,2,3,4}$. Since $b_{1,2,3,4} = b_1 + b_2 + b_3 + b_4$ and we have $b_1, b_2, b_4, b_3 = b_{1,2,3,4} - b_1 - b_2 - b_4$, in this way we can decode and generate b3.

Compare to DS-TP, NC-RTP could compensate any single lost bundle in 1/e - 1 original bundles, and transmit coded bundles (N presents the number of bundles of the file) could make all bundles in the file have the chance to be compensated. In this way, NC-RTP can reduce the number of lost bundles during the transport procedure and enhance the transport reliability. Moreover, we generate coded bundle using bundles, the coding and decoding delay is neglectable comparing to the propagation delay of interplanetary links. Because of the long propagation delay of interplanetary links, using NC-RTP can enhance the transmit reliability and complete file transfers faster.

In the next part of this section, we describe in detail the operation process of Network Coding based convergence layer Reliable TransPort mechanism (NC-RTP),followed by pseudo-code of the algorithm.

The initial value of PER e can be set by previous transmit experience. In the transmit process, the sender can compute the value of e depending on the acknowledgement information from the receiver and change it. In this study ,we consider the PER is fixed during the transmit process to simplify the analysis.

Sender procedure: The sender transmit bundles in order and transmit a redundant coded bundle every 1/e - 1 original bundles. The coded bundle is a liner combination of previous 1/e - 1 original bundles ,all coefficients are set to 1. When the sender get SNACK, retransmit lost bundle according to the acknowledgement information. In order to compensate lost bundle, use NC-RT in the retransmit operation, send a redundant coded bundle every 1/e - 1 original bundles.

Receiver procedure: When the receiver get bundles , it put original and coded bundles into its receiving queue. When it get accumulative ACK_DELAY (ACK_DELAY is set to the ratio of forward and reverse bind with of the link .When the file to be transmitted has few bundles, we can reduce ACK_DELAY to provoke send SNACK procedure) bundles ,send SNACK procedure is provoked.

Send SNACK procedure: When send SNACK procedure is provoked ,the sender check original and coded bundles in its receiving queue, decode and generate lost bundle using original and coded bundles, put generated bundle into receiving queue. Check the receiving queue again, if there is any lost bundle, write the lost bundle information into the SNACK bundle and send the SNACK bundle.

The sender side algorithm has to respond to two types of events: the arrival of a bundle from the upper layer, and get SNACK bundle from the receiver. The receiver side algorithm has to respond to the arrival of a bundle.

Source side:

1)Set SND COUNT to 0.

2)Wait state: If any following events occurs, respond as follows; else, wait.

3)Bundle arrives from upper layer:

a)Send the bundle to the receiver. SND COUNT++. b) if (SND COUNT = 1/e - 1)Generate a coded bundle with previous 1/e - 1 bundles. Send the coded bundle. $SND_COUNT=0.$ 4)Get a SNACK from the receiver: a)Read the SNACK information, find out which bundle is lost. b)retransmit the lost bundle, go to 3). Receiver side: 1)Set RCV COUNT to 0. 2)Set ACK DELAY properly (according to the ratio of the forward and reverse bandwith and the file size). 3) Wait state: If get a bundle from lower layer, respond; else, wait. 4)Bundle arrives from lower layer: a) $RCV \quad COUT + +;$ Add the bundle into the receiving queue, according to the bundle's ID. b) If $(RCV_COUNT = ACK_DELAY)$ Using the original bundles and the coded bundle to encode the lost bundles. Enqueue the generated bundle in to the receiving queue. After the encode operation, if there is any lost bundle, write the lost bundle information into the SNACK bundle and send the SNACK bundle.

4 Theoretical Evaluation

We attempt to evaluate, theoretically, the performance of NC-RTP.In order to derive some initial evaluation results regarding the performance of the NC-RTP, we assume that the link error rate remains constant during the file transfer.

In particular, we consider the Fixed-Rate Transport Protocol (FRTP), whose main functionality is summarized as follows: the FRTP sender sends data on a fixed rate according to the pre-scheduled line rate. The FRTP receiver, responds with SNACKs in order to signal for lost bundle in the receiving queue.

To simplify the analysis, we consider that SNACKs are sent back to the sender without lost, all acknowledgement information are sent back to the receiver side successfully. We evaluate the performance of the aforementioned protocols over a simple one-hop topology. Such topology represents a point-to-point links between DTN Peers in interplanetary networks. Obviously, the primary metric of interest is the time required for the whole file to be delivered at the receiver side.

We define a Round to be the end-to-end transmission of a specific amount of data. A file transfer consists of several Rounds, during the first Round the original file is transmitted, while during the rest of the Rounds, the sender retransmits bundles lost in previous Rounds.

Assume the file consists of N bundles. The link PER is e. The FRTP sender will begin the transmission of the file at the channel rate. After completion of the first round the sender will have transmitted N bundles. During the first round, $N \cdot e$ bundles are lost and will need to be retransmitted during the second round. Similarly, $N \cdot e^2$ bundles are lost during the second round and need to be retransmitted during the third round. During the n_{th} round, the FR-TP sender will need to retransmit $N \cdot e^n$ bundles. We assume that when $N \cdot e^n < 1$ the file transfer is complete.

Therefore, FRTP needs n_{FRTP} rounds in order to complete the file transfer:

$$n_{FRTP} = \log_e^{e^n} = \log_e^{\frac{1}{N}} = \frac{\ln \frac{1}{N}}{\ln e} \tag{1}$$

Besides FRTP,we also consider DSTP to evaluate NC-RTP's performance. According to its operational properties, DSTP will transmit both original and redundant bundles at the line rate. More precisely,one redundant bundle is transmitted every 1/e - 1 original bundles. Let r_1 be the number of transmitted redundant bundles during first round, $r_1 = N/(1/e - 1)$. It is to say, there are r_1 bundles that were sent twice, and $N - r_1$ bundles that were send only once. We assume that the number of bundles lost during the first round (and need to be retransmitted during the second round) equals a_1 where:

$$a_1 = (N - r_1) \cdot e + r_1 \cdot e^2 \tag{2}$$

Substituting r_1 into Equation 2, we get that:

$$a_1 = N \cdot e \cdot (1 - e) \tag{3}$$

Assume that the number of bundles lost during the second round equals a_2 where:

$$a_2 = N \cdot e^2 \cdot (1 - e)^2 \tag{4}$$

Generalizing Equations 3 and 4, we assume that during the n_{th} round, the DS-TP sender will need to retransmit a_n bundles, where

$$a_n = N \cdot e^n \cdot (1 - e)^n \tag{5}$$

The file transfer is complete, once the following equation holds:

$$N \cdot e^n \cdot (1-e)^n < 1 \tag{6}$$

$$n_{DSTP} = \log_{e \cdot (1-e)}^{(e \cdot (1-e))^n} = \log_{e \cdot (1-e)}^{\frac{1}{N}} = \frac{\ln \frac{1}{N}}{\ln (e \cdot (1-e))}$$
(7)

We add NC-RTP on FRTP, called FRTP-NC. FRTP-NC transmits redundant coded bundles every 1/e - 1 original bundles. The coded bundle is a liner combination of previous 1/e - 1 original bundles. According to its operational properties, FRTP-NC will transmit both original and coded bundles at the line rate. In ideal condition, every 1/e – bundles (including 1/e - 1 original bundles and one coded bundles) will lost only one bundle. In this way, the redundant coded bundle will compensate the lost bundle, the file transfer will complete in just one round. However, when more than one bundles are lost in every 1/e bundles, the redundant coded bundle can not compensate all the lost bundles, the sender needs to retransmit lost bundles. In this case, the coded bundle can replace one lost bundle, and the sender needs to retransmit less bundles.

Let a = 1/e, we divide the whole file into N/(a - 1) parts. All parts are independent of each other. In a single part, if only one bundle is lost(including original and coded bundles), no retransmission is needed. If $k(k \leq a)$ bundles are lost, only k - 1 bundles are needed to retransmit.

Let p_k presents the probability of k bundles retransmission in a single part(lost k+1 bundles).

$$p_k = C_a^{k+1} \cdot e^{k+1} \cdot (1-e)^{a-k-1} \tag{8}$$

The expectation of k presents the number of bundles needs to retransmit in a single part.

$$E(k) = \sum_{k=1}^{a-1} k \cdot p_k = \sum_{k=1}^{a-1} k \cdot C_a^{k+1} \cdot e^{k+1} \cdot (1-e)^{a-k-1}$$
(9)

Let k+1 = j, we get

$$E(k) = \sum_{j=2}^{a} (j-1) \cdot C_a^j \cdot e^j \cdot (1-e)^{a-j} = \sum_{j=2}^{a} j \cdot C_a^j \cdot e^j \cdot (1-e)^{a-j} - \sum_{j=2}^{a} C_a^j \cdot e^j \cdot (1-e)^{a-j}$$
(10)

Because $\sum_{j=0}^{a} j \cdot C_a^j \cdot e^j \cdot (1-e)^{a-j} = a \cdot e = 1$, $\sum_{j=0}^{a} C_a^j \cdot e^j \cdot (1-e)^{a-j} = (e+1-e)^a = 1$, and a = 1/e,

$$E(k) = 1 - (1 - e)^{1/e - 1} - [1 - (1 - e)^{1/e} - (1 - e)^{1/e - 1}] = (1 - e)^{1/e}$$
(11)

Consider the whole file consists of N/(a-1) parts, in the first transmission, $E(k) \cdot N/(a-1)$ bundles need to retransmit. Let y_1 be the number of transmitted redundant bundles during first round.

$$y_1 = E(k) \cdot N/(a-1) = (1-e)^{1/e} \cdot N/(1/e-1) = N \cdot e \cdot (1-e)^{1/e-1}$$
(12)

Generalizing Equation 12, we assume that during the n_{th} round, the FRTP-NC sender will need to retransmit y_n bundles, where

$$y_n = N \cdot e^n \cdot (1 - e)^{n \cdot (1/e - 1)} \tag{13}$$

The file transfer is complete, once the following equation holds:

$$y_n = N \cdot e^n \cdot (1 - e)^{n \cdot (1/e - 1)} < 1 \tag{14}$$

$$n_{FRTP-NC} = \log_{e \cdot (1-e)^{(1/e-1)}}^{e^n \cdot (1-e)^{n \cdot (1/e-1)}} = \frac{\ln \frac{1}{N}}{\ln e \cdot (1-e)^{(1/e-1)}}$$
(15)

We use two figures to evaluate the performance of FRTP, DSTP and FRTP-NC. Figure 1 presents the numbers of rounds of FRTP, DSTP and FRTP-NC in different PERs(N = 100000). Figure 2 presents the numbers of rounds of FRTP, DSTP and FRTP-NC in different file sizes(PER = 0.25).

From Figure 1 we observe that , when file size is fixed, for small error rates , the three protocols perform similarly. As the link error rate increases, the rounds that FRTP needs grows sharply, and the rounds that FRTP-NC and DSTP need grow slowly, and FRTP-NC always needs less rounds than DSTP. Only when PER = 0.5, DSTP and FRTP-NC have the same performance. At this time ,DSTP and FRTP-NC operate in the same way(transmit each bundle twice).

From Figure 2 we observe that ,when the link error rate is fixed, the growth of the file size results in the increment of the rounds that three protocols need, but FRTP-NC is the lowest.

Our theoretical evaluation revealed that, FRTP-NC that using NC-RTP needs less time to complete file transmission and it have a better ability to handle high PER and large files in interplanetary networks. NC-RTP can be used in convergence layer to achieve efficient transmission between peer nodes .

5 Experimental Evaluation

In this section, we evaluate the performance of NC-RTP experimentally. We use the network simulator "ns-2" [14] for all experiments and simulate. We use a simple one-hop topology, whose propagation delay is 1200s. Such topology represents a point-to-point link between DTN Peers in interplanetary networks. Since Congestion control can not respond to the link status in time because of the long propagation delay of deep-space links, the transport protocols we use in the experiment are without congestion control, using fixed rate data transmission. In our experiment

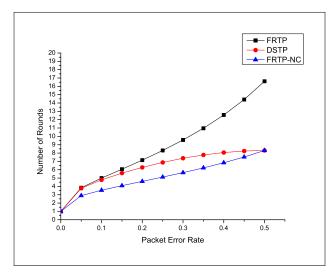


Figure 1: RTT=2400s, File size=100000 bundles, Rounds that three protocols need in different PERs

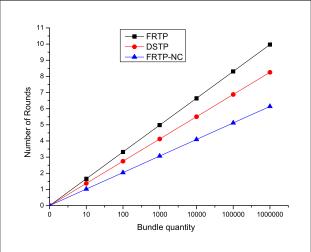


Figure 2: RTT=2400s, File size=100000 bundles, Rounds that three protocols need in different file size

, we use FRTP, DSTP in[12], and FRTP-NC that using NC-RTP to evaluate the performance of NC-RTP. The transmit rate is 1000 bundles per second.

In Scenario 1, we evaluate the performance of NC-RTP in different PERs. The propagation delay is set to 1200s, file size is 100000 bundles and PER changes from 0 to 0.5. In scenario 2, we evaluate the performance of NC-RTP in different file sizes. The propagation delay is still 1200s, PER is 0.25 and file size changes from 0 to 1000000 bundles.

From Figure 4 and Figure 5 we can get, the experimental result generally tally with the theoretical evaluation. FRTP-NC that using NC-RTP needs less time to complete file transmission than FRTP and DSTP. And with the increment of file size , the file transfer complete time of FRTP and DSTP grow sharply, FRTP-NC's transfer complete time grows slowly. So we can get a conclusion that, NC-RTP can cope with high PER in interplanetary networks and more suitable for large files, using NC-RTP can enhance transmit reliability and complete file transfers faster.

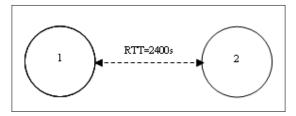


Figure 3: Topology of a point-to-point links between DTN Peers

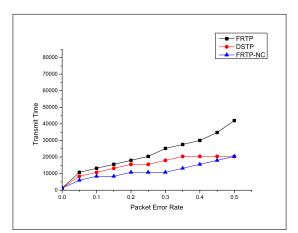


Figure 4: RTT=2400s, File size=100000 bundles, transmit complete time change with PERs

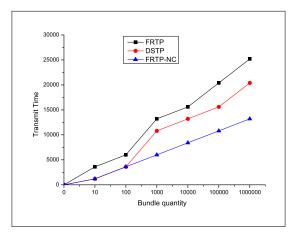


Figure 5: RTT=2400s,PER=0.25 transmit complete time change with file size

6 Conclusions and Future Works

In this study, we present a Network Coding based convergence layer Reliable TransPort mechanism (NC-RTP), whose main advantage is to compensate lost bundle forwardly .NC-RTP uses fixed rate transmission and generates a redundant coded bundle periodically. The coded bundle is a random linear combination of previous bundles. Using original bundles and the coded bundle, we can decode and generate the lost original bundle with Gaussian elimination. Because of the long propagation delay of interplanetary links, using NC-RTP can enhance the transmit reliability and complete file transfers faster.

Our theoretical and simulation performance evaluation results reveal that NC-RTP presents high potential for deployability. Protocols that use NC-RTP complete file transfers faster than without NC-RT. we can get a conclusion that, NC-RTP can cope with high PER in interplanetary networks and more suitable for large files, using NC-RTP can enhance transmit reliability and complete file transfers faster.

Future works includes, specifically implementation on how to react to change PERs and further investigation and evaluation of the performance of NC-RTP when intermittent connectivity events happen on the transmission link.

Bibliography

- I.F., Akyildiz, et al., The state of the art in interplanetary Internet Communications Magazine, IEEE Volume 42, Issue 7, July 2004 Page(s): 108 118.
- [2] Ian F. Akyildiz, Ozgur B. Akan, Chao Chen, Jian Fang, and Weilian Su, InterPlaNetary Internet: State-of-the-Art and Research Challenges, *Computer Networks Journal (Elsevier)*, Vol. 43, Issue 2, pp. 75-113, October 2003.
- [3] R.C. Durst, P.D. Feighery, K.L. Scott, Why not use the Standard Internet Suite for the InterPlaNetary Internet? Available from http://www.ipnsig.org/techinfo.htm.
- [4] R.C. Durst, G.J. Miller, E.J. Travis, TCP extensions for space communications, Wireless Networks, Vol. 3, Issue 5, pp. 389-403, October 1997.
- [5] Akan, J. Fang and I.F. Akyildiz, Performance of TCP protocols in deep space communication networks, *IEEE Communications Letters*, Vol. 6, Issue 11, pp.478-480, November 2002.
- [6] V. Cerf, S. Burleigh et al., elay Tolerant Network Architecture, IETF RFC 4838, April 2007.
- [7] K. Scott and S. Bureeigh, undle Protocol Specification IRTF DTNRG Internet Draft revision 10, July 2007
- [8] O.B. Akan, J. Fang and I.F. Akyildiz, Performance of TCP protocols in deep space communication networks, *IEEE Communications Letters*, Vol. 6, Issue 11, pp.478-480, November 2002.
- [9] V. Tsaoussidis and I. Matta, pen issues on TCP for Mobile Computing, The Journal of Wireless Communications and Mobile Computing, WCMC John Wiley and Sons, Vol. 2, Issue 1, pp. 3-20, February 2002.
- [10] L. Wood, J. McKim, W. Eddy, W. Ivancic and C. Jackson, aratoga: A Convergence Layer for Delay Tolerant Networking, *IETF July 2007 meeting*.
- [11] M. Ramadas et al., "Licklider Transmission Protocol " Specification, IETF internet draft, April 2007
- [12] Ioannis Psaras, Giorgos Papastergiou, Vassilis Tsaoussidis, and Nestor Pec-cia. DS-TP:Deep-Space Transport Protocol. *IEEE Aerospace Conference*, 2008, Big Sky, Montana, USA.
- [13] T. Ho, "Networking from a network coding perspective," it PhD Thesis, Massachusetts Institute of Technology, Dept. of EECS, May 2004.
- [14] "ns-2 Network Simulator," http://www.isi.edu/nsnam/ns/

Multi-Objective Optimization for the m-PDPTW: Aggregation Method With Use of Genetic Algorithm and Lower Bounds

I. Harbaoui Dridi, R. Kammarti, M. Ksouri, P. Borne

Imen Harbaoui Dridi

LAGIS : Ecole Centrale de Lille, Villeneuve d'Ascq, France LACS : Ecole Nationale des Ingénieurs de Tunis, Tunis - Belvédère. Tunisie E-mail: imenharbaoui@gmail.com

Ryan Kammarti, Mekki Ksouri

LACS : Ecole Nationale des Ingénieurs de Tunis, Tunis - Belvédère. Tunisie E-mail: kammarti.ryan@planet.tn, Mekki.Ksouri@insat.rnu.tn

Pierre Borne

LAGIS : Ecole Centrale de Lille, Villeneuve d'Ascq, France E-mail: p.borne@ec-lille.fr

> **Abstract:** The PDPTW is an optimization vehicles routing problem which must meet requests for transport between suppliers and customers in purpose to satisfy precedence, capacity and time constraints. We present, in this paper, a genetic algorithm for multi-objective optimization of a multi pickup and delivery problem with time windows (m-PDPTW), based on aggregation method and lower bounds. We propose in this sense a brief literature review of the PDPTW, present our approach to give a satisfying solution to the m-PDPTW minimizing the compromise between total travel cost and total tardiness time. **Keywords:** PDPTW, multi-objective, aggregation method, lower bounds.

1 Introduction

The vehicle routing planning and traffic management are considered a major logistical challenge in terms of supply, inter-plant transport or distribution transport. [1] Many studies have been directed mainly towards solving the vehicle routing problem (VRP). It's an optimization vehicle routing problem to meet travel demands. Other researchers became interested on an important variant of VRP which is the PDPTW (Pickup and Delivery Problem with Time Windows) with capacity constraints on vehicles.

The PDPTW is divided into two categories: 1-PDPTW (single-vehicle) and m-PDPTW multi-vehicle).

In the m-PDPTW problem which we are interested in, we consider a vehicles fleet V_k of capacity Q_k and a set of goods to transport providers to different destinations. The goal is to provide a set of customers under certain constraints concerning vehicles and their capacity, precedence between nodes, and this by minimizing the compromise between the total travel cost and total tardiness time. In this paper we present a literature review of the PDPTW followed by the proposed approach for the optimization of pick-up and delivery problem with time window, using the genetic algorithms, aggregation method and lower bounds.

2 Litterature Review

2.1 Vehicle routing problem

The Vehicle Routing Problem (VRP) represents a multi-goal combinatorial optimization problem which has been the subject of many works and variations in the literature [2] [3]. The theory of the VRP is formulated as follows: given a depot D and a set of customers orders C =(c1, ..., Cn), to build a package routing, for a finite number of vehicles, beginning and ending at the depot. In this routing, a customer must be served only once by a single vehicle and vehicle capacity transport for a routing should not be exceeded [4] [5]. The Meta heuristics were also applied to solve the vehicle routing problem. Among these methods, we can include ant colony algorithms, which were used by Montamenni, R et al for the resolution of DVRP [6], and by Sorin C. Negulescu et al to solve the Vehicle Route Allocation Problem (VRAP) [7]. Kaoru Hirota et al have presented a computational intelligence approach to VRSDP (Vehicle Routing, Scheduling, and Dispatching Problems). The objective of the VRSDP is to produce a delivery schedule for a group of vehicles, with respect to multiple users, so that while satisfying constraints delivery cost corresponding to users' order is minimized [8]. Savelsbergh et al have shown that the VRP is a NP-hard problem [9].

2.2 The PDPTW: Pickup and Delivery Problem with Time Windows

The PDPTW is a variant of VRPTW where in addition to the existence of time constraints, this problem implies a set of customers and a set of suppliers geographically located. Every routing must also satisfy the precedence constraints to ensure that a customer should not be visited before his supplier. [10] A dynamic approach for resolving the 1-PDP without and with time windows was developed by Psaraftis, H.N considering objective function as a minimization weighting of the total travel time and the non-customer satisfaction [11]. Jih, W et al have developed an approach based on the hybrid genetic algorithms to solve the 1-PDPTW, aiming to minimize combination of the total cost and total waiting time. [12] Another genetic algorithm was developed by Velasco, N et al to solve the 1-PDP bi-objective in which the total travel time must be minimized while satisfying in priority the most urgent requests. In this literature, the method proposed to resolve this problem is based on a No dominated Sorting Algorithm (NSGA-II). [13] Kammarti, R et al treat the 1-PDPTW, minimizing the compromise between the total travel distance, total waiting time and total tardiness time, using an evolutionary algorithm with special genetic operators, tabu search to provide a set of viable solutions. [14] [15]. This work has been extended, by proposing a new approach based on the use of lower bounds and Pareto dominance method, to minimize the compromise between the total travel distance, total waiting time and total tardiness time. [16] About the m-PDPTW, Sol, M et al have proposed a branch and price algorithm to solve the m-PDPTW, minimizing the vehicles number required to satisfy all travel demands and the total travel distance. [17] Quan, L et al have presented a construction heuristic based on the integration principle with the objective function, minimizing the total cost, including the vehicles fixed costs and travel expenses that are proportional to the travel distance. [18] A new metaheuristic based on a tabu algorithm, was developed by Li, H et al to solve the m-PDPTW. [19] Li, H et al have developed a "Squeaky wheel" method to solve the m-PDPTW with a local search. [20] A genetic algorithm was developed by Harbaoui Dridi I et al treating the m-PDPTW to minimize the total travel distance and the total transport cost [21]. This work has been extended, by proposing a new approach based on the use of Pareto dominance method to give a set of satisfying solutions to the m-PDPTW minimizing total travel cost, total tardiness time and the vehicles number. [22] [23]

3 Mathematical Formulation

Our problem is characterized by the following parameters:

- N: Set of customers, supplier and depot vertices,
- N': Set of customers and supplier vertices,
- N^+ : Set of supplier vertices,
- N^- : Set of customers vertices,
- n: Size of the initial population,
- K: Vehicle number,
- d_{ij} : Euclidian distance between the vertex i and the vertex j. If $d_{ij} = \infty$ then the road between i and j doesn't exist,
- t_{ijk} : Time used by the vehicle k to travel from the vertex I to the vertex j,
- $[e_i, l_i]$: Time window of the vertex i,
- s_i : Stopping time at the vertex i,
- q_i : Goods quantity of the vertex i request. If $q_i > 0$, the vertex i is a supplier; if $q_i < 0$, the vertex i is a customer and if $q_i = 0$ then the vertex was served.
- Q_k : Capacity of vehicle k,
- i = 0...N : Predecessor vertex index,
- j = 0...N : Successor vertex index,
- k: 1...K: Vehicle index,
- $Xijk = \begin{cases} 1 \text{ If the vehicle travel from the vertex i to the vertex j} \\ 0 \text{ Else} \end{cases}$
- A_i : Arrival time of the vehicle to the vertex i,
- D_i : Departure time of the vehicle from the vertex i,
- y_{ik} : The goods quantity in the vehicle k visiting the vertex i,
- C_k : Travel cost associated with vehicle k,
- A vertex is served only once,
- A vertex is served only once,
- The capacity constraint must be respected,
- The depot is the starting and finishing vertex for the vehicle,
- The vehicle stops at every vertex for a period of time to allow the request processing,
- If the vehicle arrives at a vertex i before its time windows beginning date it waits.

The function to minimize is given as follows:

$$Minimizef = \begin{pmatrix} \lambda_1 c_1 \sum_{k \in K} \sum_{i \in N} \sum_{j \in N} C_k d_{ijk} X_{ijk} + \\ \lambda_2 c_2 \sum_{k \in K} \sum_{i \in N} \sum_{j \in N} \max(0, D_i - l_i) X_{ijk} \end{pmatrix}$$
(1)

Where λ_i and c_i are weights and scaling coefficients. Subject to:

$$\sum_{i=1}^{N} \sum_{k=1}^{K} x_{ijk} = 1, j = 2, \dots N$$
(2)

$$\sum_{j=1}^{N} \sum_{k=1}^{K} x_{ijk} = 1, i = 2, \dots N$$
(3)

$$\sum_{i \in N} X_{i0k} = 1, \forall k \in K$$
(4)

$$\sum_{j \in N} X_{0jk} = 1, \forall k \in K$$
(5)

$$\sum_{i \in N} X_{iuk} - \sum_{j \in N} X_{ujk} = 0, \forall k \in K, \forall u \in N$$
(6)

$$X_{ijk=1\Rightarrow y_{jk}=y_{ik}+q_i,\forall i,j\in N;\forall k\in K}$$
(7)

$$y_{0k=0,\forall k\in K} \tag{8}$$

$$Q_k \ge y_{ik} \ge 0, \forall i \in N; \forall k \in K$$
(9)

$$D_w \le D_v, \forall w \in N^+; \forall v \in N^-$$
(10)

$$D_0 = 0 \tag{11}$$

$$X_{ijk} = 1 \Rightarrow D_i + t_{ijk} \le D_j \forall i, j \in N; \forall k \in K$$
(12)

The constraint (2) and (3) ensure that each vertex is visited only once by a single vehicle. The constraint (4) and (5) ensure that the vehicle routing is beginning and finishing in the depot. The constraint (6) ensures the routing continuity by a vehicle. (7), (8) and (9) are the capacity constraints. The precedence constraints are guaranteed by (10), (11) and (12).

4 Genetic Algorithm For Optimization Of The m-PDPTW

4.1 Generation of the initial populations

In our case, we generate two types of populations. A first population noted P_{node} (Figure 1), which represents all nodes to visit with all vehicles, according to the permutation list coding. The second population noted $P_{vehicle}$ (Figure 2) indicates nodes number visited by each vehicle.

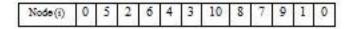


Figure 1: The permutation list coding

Иį	1/2	V3 V4		VS	
6	4	0	0	0	

Figure 2: Example of Individual from the population $P_{vehicle}$

4.2 Correction procedures

Before beginning the construction of the population $P_{node/vehicule}$, we proceed to the correction procedures of precedence and capacity between nodes. We consider the following couples customer / supplier: (1,5), (2,8), (9,7), (10,3) and (4,6), noting that $Q_{k \max} = 60$ and q = 20, we present, respectively, in figures 3 and 4 the principle of correction precedence and capacity.

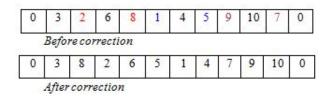


Figure 3: Correction precedence

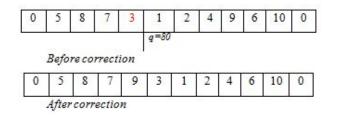


Figure 4: Correction capacity

4.3 Computation procedure

Taking into account the population P_{node} , correction procedures and $P_{vehicle}$ we illustrate in Figure 5 an example of an individual from the population $P_{node/vehicule}$. Knowing that it is necessary to verify that a couple is visited by only one vehicle. [24]

with: N' = 10 and K = 2

Figure 6 represents the process to determine the population $\mathrm{P}_{\mathrm{node/vehicule}}.$

The principles of different genetic operations such as crossover and mutation operator are detailed in our work [21].

Vi	c_{l}	0	5	1	8	2	6	4	0
V2	<i>c</i> ₂	0	3	7	10	9	0		

Figure 5: Example of an individual from the population $P_{node/vehicule}$

4.4 Multi-criteria evaluation

A multi-objective problem is defined as an optimization vector problem, which seeks to optimize several components of a vector function cost.

Pareto dominance method

A multi-criteria problem P is composed of n variables, m inequality constraints, p equality constraints and k criteria that can be formulated as follows:

$$P \Rightarrow \begin{cases} \min f(x) = [f_1, f_2, f_3, \dots, f_k(x)] \\ g_{i(x) \le 0i=0...m} \\ g_{j(x)=0j=0...p} \end{cases}$$
(13)

However, it is necessary to find solutions representing a possible compromise between the criteria. The Pareto optimality concept introduced by the economist V. Pareto in the twentieth century is frequently used [25]. V. Pareto formulated the following concept: in a multi-criteria problem, there is equilibrium so that we can not improve one criterion without deteriorating at least one other. This equilibrium has been called Pareto optimal A solution is noted Pareto optimal if it is dominated by any other point in solutions space. These points are noted non dominated solutions.

A point $X \in E$ dominates $Y \in E$ if:

$$\begin{cases} \forall i, f_i(x) \le f_i(y) \\ \text{and } \exists j, \text{such as } f_j(x) < f_j(y) \end{cases}$$
(14)

Figure 7 shows an example where we seek generations of the initial populations to minimize fland f2. The points 1, 3 and 5 are not dominated. On the contrary point 2 is dominated by point 3, and point 4 is dominated by point 5.

Aggregation method

In the resolution of MOP (Multi objectives Problem), several traditional methods are transforming the MOP into a single objective problem. Among these methods we find the aggregation method. This is one of the first methods used to generate Pareto optimal solutions. It is to transform the problem (MOP) in a problem (PMO_{λ}) which combines the different cost functions of the problem into a single objective function F generally linear [26]:

$$F(x) = \sum_{i=1}^{n} c_i \lambda_i f_i(x) \tag{15}$$

Where λ_i and c_i are weights and scaling coefficients, according to the application, that the different objectives are not necessarily commensurable. The constants c_i are usually initialized to $\frac{1}{f_i(x^*)}$ where $f_i(x^*)$ is the optimal solution associated to the objective function f_i considered

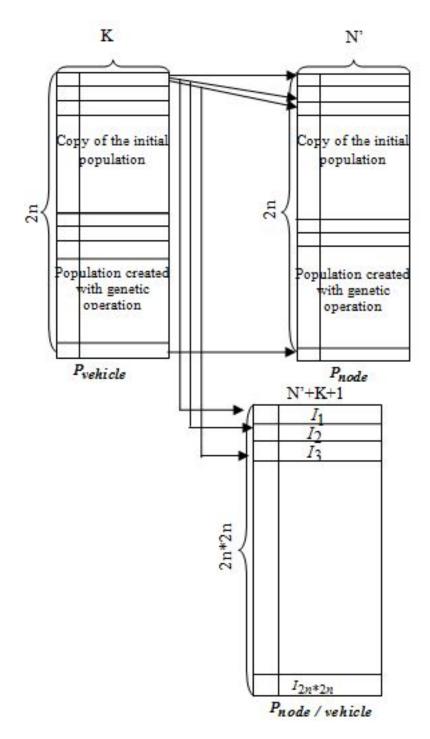


Figure 6: Computation procedure

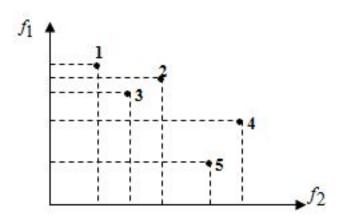


Figure 7: Dominance example

separately. The idea of the aggregation method (Figure 8) is to fix a weight vector i.e. find a hyper-plane in the objective space (a line for a bi-criteria problem) with a fixed orientation. The Pareto optimal solution is the point where the hyper-plane has a common tangent with a feasible space. The advantage of the aggregation method is to produce a single solution and thus do not require interaction with the decision maker.

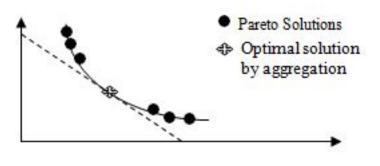


Figure 8:

Multi-objective optimization and computing of lower bounds

The computing of lower bounds has been studied in literature for several scheduling problems which we quote: the problems with a machine [27], with parallel machines [28] and [29], the hybrid flow-shop problems [30] and the flexible job shop problems [31]. These proposed methods for computing of the lower bounds are generally based on the relaxation of constraints (preemption of tasks, constraints related resources ...) to minimize one or more criteria for optimal scheduling. Based on these methods and seen that we don't have information on the optimal solutions associated with different cost functions f_i for our problem we should compute the minimum value to determine the scaling constants c_i . For this objective, we use the relaxation of various constraints. To find a minimum value associated with the criterion of total travel cost $f_1 = \sum_{k \in K} \sum_{i \in N} \sum_{j \in N} d_{ij}X_{ijk}$, we have treated this problem to the travelling salesman problem when we try to minimize $\sum_{i \in N} \sum_{j \in N} d_{ij}X_{ij}$. We subsequently determine the routing crossed by a single vehicle, minimizing the total travel distance by incorporating the constraints and capacity precedence. What gives us $d_{\min.p.c}$.

$$d_{\min.p.c} = \min(\sum_{i \in N} \sum_{j \in N} d_{ij} X_{ij})$$
(16)

By setting K, the number of vehicles used, we get:

$$d_{\min k.p.c} = \frac{d_{\min .p.c}}{K} \tag{17}$$

Consequently, we acquire a value f_{1b} that represents a minimum value of the total travel cost.

$$f_{1b} = \sum_{k \in K} C_k d_{\min k.p.c} \tag{18}$$

To determine the minimum value of the total tardiness time, we fix K the vehicles number and find out the population $P_{node/vehicule}$. Thus, we calculate the total tardiness time for each individual and determine the minimum.

$$f_{2b} = \sum_{k \in K} \sum_{i \in N} \sum_{j \in N} \max(0, D_i - l_i) X_{ijk}$$
(19)

Knowing that for the criterion of tardiness, a better lower bound is zero. We will have therefore:

$$\begin{cases} c_2 = \frac{1}{f_{2b}} \\ s.c.f_{2b} \neq 0 \end{cases}$$
(20)

4.5 Computational results

To test our approach, we use benchmark problem instances generated by Li and Lim [19] from Solomon's ones [32]. Corresponding to Solomon's classification of C1, C2, R1, R2, RC1 and RC2, their data sets were also generated in six classes: LC1, LC2, LR1, LR2, LRC1 and LRC2. The LC problems are clustered whereas in the LR problems, providers and customers are randomly generated. Therefore in the LRC problems the providers and the customers are partially clustered and partially randomly distributed. While LC1, LR1 and LRC1 problems have a short scheduling horizon, LC2, LR2 and LRC2 have longer scheduling one. [33] In our work, we consider a vehicle number k ranging between 1 and 25. Table 1 shows the results of our simulation using the parameters of the problem LRC1. Of course, for every given solution, we note the corresponding routing, crossed by each vehicle.

 N_{sol} : represents the number of non dominated solutions.

 N_k : represents the vehicles number used.

We observe that our approach generates a multiple number of solutions that give flexibility of choice for the decision maker and that by using two different methods to determine the vehicles number used, minimizing the compromise between the total travel cost and the total tardiness time. We also observe that we obtain a total tardiness equal to zero with a tolerable cost.

5 Conclusion

In this paper, we have presented our approach to solve the m-PDPTW, based on Pareto dominance method, with use of genetic algorithm and lower bounds. Our purpose was in a first part a brief literature review on the VRP, 1-PDPTW and m-PDPTW. The mathematical formulation of our problem is detailed in second part. Then, we have detailed the use aggregation

LRC1	N_{sol}	N_k	f_{1b}	f_{2b}	f_1	f_2	F(x)
LRC101	4	11	67971,97	39,08	187085,73	2,173	1,09
					192150,78	$2,\!09$	
					$192843,\!93$	$1,\!32$	
					194911,78	$0,\!651$	
LRC101	1	25	67971,97	0	216261,26	0	2,49
LRC103	2	11	645,1144	0	1946,2465	0	1,09
					$1711,\!304$	$0,\!429$	
LRC103	1	25	645,1144	0	2081,1978	0	2,49
LRC105	3	9	697,0005	0,823	1862,0582	4,28	0,89
					$1965,\!2465$	0,823	
					1947,6537	$4,\!27$	
LRC105	2	25	697,0005	0	2171,1193	0,107	2,49
					$2246,\!6729$	0	
LRC107	1	11	647,345	0	1803,9515	0	1,09
LRC107	1	25	647,345	0	2171,1006	0	2,49

Table 1: Results for the LRC1 problem

method and lower bounds to determine a set of solutions, minimizing our objective functions. Simulation was presented in a last part by using benchmark's data.

Bibliography

- H. Ezzedine T. Bonte C. Kolski and C. Tahon. Integration of traffic management and traveller information systems: basic principles and case study in intermodal transport system management. International Journal of Computers, Communications & Control (IJCCC), ISSN 1841 9836, E-ISSN 1841-9844 Vol. III, No. 3, pp. 281-294, 2008.
- [2] Christofides N. Mingozzi A. and Toth P. The vehicle routing problem. In Combinatorial Optimization, volume 11, pages 315-338. John Wiley, 1979.
- [3] Lenstra J. and Rinnooy Kan A. Complexity of the vehicle routing and scheduling problems. In Networks, volume 11, pages 221-228. Springer, 1981.
- [4] Nabaa M. Zeddini B. and Tranouez P. Approche décentralisée pour résoudre le problème du transport r la demande. Schedae, prépublication n° 23, (fascicule n° 2, p. 133-140), 2007.
- [5] Toth P. and Vigo D. The vehicle routing problem, SIAM Monographs on Discrete Mathematics and Applications, 2002, ISBN 0-89871-498-2.
- [6] Montamenni R. Gambardella L.M. Rizzoli A.E. and Donati A.V. A new algorithm for a Dynamic Vehicle Routing Problem based on Ant Colony System. IDSIA, Switzerland, 2002.
- [7] Sorin C. Negulescu, Claudiu V. Kifor, and Constantin Oprean. Ant Colony Solving Multiple Constraints Problem: Vehicle Route Allocation. Int. J. of Computers, Communications & Control (IJCCC), ISSN 1841-9836, E-ISSN 1841-9844, Vol. III, No. 4, pp. 366-373, 2008.

- [8] Kaoru Hirota, Fangyan Dong and Kewei Chen. A Computational Intelligence Approach to VRSDP (Vehicle Routing, Scheduling, and Dispatching Problems). ICCCC, Baile Felix-Oradea, Romania pp. 53-54, 2006.
- [9] Savelsbergh M.P.W. and SOL M., "The general pickup and delivery problem", Transportation Science 1995.
- [10] Psaraftis H.N. An exact algorithm for the single vehicle many to many immediate request dial a ride problem with time windows. Transportation Science, 17, 351-357, 1983.
- [11] Psaraftis H.N. A dynamic programming solution to the single vehicle many to many immediate request dial a ride problem. Transportation Science, 14(2):130-154, 1980.
- [12] Jih. W. and Hsu J. Dynamic vehicle routing using hybrid genetic algorithms. International Conference on Robotics and Automation, pages 453-458, 1999.
- [13] Velasco N. Dejax P. Gueret C. and Prins C. Un algorithme génétique pour un problème de collecte bi-objectif. MOSIM 2006.
- [14] Kammarti. R. Hammadi. S. Borne. P. and Ksouri. M. A new hybrid evolutionary approach for the pickup and delivery problem with time windows. IEEE International Conference on Systems, Man and Cybernetic. 2004. Volume 2, P 1498-1503, Oct 2004.
- [15] Kammarti. R. Hammadi. S. Borne. P. and Ksouri. M. "Improved tabu search in an hybrid evolutionary approach for the pickup and delivery problem with time windows", Intelligent Transportation Systems, 2005. Proceeding 2005 IEEE, p 148-153, 2005.
- [16] Kammarti. R. Hammadi. S. Borne P. and Ksouri. M. "Solving the Real Time dynamic Pickup and Delivery Problem with an hybrid evolutionary approch". Multiconference on Computational Engineering in Systems Application. Volume 2, P 1520-1525, Oct 2006.
- [17] Sol M. and Savelsbergh M. A branch-and-price algorithm for the pickup and delivery problem with time windows. Memorandum COSOR 94-22, Dept. of Mathematics and Computing Science, Eindoven University of Technology, Eindoven, The Netherlands, 1994.
- [18] Quan L. and Maged M. A new insertion-based construction heuristic for solving the pickup and delivery problem with time windows. Science direct, European Journal of Operational Research 2003.
- [19] Li H. and Lim A. A metaheuristic for the pickup and delivery problem with time windows. In IEEE International Conference on Tools with Artificial Intelligence, volume 13, pages 160-167, 2001.
- [20] Li H. Lim A. and Rodrigues B.. Solving the pickup and delivery problem with time windows using squeaky wheel" optimization with local search. SMU Business Conference Paper Series, 2002.
- [21] Harbaoui Dridi I. Kammarti R. Borne P. and Ksouri M. Un Algorithme génétique pour le problème de ramassage et de livraison avec fenetres de temps r plusieurs véhicules. CIFA 2008, Bucarest (Roumanie), Septembre 2008 Proc. Article 176.
- [22] Harbaoui Dridi I. Kammarti R. Borne P. and Ksouri M. Approche multicritère pour le problème de ramassage et de livraison avec fenetres de temps r´ plusieurs véhicules. Symposium LT'2009, Tunisie, (Sousse), LT09-006, Mars 2009.

- [23] Harbaoui Dridi I. Kammarti R. Borne P. and Ksouri M. Genetic Algorithm for Multicriteria Optimization of a Multi-Pickup and Delivery Problem with Time Windows. 13th INCOM IFAC symposium, Moscow (Russia), pp 1521-1526, June 2009.
- [24] Harbaoui Dridi I. Kammarti R. Ksouri M and Borne P. A Genetic Algorithm for the Multi-Pickup and Delivery Problem with time windows. Studies in Informatics and Control (SIC), Vol 18, N°2, pp 173-180, Juin 2009.
- [25] Pareto. V. Cours d'économie politique, Lausanne : Rouge, 1896-7, reproduit in Vilfredo Pareto, Oeuvres complètes, Genève : Librairie Droz, 1964.
- [26] Hwang, C. and Masud, A. Multiple objective decision making methods and applications. In Lectures Notes in Economics and Mathematical Systems, volume 164. Springer-Verlag, Berlin, 1979.
- [27] Jouglet. A, Baptiste. B et Carlier. J. Exact Procedures for Single Machine Total cost Scheduling. Proceeding of the IEEE / SMC'02 Conf. 6-9 October, Hammamet, Tunisia, 2002.
- [28] Jurich, B. Scheduling Jobs in Shops with Multi-purpose Machines. Ph.D thesis, Fachbereich Mathematik / Informatik, Universitat Osnabruck., 1992.
- [29] Carlier, J. Scheduling jobs with release dates and tails on identical machines to minimize the makespan. European Journal of Operational Research, 29 (1987) 298-306, North-Holland.
- [30] Billaut, J.C, Carlier, J, Néron, E. Ordonnancement d'ateliers r ressources multiples. Chapitre dans l'ouvrage : Ordonnancement de la production, Edition Hermès, France, 2001.
- [31] I. Kacem. Ordonnancement multicritères des job shops flexibles : formulation, bornes inferieures et approche éévolutionniste coopérative. Thèse en Automatique et informatique industrielle r´ l'université de Lille 1, 2003.
- [32] Solomon M.M., "Algorithms for the vehicle Routing and Scheduling Problem with Time Window Constraints", Operations Research, 41, 469-488, 1987.
- [33] Kammarti. R. Hammadi. S. Borne. P. and Ksouri. M. "Lower Bounds In An Hybrid Evolutionary Approach For The Pickup And Delivery Problem With Time Windows". IEEE International Conference on Systems, Man and Cybernetics. 2005. Volume 2, P 1156-1161, Oct 2005.

An Authenticated Key Agreement Protocol Using Isogenies Between Elliptic Curves

D. He, J. Chen, J. Hu

Debiao He, Jianhua Chen, Jin Hu Wuhan University

School of Mathematics and Statistics Wuhan, Hubei, 430072, China E-mail: hedebiao@163.com, chenjh-ecc@163.com, hujin-ecc@163.com

> **Abstract:** All the current public-key cryptosystems will become insecure when size of a quantum register is sufficient. An authenticated key agreement protocol, which is against the attack of quantum computer, is proposed. The proposed protocol can provide the security properties known session key security, forward security, resistance to key-compromise impersonation attack and to unknown key-share attack, key control. We also prove its security in a widely accepted model.

> **Keywords:** public-key cryptosystem; quantum computer; isogeny; elliptic curve; key agreement protocol.

1 Introduction

Key agreement is one of the fundamental cryptographic primitive after encryption and digital signature. Such protocols allow two or more parties to exchange information among themselves over an adversarially controlled insecure network and agree upon a common session key, which may be used for later secure communication among the parties. Thus, secure key agreement protocols serve as basic building block for constructing secure, complex, higher-level protocols. The first modern key agreement protocol is the Diffie-Hellman protocol given by Diffie and Hellman in the seminal paper [1] in 1976. Its security is based on the difficulty of solving discrete logarithm problems. As the first practical key agreement protocol without authentication, it is not secure to the man-in-the-middle attack. After that, lots of protocols have been published and some of them use certificates generated by the trusted third parties (public key infrastructures, PKIs) to prevent attacks such as the man-in-the-middle attack. Most of the systems based on PKI are complex and expensive for the cost of the authentication, refresh and revocation of certificates. Security of the known key agreement protocols is based on two general mathematical problems: determination of order and structure of a finite Abelian group, and discrete logarithm computation in a cyclic group with computable order. Both of the problems can be solved in polynomial time using Shor's algorithm for a quantum computer [2]. Thus, most of the current public-key cryptosystems will become insecure when size of a quantum register is sufficient. Development of key agreement protocols, which would be strong against a quantum computer, is necessary.

The rest of our paper is organized as follows. Section 2 describes theoretical background and a public-key encryption technique. Section 3 analyses the complicity of the problem of searching for an isogeny between elliptic curves. We give the proposed key agreement protocol in Cection4, and analyse the security of the proposed protocol in Section 5. We conclude this paper in Section 6.

2 Elliptic Curves over F_p and Isogeny Star

Let E be a elliptic curve, defined on the finite fields F_p , and it's equation is

$$y^2 = x^3 + ax + b, a, b \in F_p$$
 (1)

Then the map

$$\pi: (x, y) \to (x^p, y^p) \tag{2}$$

specifes the Frobenius endomorphism of the curve E. A Frobenius map satisfies its characteristic equation

$$\pi^2 - T\pi + p = 0 \tag{3}$$

where $T = p - a - \sharp E(F_p)$ is the Frobenius trace. Through the Hasse's theorem, we know that

$$|T| < 2\sqrt{p}).\tag{4}$$

So the discriminant D_{π} of the Frobenius equation (3) satisfies

$$D_{\pi} = T^2 - 4p < 0 \tag{5}$$

Theorem 1 Elliptic curves are isogenous over F_p if and only if they have equal number of points.

Proof. See[3].

Theorem 2 Let an elliptic curve $E(F_p)$ have the Frobenius discriminant D_{π} and $(\frac{D_{\pi}}{l})$ be a Kronecker symbol for some *l*-degree isogeny. If $(\frac{D_{\pi}}{l}) = -1$, then there are no *l*-degree isogenies; if $(\frac{D_{\pi}}{l}) = 1$, then two *l*-degree isogenies exist; if $(\frac{D_{\pi}}{l}) = 0$, then 2 or l + 1 *l*-degree isogenies exist and *l* is called Elkies prime number.

Proof. See[3].

Let $U = E_i(F_p)$ be a set of elliptic curves with equal number of points, so that each element of U is uniquely determined by a j-invariant of an elliptic curve. According to the theorem 1 and the equation (4), we can consider U as a category, and the set of isogenies between elements of U as a set of morphisms of this category. We can compute $\sharp U = h_{D_{\pi}}$, where $h_{D_{\pi}}$ is the degree of Hilbert polynomial[3].

Let l is Elkies prime number, we can get that there are two isogenous elliptic curves for any elliptic curves of U, from theorem 2. It is practically determined that, when $\sharp U$ is prime, all the elements of U form a single isogeny cycle.

Let $l_1 \neq l$ be one more prime isogeny degree with the property that $\left(\frac{D_{\pi}}{l_i}\right) = 1$. In this case, l_i -degree isogenies form a cycle over U as well. Then we can put the l and l_i degree isogeny cycles over each other. Same can be done for other isogeny degrees of such kind.

Definition 1. A graph, consisted of prime number of elliptic curves, connected by isogenies of degrees satisfying $\left(\frac{D_{\pi}}{l_i}\right) = 1$, is an isogeny star.

If an isogeny star is wide enough, we can use it for crypto algorithm constructing. For that purpose, it is necessary to specify a direction on a cycle and route of isogeny stat. The method for direction determination on an isogeny cycle is mentioned in [4], we don't give the detail here.. It uses impact of Frobenius endomorphism on an isogeny kernel. The definition of isogeny stat is following.

Let S be an isogeny star, $L = \{l_i\}$ -a set of Elkies isogeny degrees being used and $F = \{\pi_i\}$ -a set of Frobenius eigenvalues, which specify positive direction for every $l_i \in L$.

Definition 2. A set $R = r_i$, where r_i is number of steps by the l_i -isogeny in the direction $F = \pi_i$, is a route on the isogeny star.

We can define composition[3] of routes $A = \{a_i\}$ and $B = \{b_i\}$ as $AB = \{a_i + b_i\}$. Routes are commutative: AB = BA.

The computation of iosgeny between elliptic curve can be done using the method in [5, 6, 7], we don't give the detail here.

3 Complexity of Isogeny Search

The are several techniques can be used for isogeny search[3]:

- Brute-force: Complexity of these attacks is estimated at isogeny computations.
- Meet-in-the-middle: Complexity of the attack is estimated at isogeny computations.
- Method described in [8]: Complexity of the attack is estimated at isogeny computations.

The reason of the problem of searching for an isogeny between elliptic curves can against the attack of quantum computer is following[3].

In order to computer the isogenies between elliptic, we must solve the equation

$$\phi(X,j) = 0,\tag{6}$$

and the process of computing the isogeny cycle is following

$$E_1 \to \phi_{l1}(X, j_{E1}) = 0 \to j_{E2} \to E_2 \to \phi_{l2}(X, j_{E2}) = 0 \to j_{E3} \to \dots$$
(7)

To compute a chain of q isogenies, one should consecutively solve these q equations, because of the equation parameter (j-invariant) is changed with every step. So one can't parallelize and the problem is against the attack of quantum computer.

Then, we conclude that the complexity of searching for an isogeny between elliptic curves is $O(\sqrt{n}) \approx O(\sqrt[4]{p})$, and the problem can be against the attack of quantum computer.

From the above discussion, the decisional Diffie-Hellman assumption can be easily extent to the isogenies through the property of isogenies between the elliptic curves.

Definition 3. The decisional Diffie-Hellman assumption over isogeny star (DDHA-IS): DDHA-IS is that it is difficult to distinguish the following real Diffie-Hellman distribution

$$\Gamma_{real} = \{R_1(E_{init}), R_2(E_{init}), R_1R_2(E_{init}) | R_1, R_2 \in G\}$$
(8)

and random Diffie-Hellman distribution

$$\Gamma_{rand} = \{R_1(E_{init}), R_2(E_{init}), R_3(E_{init}) | R_1, R_2, R_3 \in G\}$$
(9)

More formally, if we define the advantage function $Adv_G^{DDH-IS}(A)$ as

$$Adv_G^{DDH-IS}(A) = |Pr[A(X) = 1|X \in \Gamma_{real}] - Pr[A(Y) = 1|Y \in \Gamma_{rand}]|,$$
(10)

we say that the *DDHA-IS* holds in set G if $Adv_G^{DDH-IS}(A)$ is negligible for any probabilistic polynomial time adversary A.

4 Key Agreement Protocol Based on Isogeny

In this section we describe the proposed key agreement protocol which is specified by the key generation and the protocol description.

4.1 Key Generation

In this paper we use an elliptic curve E defined over a finite field F_p , the parameters is following.

1) F_p : the finite field;

2) E_{init} : an initial elliptic curve, its equation is $y^2 = x^3 + a_{init}x^3 + b_{init}, a_{init}, b_{init} \in F_p$ Lť 3)d:number of isogeny degrees being used;

 $4L = l_i, 1 \le i \le d$: a set of Elkies isogeny degrees being used;

 $5)F = \pi_i, 1 \leq i \leq d$: a set of Frobenius eigenvalues, which specify the positive direction for every $l_i \in L$;

6)k:a limit for number of steps by one isogeny degree in a root. For any root $\{r_i\}$, numbers of steps are selected in $-k \leq r_i \leq k$;

7)H:SHA-1;

8)Select random routes R_{privA} and R_{privB} . The value R_{privA} is a secret key of the user A, and R_{privB} is the secret key of the user B;

9)Compute the curves $E_{pubA} = R_{privA}(E_{init})$ and $E_{pubB} = R_{privB}(E_{init})$, which are the public key of a user A and B, respectively;

4.2 Our Key Agreement Protocol

Let E be the elliptic curve, defined on the finite field, with the equation (1), and let A_E and B_E be its parameter and j be its j-invariant. The proposed protocol is following.

1) A generate random route R_A and computes $E_A = R_A(E_{init}), e_A = H(A_{E_A}, B_{E_A})$. At last, A sends $M_1 = \{A_{E_A}, B_{E_A}, e_a\}$ to B.

2) Upon receiving M_1 , B checks whether e_A equals $H(A_{E_A}, B_{E_A})$. If not, B stops the session. Otherwise, B generate random route R_B and computes $E_B = R_B(E_{init}), E'_B = R_B(E_A), E''_B = R_{privB}(E_A), e_B = H(A_{E_B}, B_{E_B}, A_{E''_B}), B_{E''_B})$. At last, B sends $M_2 = \{A_{E_B}, B_{E_B}\}$ to A.

3) Upon receiving M_2 , A computes $E'_A = R_A(E_B)$, $E''_A = R_A(E_{pubB})$, $E''_A = R_A(E_{pubB})$ and checks whether the equation $e_B = H(A_{E_B}, B_{E_B}, A_{E''_A})$, $B_{E''_A})$ holds or not. If it does not hold, then A terminates the execution. Otherwise, A computes $e'_A = H(A_{E''_A}), B_{E''_A})$ the session key $sk_{AB} = H(A_{E'_A}) \oplus B_{E'_A})$. At last, A sends $M_3 = e'_A$ to B.

4) Upon receiving \tilde{M}_3 , B computes $E''_B = R_B(E_{pubA})$ and checks whether the equation $e'_A = H(A_{E''_B}), B_{E''_B})$ holds or not. If it does not hold, then B terminates the execution. Otherwise, computes the session key $sk_{BA} = H(A_{E'_B}) \oplus B_{E'_B})$.

As a result, A and A achieve the same shared secret key:

$$E'_{A} = R_{A}(E_{B}) = R_{A}(R_{B}(E_{init})) = R_{B}(R_{A}(E_{init})) = R_{B}(E_{A})$$
(11)

$$sk_{AB} = H(A_{E'_A}) \oplus B_{E_A}) = H(A_{E'_A}) \oplus B_{E'_A}) = sk_{AB}$$
 (12)

and authenticate each other.

5 Security analysis

5.1 Security model

In this work we shall use a modified Bellare-Rogaway key exchange model [9, 10] to analyse the protocol security. In the model, each party involved in a session is treated as an oracle, and an adversary can access the oracle by issuing some specified queries (defined later). An oracle $\Pi_{i,j}^s$ denotes the *s*-th instance of party *i* involved with a partner party *j* in a session.

The security of a protocol is defined by a game with two phases. In the first phase, an adversary E is allowed to issue the following queries in any order.

1)Send($\Pi_{i,j}^s, m$). Upon receiving the message m, oracle $\Pi_{i,j}^s$ executes the protocol and responds with an outgoing message m or a decision to indicate accepting or rejecting the session. If the oracle $\Pi_{i,j}^s$ does not exist, it will be created as initiator if $m = \lambda$, or as a responder otherwise.

 $2)Reveal(\Pi_{i,j}^s)$. If the oracle has not accepted, it returns \perp ; otherwise, it reveals the session kev.

3)Corrupt(i). The party responds with its private key.

Once the adversary decides that the first phase is over, it starts the second phase by choosing a fresh oracle and issuing a query, where the fresh oracle and query are defined as follows.

Definition 4 (fresh oracle) An oracle $\Pi_{i,j}^s$ is fresh if $(1)\Pi_{i,j}^s$ has accepted; $(2)\Pi_{i,j}^s$ is unopened (not been issued the query); (3) party $j \neq i$ is not corrupted (not been issued the *Corrupt* query); (4) there is no opened oracle $\Pi_{j,i}^t$, which has had a matching conversation to . The above fresh oracle definition is particularly defined to cover the key-compromise impersonation resilience property since it implies that the user could have been issued a query.

4) $Test(\Pi_{i,j}^s)$.Oracle $\Pi_{i,j}^s$ which is fresh, as a challenger, randomly chooses $b \in \{0,1\}$ and responds with the session key, if b = 0, or a random sample from the distribution of the session key otherwise.

After this point the adversary can continue querying the oracles except that it cannot reveal the test oracle $\Pi_{i,j}^s$ or its partner $\Pi_{j,i}^t$ (if it exists), and it cannot corrupt party j. Finally, the adversary outputs a guess b'' for b. If b' = b, we say that the adversary wins. The adversary's advantage is defined as

$$Adv^{E}(k) = |2Pr[b' = b] - 1|.$$
(13)

We use the session ID which can be the concatenation of the messages in a session to define matching conversations, i.e., two oracles $\Pi_{i,j}^s$ and $\Pi_{j,i}^t$ have matching conversations to each other if they have the same session ID.

A secure authenticated key (AK) agreement protocol is defined as follows.

Definition 5 Protocol Π is a secure AK if:

1). In the presence of a benign adversary, which faithfully conveys messages, on $\Pi_{i,j}^s$ and $\Pi_{j,i}^t$, both oracles always accept holding the same session key, and this key is distributed uniformly on $\{0,1\}^k$;

2). For any polynomial time adversary E, $Adv^{E}(k)$ is negligible.

5.2 Security analysis

Using the above security definitions, we have the following Theorem 1.

Theorem 3. In the random oracle, if DDHA-IS is hard, our proposed protocol is a secure AK protocol.

Proof: The first two conditions follow immediately from the description of our proposed protocol and the assumption that H is random oracle.

Let's turn to the second condition. We use the method proposed by Pan et al.[10] to analyze the security. Consider there exists an adversary E and $Adv^{E}(k)$ is non-negligible. Suppose its running time is t. We will use E to construct an algorithm F which can solve DDHA-IS. Let l, N and q_H be the number of sessions related to E, the number of entitys and the queries' numbers of $H(\cdot)$ made by E.

Given (E_1, E_2, E_3) , where $E_1 = R_1(E_{init}), E_2 = R_2(E_{init}), E_3 = R_3(E_{init})$. First, F selects randomly i, j, r, s and generates the public/private key pair for every entity. Then, F starts E and answers all queries made by E.

H query: If the input exists, answers by its corresponding value. Otherwise, F picks a random number as the answer to the new query, and adds the input, output pair at the end of the H-string;

Corrupt query: Because F knows all entities' private keys, F can answers by the corresponding private key;

Reveal query: If the input is $\Pi_{i,j}^s$ or $\Pi_{j,i}^t$, F selects b' as its output and halts. Suppose the input is $\Pi_{i',j}^{s'}$ or $\Pi_{j,i'}^{t'}$, where $(i',s') \neq (i,t)$ and $(i',t') \neq (i,t)$. Because F knows all entity's private keys and simulates the run of i' and j according to the protocol, F can get random numbers selected by i' and j. So, F knows the session key of $\Pi_{i',j}^{s'}$ or $\Pi_{j,i'}^{t'}$, and answers by this session key.

Send query: If the input is $\Pi_{i,j}^s$ and an empty string, F computes $e_A = H(A_{E_1}, B_{E_1})$ and answers by A_{E_1}, B_{E_1}, e_A . If the input is $\Pi_{i,j}^s$ and a string that is not none, F computes $E''_A = R_{privA}(E_2)$ and answers by $\{e'_A\}$. If the input is $\Pi_{j,i}^t$ and a string that is not none, computes $E''_B = R_{privB}(E_1), e_B = H(A_{E_2}, B_{E_2},), A''_{E_2}, B''_{E_2}$ and answers by A_{E_2}, B_{E_2}, e_B . Else, answers by random numbers according to the protocol.

Test query: If the input is not $\Pi_{i,j}^s$, F outputs $b' \in \{0,1\}^k$ and halts. Else, if A_{E_3} and B_{E_3} exist in H-string, let the corresponding value be r; Else, F selects r randomly and appends $\{A_{E_3}, B_{E_3}, r\}$ to the H-string. F answers by r.

When E halts, F outputs E's output b' and halts.

1)Suppose $\Pi_{i,j}^s$, selected by F, is the input of the *Test* oracle, and $\Pi_{i,j}^s$ and $\Pi_{j,i}^t$ have matching conversations. If E_3 is the DH value of E_1, E_2, r is the session key. Else r is a random number. Because E_3 is the DH value of E_1, E_2 with the probability $\frac{1}{2}$ and F answers all queries made by E correctly, the probability of the event that distinguishes r and the session key correctly is equivalent to the probability of the event that F decides whether (E_1, E_2, E_3) is a DH triple.

2) Suppose $\Pi_{i,j}^s$, selected by F, is the input of the *Test* oracle, or $\Pi_{i,j}^s$ and $\Pi_{j,i}^t$ have not matching conversations. F outputs correctly with the probability $\frac{1}{2}$.

Suppose the success probability of E is $\frac{1}{2} + \epsilon$, where ϵ is non-negligible. Because the first event happens with the probability $\frac{1}{l^2}$, the success probability of F in this condition is $\frac{1}{2l^2} + \frac{\epsilon}{l^2}$. The second event happens with the probability $1 - \frac{1}{l^2}$. So the success probability of in this condition is $\frac{1}{2l^2} - \frac{1}{2l^2}$. From the above discussion, we know the success probability of F is $\frac{1}{2} + \frac{\epsilon}{l^2}$. Then we can solve the DDHP-IS with non-negligible probability. This is a contradiction. So our proposed protocol satisfies the second condition. i.e., the protocol is a secure AK protocol.

5.3 Other discussion

A number of desirable attributes of key agreement protocols have also been identified [9] and nowadays most protocols are analyzed with such attributes. Here, the following six security properties must be considered for the proposed protocol: a known-key security, perfect forward secrecy, a key-compromise impersonation attack, a unknown key-share security, a key-control security. Regarding the above mentioned definitions, the following theorems are used to analyze the six security properties of the proposed protocol. Our protocol also satisfies the following security notions which are often used to judge the security of key agreement protocols.

Known session key security: A protocol is called Known session key security, if an adversary, having obtained some previous session keys, cannot get the session keys of the current run of the key agreement protocol. In our scheme, the agreed session key relies on the one-way hash function and session secrets. The output of hash function is distributed uniformly in $\{0, 1\}^k$, thus one session key which is the output of hash function has no relation with the others. Besides, the session key is generated with the session secrets which are computed from the random ephemeral key, thus even one session's session secrets are revealed, the other session secrets will still remain safe.

Perfect forward secrecy: A protocol is called Perfect forward secrecy, if compromise of the three private keys of the participating entities does not affect the security of the previous session keys. Even if an attacker gets the value the secret key R_{privA} and R_{privB} in our scheme, he can't deduce E'_A or E'_B , without the knowledge of the two random numbers R_A and R_B . Therefore, our scheme can provide perfect forward secrecy.

No key-compromise impersonation: The compromise of one entity's static private key does not imply that the private keys of other entities will also be compromised in our protocol. The adversary may impersonate the compromised entity in subsequent protocols, but he cannot impersonate other entities. This property is called no key-compromise impersonation.

First, suppose an attacker C obtains the long-term private key R_{privA} from the compromised user A. In order for the key-compromise impersonation attack to succeed, C must know A's ephemeral keys. In this case, C would also have to extract from 's ephemeral public value R_A , so as to generate the same session key with A. C, however, will face the problem of searching for an isogeny between elliptic curves. Therefore, the proposed protocol is secure against a key-compromise impersonation attack.

No unknown key-share: If the adversary convinces a group of entities that they share some session key with the adversary, while in fact they share the key with another entity, we call the protocol as suffering from unknown key-share attack. To implement such an attack on our protocol, the adversary is required to learn the private key of some entity. Otherwise, the attack hardly works. Hence, we claim that our protocol has the attribute of no unknown key-share.

No key control: No key-control security means that neither entity can't force the session key to a preselected value. From the execution of the proposed protocol, we know that the only possibility of key-control attack may be brought out by the participant of the protocol B. But, for the party B to make the party A generate the session key K_B which is preselected value by B, B should solve the equation $E'_B = R_B(E_A)$. This is the problem of searching for an isogeny between elliptic curves. Therefore, the proposed protocol provides a no key-control security.

6 Conclusion

In this paper, we propose a secure and efficient authenticated key agreement, which works on the isogeny star. We prove that our protocol meets the security attributes under the assumption that the problem of searching for an isogeny between elliptic curves is secure.

7 Acknowledgments

The authors would like to thank the anonymous reviewers for their constructive comments.

Bibliography

- W. Diffie and M. Hellman, New directions in cryptography, *IEEE Trans. Info. TH*, vol. 22, pp.644-654, 1976.
- [2] Boneh D., Lipton R. Quantum cryptanalysis of hidden linear functions. Proceedings of the 15th Annual International Cryptology Conference on Advances in Cryptology (LNCS 963), 1995:424-437.
- [3] Rostovtesv A. and Stolbunov A., Public-key cryptosystem based on isogenies. Cryptology ePrint Archive, Report 2006/145, 2006. http://eprint.iacr.org/.
- [4] Couveignes J. M., Dewaghe L., Morain F. Isogeny cycles and the schoof-elkies-atkin algorithm. Ecole polytechnique, France, 1996.
- [5] Elkies N., Elliptic and modular curves over finite fields and related computational issues, Proceedings of a Conference in Honor of A.O.L. Atkin, AMS International Press, 1998, pp.21-76.
- [6] Muller V., Ein Algorithmus zur Bestimmung der Punktanzahl elliptisher Kurven uber endlichen Korpern der Charakteristik groser drei, 1995. http://www.informatik.tudarmstadt.de/ti/forschung/ecc.
- [7] F.Morain, E.Schost, Fast Algorithms for Computing Isogenies between Elliptic Curves. http://www.lix.polytechnique.fr/ morain/jcomp.pdf, 2006.
- [8] S. Galbraith. Constructing isogenies between elliptic curves over finite fields, *Journal of Com*putational Mathematics, vol. 2, pp.118-138, 1999.
- [9] S. Blake-Wilson, D. Johnson and A. Menezes, Key Agreement Protocols and Their Security Analysis, Proceedings of Sixth IMA International Conference on Cryptography and Coding, Circnecester, UK, 1997, pp. 30-45.
- [10] H. Pan, J.-F. Li, Q.-S. Zheng, A Provable-Security Mutual Authenticated Key Agreement Protocol for Mobile Communication, *The 4th International Conference on Wireless Communications*, Networking and Mobile Computing, 2008, pp.1-4.

On the Power of Small Size Insertion P Systems

A. Krassovitskiy

Alexander Krassovitskiy University Rovira i Virgili Research Group on Mathematical Linguistics Av. Catalunya, 35, Tarragona 43002, Spain E-mail: alexander.krassovitskiy@estudiants.urv.cat

Abstract: In this article we investigate insertion systems of small size in the framework of P systems. We consider P systems with insertion rules having one symbol context and we show that they have the computational power of context-free matrix grammars. If contexts of length two are permitted, then any recursively enumerable language can be generated. In both cases a squeezing mechanism, an inverse morphism, and a weak coding are applied to the output of the corresponding P systems. We also show that if no membranes are used then corresponding family is equal to the family of context-free languages. **Keywords:** P systems, insertion-deletion systems, context-free languages, matrix grammars, computational completeness

1 Introduction

The study of insertion-deletion operations on strings has a long history; we just mention [2,5,11,16,20]. Insertion-deletion systems motivated from molecular computing have been studied in [1,3,10,19,21]. With some linguistic motivation they may be found in [9].

In general form, an insertion operation means adding a substring to a given string in a specified (left and right) context, while a deletion operation means removing a substring of a given string from a specified (left and right) context. A finite set of insertion/deletion rules, together with a set of axioms provide a language generating device: starting from the set of initial strings and iterating insertion-deletion operations as defined by the given rules we get a language. The number of axioms, the length of the inserted or deleted strings, as well as the length of the contexts where these operations take place are natural descriptional complexity measures of the insertion-deletion systems.

Inspired by the structure and the functioning of a living cell, especially by the local information processing, P systems are brought to light few years ago [15]. This is a highly distributed computational model which combines local processing (in membranes) and communication (between them). It is natural to consider insertion and deletion operations in the framework of P systems and it was firstly done by Gh. Paun in [17]. Such a combination permits do define programmed-like insertion-deletion systems, which, as expected, increase their computational power. Some combinations of parameters for pure insertion-deletion lead to systems which are not computationally complete [6, 12] or even decidable [22], while in [7,8] it was shown that P systems framework can easily increase the computational power comparing to ordinary insertiondeletion systems.

Traditionally, language generating devices having only insertion or only deletion rules were studied. Early computational models based only on insertion appear already in [9], and are discussed in [19] and [17] (with membrane tree structure). It was proved that pure insertion systems having one letter context are always context-free. Yet, there are insertion systems with two letter context which generate nonsemilinear languages (see Theorem 6.5 in [19]). On the other hand, it appears that by using only insertion operations the obtained language classes with contexts greater than one are incomparable with many known language classes. For example, there is a simple linear language $\{a^n b a^n \mid n \ge 1\}$ which cannot be generated by any insertion system (see Theorem 6.6 in [19]).

In order to overcome this obstacle one can use some codings to "interpret" the generated strings. The computational power of insertion systems with morphisms and intersection with special languages was investigated in [14] and [18]. In [11] there were used two additional mapping relations: a morphism h and a weak coding φ . The strings of a language are considered being the products $h^{-1} \circ \varphi$ over the generated strings. More precisely, the squeezing mechanism selects only those output strings of the corresponding (P) systems to which h^{-1} and φ are being applied. As expected, the language generating mechanisms have greater expressivity, and the corresponding language family is larger. It appears that with the help of morphisms and codings one can obtain every *RE* language if insertion rules have sufficiently large context. It is proved in [11] that for every recursively enumerable language *L* there exists a morphism *h*, a weak coding φ and a language *L'* generated by an insertion system with rules using the length of the contexts at most 7, such that, $L = h(\varphi^{-1}(L'))$. The result was improved in [13], showing that rules having at most 5 letter contexts are sufficient to encode every recursively enumerable language. Recently, in [4] it was shown that the same result can be obtained with the length of contexts equal to 3.

Our aim is to reduce the length of the contexts in insertion rules by regulating derivations in terms of membranes. Unlike the previous works, our article considers the encoding as a part of insertion P systems. The obtained model is quite powerful and has the power of matrix languages if contexts of length one are used. We also show that if no encoding is used, then the corresponding family is strictly included into the family of matrix languages and is equal to the family of context-free languages if no membranes are used. If an insertion of two symbols in two letters contexts is used, then all recursively enumerable languages can be generated (using of course the inverse morphism, the weak coding, and the squeezing mechanism).

2 Prerequisites

Here we use standard formal language theoretic notions and notations. The reader can consult any of the many monographs in this area (see, e.g., in [20] for the unexplained details).

We denote by |w| the length of a word w and by card(A) the cardinality of the set A, while ε denotes the empty string.

By CF and RE we denote the classes of context-free and recursively enumerable languages, respectively. A language L is context-free if there exists a context-free grammar G such that L(G) = L. A context-free grammar is a construct G = (N, T, S, P), where N and T are disjoint alphabets, $S \in N$, and P is a finite set of context-free rules of the form $A \longrightarrow v$, where $A \in N$ and $v \in (N \cup T)^*$. We say that a context-free grammar G = (N, T, P, S) is in Chomsky normal form, if each rule in P has one of the forms $A \longrightarrow \alpha$, or $A \longrightarrow BC$, for $A, B, C \in N, \alpha \in T \cup \{\varepsilon\}$.

A language L is recursively enumerable if there exists a type 0 grammar G such that L(G) = L. A type 0 grammar is a construct G = (N, T, S, P), where N and T are disjoint alphabets, $S \in N$ and P is a finite set of rules of the form $u \longrightarrow v$, where $u \in (N \cup T)^* N(N \cup T)^*$ and $v \in (N \cup T)^*$.

We say that a type 0 grammar G = (N, T, P, S) is in Penttonen normal form, if each rule in P can be written either as $A \longrightarrow \alpha$, or $AB \longrightarrow AC$, for $A, B, C \in N, \alpha \in (T \cup N)^*, |\alpha| \le 2$. It is well known that for every type 0 language there is a modified Penttonen normal form where each rule can be written as either: $A \longrightarrow \alpha$, or $A \longrightarrow AC$, or $A \longrightarrow CA$, or $AB \longrightarrow AC$, or $AB \longrightarrow CB$, for $A, B, C \in N, \alpha \in T \cup \{\varepsilon\}$.

We also recall the following definition from [17]. A context-free matrix grammar (without appearance checking) is a construct G = (N, T, S, M), where N, T are disjoint alphabets (of non-terminals and terminals, respectively), $S \in N$ (axiom), and M is a finite set of matrices, that is sequences of the form $(A_1 \longrightarrow x_1, \ldots, A_n \longrightarrow x_n), n \ge 1$, of context-free rules over $N \cup T$. For a string w, a matrix $m : (r_1, \ldots, r_n)$ is executed by applying the productions r_1, \ldots, r_n one after the other, following the order in which they appear in the matrix. Formally, we write $w \Longrightarrow_m u$ if there is a matrix $m : (A_1 \longrightarrow u_1, \ldots, A_n \longrightarrow u_n) \in M$ and the strings $w_1, w_2, \ldots, w_{n+1} \in (N \cup T)^*$ such that $w = w_1, w_{n+1} = u$, and for each $i = 1, 2, \ldots, n$ we have $w_i = w'A_iw''$ and $w_{i+1} = w'u_iw''$. If the matrix m is understood, then we write \Longrightarrow instead of \Longrightarrow_m . As usual, the reflexive and transitive closure of this relation is denoted by \Longrightarrow^* . Then, the language generated by G is $L(G) = \{w \in T^* \mid S \Longrightarrow^* w\}$. The family of languages generated by context-free matrix grammars is denoted by MAT^{λ} (the superscript indicates that erasing rules are allowed). It is well known fact that every language from MAT^{λ} can be generated by a modified binary normal form, (similarly to the binary normal form, see e.g., [17]) having each matrix m of the following form $m : (A \longrightarrow \alpha, A' \longrightarrow \alpha')$, for $A, A' \in N, \alpha, \alpha' \in (N \cup T)^*, max(|\alpha|, |\alpha'|) \le 2$.

An insertion system is a construct $\Gamma = (V, A, I)$, where V is an alphabet, A is a finite language over V, and I is finite set of triples of the form (u, α, v) , where u, α , and v are strings from V^* . The elements of A are called *axioms*, the triples in I are *insertion rules*. An insertion rule $(u, \alpha, v) \in I$ indicates that the string α can be inserted in between u and v. Stated otherwise, $(u, \alpha, v) \in I$ corresponds to the rewriting rule $uv \to u\alpha v$. We denote by \Longrightarrow the relation defined by an insertion rule (formally, $x \Longrightarrow y$ iff $x = x_1 uv x_2, y = x_1 u\alpha v x_2$, for some $(u, \alpha, v) \in I$ and $x_1, x_2 \in V^*$). We denote by \Longrightarrow^* the reflexive and transitive closure of \Longrightarrow (as usual, \Longrightarrow^+ is its transitive closure).

The language generated by Γ is defined by

$$L(\Gamma) = \{ w \in V^* \mid x \Longrightarrow^* w, x \in A \}.$$

We say that an insertion system (V, A, I) has weight (n, m, m') if

$$m = \max\{|u| \mid (u, \alpha, v) \in I\}, \ m' = \max\{|v| \mid (u, \alpha, v) \in I\},\$$
$$n = \max\{|\alpha| \mid (u, \alpha, v) \in I\}.$$

We denote by $INS_n^{m,m'}$ the corresponding families of languages generated by insertion systems.

An *insertion P system* is a construct:

$$\Pi = (V, \mu, M_1, \ldots, M_k, R_1, \ldots, R_k),$$

where

- V is a finite alphabet,
- μ is a (cell-like, i.e., hierarchical) membrane structure with k membranes. This structure will be represented by a word containing correctly nested marked parentheses, i.e., by a word from the Dyck language. The skin membrane is labeled with "1".
- M_i , for each $1 \le i \le k$ is a finite language associated to the membrane *i*.
- R_i , for each $1 \leq i \leq k$ is a set of insertion rules with target indicators associated to membrane *i* and which have the following form: (u, x, v; tar), where (u, x, v) is an insertion rule, and tar, called the target indicator, is from the set {here, in_j , out}, $1 \leq j \leq k$.

A configuration of Π is a k-tuple (N_1, \ldots, N_k) of finite languages over V. For two configurations (N_1, \ldots, N_k) and (N'_1, \ldots, N'_k) of Π we write $(N_1, \ldots, N_k) \Longrightarrow (N'_1, \ldots, N'_k)$ if we can pass from (N_1, \ldots, N_k) to (N'_1, \ldots, N'_k) by applying nondeterministically the insertion rules, to all strings which can be rewritten from the corresponding regions, and following the target indications associated with the rules. More specifically, $w' \in N'_j$ if either $w' \in N_j$ or there is a word $w \in N_i$ and $w \Longrightarrow_r w'$, where $r = (u, x, v; tar) \in R_i$. Moreover, the membrane labeled by j is immediately outside the membrane labeled by i if tar = out; the membrane labeled by j is immediately below the membrane labeled by i if $tar = in_j$; and i = j if tar = here. No other words are present in N'_j , $1 \le j \le k$. We say that a word w' is sent out of the system if there is a configuration (N_1, \ldots, N_k) , a word $w \in N_1$, and $w \Longrightarrow^r w'$ where $r = (u, x, v; out) \in R_1$.

We use the definition of the language generated by Π according to [17]. This language is denoted by $L(\Pi)$ and it is defined as follows: we start from the initial configuration (M_1, \ldots, M_k) of the system and proceed iteratively, by transition steps performed by applying the rules in parallel, to all strings which can be rewritten. All strings over the alphabet V sent out of the system (*i.e.* sent from the skin membrane) during any step of any computation form the language $L(\Pi)$.

Insertion tissue P systems are defined in an analogous manner. As the tissue P systems use arbitrary graph structures we write the target indicator in the form $tar = go_j, j = 1, ..., k$. We remark, that the result of a computation consists of all strings over V which are sent to one selected output cell.

The weight of insertion rules (n, m, m') and the membrane degree k describe the complexity of the system. We denote by $LSP_k(ins_n^{m,m'})$ (see [17]) the family of languages $L(\Pi)$ generated by insertion P systems of degree at most $k \ge 1$, having weight at most (n, m, m'). If some of the parameters n, m, m', or k is not specified we write "*" instead.

We say that a language L' is from $MorINS_n^{m,m'}$ (from $MorLSP_k(ins_n^{m,m'})$, respectively) if there exist a morphism h, a weak coding φ and $L \in INS_n^{m,m'}$ ($L \in LSP_k(ins_n^{m,m'})$) such that $\varphi(h^{-1}(L)) = L'$.

For every language $L \in MorLSP_k(ins_n^{m,m'}), L = \varphi(h^{-1}(L(\Pi)))$, we add h and φ to the system, and we write Π in the form

$$\Pi = (V, \mu, M_1, \dots, M_n, R_1, \dots, R_n, h, \varphi),$$

Hereafter h and φ are naturally extended to strings: $h(a_1a_2...a_t) = h(a_1)h(a_2)...h(a_t)$, and $\varphi(a_1a_2...a_t) = \varphi(a_1)\varphi(a_2)...\varphi(a_t), a_i \in V.$

We insert "t" before P to denote classes of languages corresponding to the tissue cases (e.g., LStP). We also write "[t]" (e.g., LS[t]P) if we do not distinguish between tissue and tree classes.

We say that a letter a is marked in a sentential form waw'' if it is followed by #, i.e., |w''| > 0, and # is the prefix of w''. In the following proofs we use a marking technique introduced in [17]. The technique works as follows: in order to simulate a rewriting production $A \longrightarrow B$ we add adjacently right from A the word #B specifying that letter A is already rewritten. As soon as the derivation of the simulated sentential form is completed, every nonterminal A is marked and the pair A# is subject to the inverse morphism.

3 Main results

Let us consider insertion systems (without membranes) with one letter context rules, i.e., the family $MorINS_*^{1,1}$. Applying the marking technique we get a characterization of context-free languages.

Theorem 3.1. $MorINS^{1,1}_* = CF.$

Proof: First we show that $CF \subseteq MorINS_3^{1,1}$.

Let G = (V, T, S, P) be a context-free grammar in Chomsky normal form. Consider the following insertion system

 $\Pi = (T \cup V \cup \{\#\}, R, \{S\}, h, \varphi),$

where $R = \{(A, \#\gamma, \alpha) \mid \alpha \in T \cup V, A \longrightarrow \gamma \in P, \gamma \in (T \cup V)^*, 1 \le |\gamma| \le 2\}$, the morphism h is defined as follows

$$h(a) = a \#$$
, if $a \in V$, and $h(a) = a$, if $a \in T$,

and the weak coding φ is defined as follows

$$\varphi(a) \longrightarrow \varepsilon$$
, if $a \in V, \varphi(a) \longrightarrow a$, if $a \in T$.

It is clear that $L(\Pi) \in MorINS_3^{1,1}$. We claim that $L(\Pi) = L(G)$. Indeed, each rule $(A, \#\gamma, \alpha) \in R$ can be applied in the sentential form wA alphaw' if A is unmarked (not rewritten). Hence, the production $A \longrightarrow \gamma \in P$ can be simulated by the corresponding derivation of G. Hence, by applying the counterpart rules we get equivalent derivations. At the end of the computation every nonterminal is marked, and no rules can be applied any more. (Indeed, if the system produces a word having some unmarked nonterminal then h^{-1} is not defined on this word.) At this point h^{-1} removes all marks, and φ removes all nonterminal symbols. Hence $L(\Pi) = L(G)$. We get $CF \subseteq MorINS_3^{1,1}$.

The equivalence of these two classes follows from Theorem 6.4 in [19] stating that $INS_*^{1,1} \subseteq CF$ and the fact that the family of context-free languages is closed under inverse morphisms and weak codings.

Now we consider insertion P systems with the left and right contexts of at most one letter. It is known from Theorem 5.5.1 in [17] that $LSP_2(ins_2^{1,1})$ contains non context-free languages. We prove that the more general family $LSP_*(ins_*^{1,1})$ is bounded by the class of languages generated context-free matrix grammars:

Lemma 3.2. $LStP_*(ins_*^{1,1}) \subset MAT^{\lambda}$.

Proof: The proof uses a similar technique as in [19], Theorem 6.4 for context-free grammars.

Let $\Pi = (V, \mu, M_1, \dots, M_n, R_1, \dots, R_n)$ be an insertion P system such that $L(\Pi) \in LStP_n(ins_*^{1,1})$ for some $n \ge 1$.

Consider the matrix grammar $G = (D \cup Q \cup \{S\}, V, S, P)$, where $Q = \{Q_i \mid i = 1, ..., n\}$, $D = \{D_{a,b} \mid a, b \in V \cup \{\varepsilon\}\}$, and P is constructed as follows:

1. For every rule $(a, b_1 \dots b_k, c, go_j) \in R_i$, $a, c \in V \cup \{\varepsilon\}, b_1, \dots, b_k \in V, k > 0$ we add to P $(Q_i \longrightarrow Q_j, D_{\overline{a},\overline{c}} \longrightarrow D_{\overline{a},b_1}D_{b_1,b_2} \dots D_{b_{k-1},b_k}D_{b_k,\overline{c}})$, where

$$\overline{a} = \begin{cases} a, & \text{if } a \in V, \\ t, \forall t \in V \cup \{\varepsilon\}, & \text{if } a = \varepsilon \end{cases} \qquad \overline{c} = \begin{cases} c, & \text{if } c \in V, \\ t, \forall t \in V \cup \{\varepsilon\}, & \text{if } c = \varepsilon. \end{cases}$$

- 2. For every rule $(a, \varepsilon, c, go_j) \in R_i$, $a, c \in V \cup \{\varepsilon\}, k > 0$ we add to P $(Q_i \longrightarrow Q_j, D_{\overline{a},\overline{c}} \longrightarrow D_{\overline{a},\overline{c}})$, where \overline{a} and \overline{c} are defined as in the previous case.
- 3. Next, for every $w = b_1 \dots b_k \in M_i$, $i = 1, \dots, n, k > 0$ we add to P the matrix $(S \longrightarrow Q_i D_{\varepsilon, b_1} D_{b_1, b_2} \dots D_{b_{k-1}, b_k} D_{b_k, \varepsilon}).$
- 4. As a special case if $\varepsilon \in M_i$ we add $(S \longrightarrow Q_i D_{\varepsilon,\varepsilon})$ to P.
- 5. Also, for every $D_{a,b} \in D, a, b \in V \cup \{\varepsilon\}$ we add $(D_{a,b} \longrightarrow a)$ to P.

6. Finally, we add $(Q_1 \longrightarrow \varepsilon)$ to P (we assume that the first cell is the output cell).

The simulation of Π by the matrix grammar is straightforward. We store the label of current cell by means of nonterminals from Q. Every nonterminal $D_{a,c} \in D, a, c \in V \cup \{\varepsilon\}$ represents a pair of adjacent letters, so we can use them as a context. A rule $(a, b_1 \dots b_k, c, go_j) \in R_i, a, c \in$ $V, b_1 \dots b_k \in V^k$ can be simulated by the grammar iff the sentential form contains both Q_i and $D_{a,c}$. As a result, the label of the current cell is rewritten to Q_j and $D_{a,c}$ is rewritten to the string $D_{a,b_1}D_{b_1,b_2}\dots D_{b_{k-1},b_k}D_{b_k,c}$. We note, that in order to simulate rules that have no context we introduce productions with every possible context symbols by writing $\overline{a} \in V \cup \{\varepsilon\}$ and $\overline{c} \in V \cup \{\varepsilon\}$. Clearly, since indexed symbols are duplicated for adjacent nonterminals, e.g., $D_{b_{i-1},b_i}D_{b_i,b_{i+1}}$, every nonterminal thus preserves one symbol right and left contexts.

The simulation of Π by the grammar starts with a nondeterministic choice of the axiom from M_1, \ldots, M_n . Then, during the derivation any rule from R_1, \ldots, R_n having the context (a, b) is applied in one to one correspondence with grammar productions having $D_{a,b}$ in the left hand side. Finally, the string over V is produced by the grammar iff Q_1 has been deleted from the simulated sentential form. The deletion of Q_1 specifies that Π reached the output cell. So, we obtain $L(\Pi) = L(G)$. Hence, $LStP_*(ins_*^{1,1}) \subseteq MAT^{\lambda}$.

The strictness of the inclusion follows from the fact there are languages from MAT^{λ} which cannot be generated by any insertion P system from $LStP_*(ins_n^{1,1})$, for any $n \geq 1$. Indeed, consider the context-free language $L_a = \{ca^k ca^k c \mid k \geq 0\}$. Since every context-free language is a matrix language we have $L_a \in MAT^{\lambda}$. On the other hand, $L_a \notin LStP_*(ins_n^{1,1})$, for any $n \geq 1$. For the contrary, assume there is such a system Π' . We note, that the system cannot delete or rewrite any letter, so every insertion is terminal. As the languages of axioms are finite we need an insertion rule of letter a. Consider the final insertion step in a derivation which has at most one step and derives a word $w = ca^k ca^k c$, for some $k \geq n+1$:

$$w_0 \Longrightarrow^* w' \Longrightarrow w,$$

where w_0 is an axiom. Since $|w_0|_c \leq 3$, c may be inserted by the last insertion. Assume, that $|w'|_c = 3$. In the latter case, let a^p be the inserted string, $p \leq n$. Because, we may insert a^p in the distinct positions of w' we get that either $ca^{k-p}ca^{k+p}c \in L(\Pi')$ or $ca^{k+p}ca^{k-p}c \in L(\Pi')$. This is a contradiction.

Now assume that c is inserted by the last insertion. We note that the insertion of two c is not possible, since $k \ge n + 1$. Consider three cases: (1) the last applied rule inserts c in the middle, (2) at the end, or (3) at the beginning of w'.

(1) Let $w_c = a^{p'}ca^{p''}$ be the inserted string, where $p' + p'' \leq n - 1$. Hence, $w' = ca^{k'+k''}c$, where k' + p' = k'' + p'' = k, and $k' + k'' = 2k - p' - p'' \geq 2n + 2 - n + 1 \geq 4$. Obviously, regardless of the contexts of the last insertion rule there are at least two positions at which w_c can be inserted. So, we get a contradiction because either $ca^{k'+p'+1}ca^{k''+p''-1}c \in L(\Pi')$, or $ca^{k'+p'-1}ca^{k''+p''+1}c \in L(\Pi')$.

(2) Let $a^q c$ be the inserted string, where $q \leq n-1$. The corresponding insertion rule has one of the following forms: $(\varepsilon, a^q c, \varepsilon, go_j)$ or $(a, a^q c, \varepsilon, go_j)$, where j is an index of the final membrane. In ether case, $a^q c$ may be inserted in w' before the last letter a. This is a contradiction. The case (3) is a mirror to the case (2) and can be treated similarly.

So we proved $L_a \notin LStP_*(ins_n^{1,1})$, for any $n \ge 1$ and, hence, $LStP_*(ins_*^{1,1}) \subset MAT^{\lambda}$. \Box

Since trees are special cases of graphs we get the following result

Corollary 3.3. $LSP_*(ins_*^{1,1}) \subset MAT^{\lambda}$.

Lemma 3.4. $MAT^{\lambda} \subseteq MorLSP_{*}(ins_{2}^{1,1}).$

We prove the lemma by direct simulation of a context-free matrix grammar Proof: G = (N, T, S, P). We assume that G is in the modified binary normal form, i.e., every matrix has the form $i: (A \longrightarrow BC, A' \longrightarrow B'C') \in P$, where $A, A' \in N; B, B', C, C' \in N \cup T \cup \{\varepsilon\}$, and i = 1, ..., n.

Consider a P insertion system Π defined as follows:

$$\Pi = (V, [1 [2 [3 [4]]_4 \dots [n+3]_{n+3}]_3]_2]_1, \{S\$\}, \emptyset, \dots, \emptyset, R_1, \dots, R_{n+3}, h, \varphi),$$

where $V = N \cup T \cup \{C_i, C'_i \mid i = 1, ..., n\} \cup \{\#, \$\}.$ For every matrix $i : (A \longrightarrow BC, A' \longrightarrow B'C')$ we add

$r.1.1: (A, \#C_i, \alpha, in_2),$		to R_1 ;
$r.2.1:(C_i,BC,\alpha,in_3),$	$r.2.2:(C_i',\#,\alpha,out)$	to R_2 ;
$r.3.1: (C_i, \#, \alpha, in_{i+3}),$	$r.3.2:(C_i',B'C',\alpha,out)$	to R_3 ;
$r.i + 3.1 : (A', \#C'_i, \alpha, out),$		to R_{i+3}

for every $\alpha \in V \setminus \{\#\}$. In addition we add $(\varepsilon, \$, \varepsilon, out)$ to R_1 .

We define the morphism h and the weak coding φ by:

$$h(a) = \begin{cases} a, & \text{if } a \in T, \\ a\# & \text{if } a \in V \setminus (T \cup \{\#\}) \end{cases} \qquad \varphi(a) = \begin{cases} a, & \text{if } a \in T, \\ \varepsilon & \text{if } a \in V \setminus T \end{cases}$$

Clearly, $L(\Pi) \in MorLSP_{n+3}(ins_2^{1,1})$. We claim that $L(\Pi) = L(G)$. To prove this it is enough to prove that $w \in L(G)$ iff $w' \in L(\Pi)$ and $w' = \varphi(h^{-1}(w))$.

First we show that for every $w \in L(G)$ there exists $w' \in L(\Pi)$ and $w' = \varphi(h^{-1}(w))$. Consider the simulation of the *i*-th matrix $(A \longrightarrow BC, A' \longrightarrow B'C') \in P$. It is controlled by letters C_i and C'_i . First, we insert $\#C_i$ in the context of unmarked A and send the obtained string to the second membrane. Then we use C_i as a context to insert adjacently right the word BC. After that, we mark the control letter C_i and send sentential form to the (i+3)-rd membrane. Here we choose nondeterministically one letter A', mark it, write adjacently right the new control letter C'_i , and, after that, send the obtained string to the third membrane. (We remark, the third membrane is immediately outside of the i + 3-rd membrane.) It is clear that it is not possible to apply the rule $r.i + 3.1 : (A', \#C'_i, \alpha; out)$ in the (i + 3)-rd membrane and to reach the skin membrane if the sentential form does not contain unmarked A'. So, this branch of computation cannot influence the result and may be omitted in the consideration. Next, in the third membrane, B'C' is inserted in the context of unmarked C'_i and the form is sent to the second membrane. Finally, we mark C'_i and send the resulting string back to the skin membrane.

At the beginning of the simulation the sentential form in the skin membrane does not contain unmarked C_i, C'_i . Hence, the insertions in the second and third membranes are deterministic. The derivation preserves this property, as after the sentential form is sent back to the skin membrane, introduced C_i , and C'_i are marked. At the end of the computation we send the resulting form out from the system by the rule $(\varepsilon, \$, \varepsilon, out)$.

Let w be a string in the skin region which contains some unmarked A and A'. If the letter A precedes A' then we can write $w = w_1 A \alpha_1 w_2 A' \alpha_2 w_3$. The simulation of the matrix is the following

$$w_1 A \alpha_1 w_2 A' \alpha_2 w_3 \xrightarrow{r.1.1, r.2.1, r.3.1} w_1 A \# C_i \# B C \alpha_1 w_2 A' \alpha_2 w_3 \xrightarrow{r.3+i.1, r.3.2, r.2.2} w_1 A \# C_i \# B C \alpha_1 w_2 A' \# C'_i \# B' C' \alpha_2 w_3,$$

where $w_1, w_2, w_3 \in V^*, \alpha_1, \alpha_2 \in V \setminus \{\#\}$. We can write a similar derivation if A' precedes A.

Hence, as result of simulation of *i*-th matrix we get both A and A' marked and BC, B'C'inserted in the correct positions. The derivation in Π may terminate by the rule $(\varepsilon, \$, \varepsilon, out)$ only in the first membrane. Hence it guaranties that the simulation of each matrix has completed. According to the definition of MorLSP the string w' belongs to the language if $w' = \varphi(h^{-1}(w))$, where w is the generated string. We may consider only the final derivations of Pi in which each nonterminal is marked. Hence, we have $L(G) \subseteq \varphi(h^{-1}(L(\Pi)))$.

The inverse inclusion is obvious since every rule in Π has its counterpart in G. Moreover the case when the derivation in Π is blocked corresponds to the case in which the simulation of a matrix cannot be completed.

Hence, we get $MAT^{\lambda} \subseteq MorLSP_{*}(ins_{2}^{1,1})$.

Remark 3.5. One can mention that a similar result can be obtained with a smaller number of membranes at the price of the maximal length of inserted words. Precisely, for any context-free matrix grammar G' there is a P insertion system Π' such that $L(\Pi') \in MorLSP_{n+1}(ins_3^{1,1})$ and $L(G') = L(\Pi')$, where n is the number of matrices in G'. To prove this we can use the same argument as in the previous theorem and replace rules in R_1, \ldots, R_{n+3} by

$$(A, \#BC, \alpha, in_{i+1})$$
to $R_1,$
(A', #B'C', α, out) to $R_{i+1},$

for every $\alpha \in V \setminus \{\#\}, i = 1, \ldots, n$.

Since trees are special cases of graphs we get the following result

Corollary 3.6. $MAT^{\lambda} \subseteq MorLStP_{*}(ins_{2}^{1,1}).$

Taking into account Lemma 3.4 and 3.2, and the fact that the class of languages generated by context-free matrix grammars is closed under inverse morphisms and weak codings we get a characterization of MAT:

Theorem 3.7. $MorLS[t]P_*(ins_*^{1,1}) = MAT^{\lambda}$.

Now we increase the maximal size of the context of insertion rules to two letters. It is known from [19] that $INS_2^{2,2}$ contains non-semilinear languages. By considering these systems with membrane regulation we obtain the following result

Theorem 3.8. $MorLSP_3(ins_2^{2,2}) = RE.$

Proof:

We prove the theorem by simulating a type 0 grammar in the modified Penttonen normal form. Let G = (N, T, S, P) be such a grammar. Suppose that rules in P are ordered and n = card(P).

Now consider the following insertion P system,

 $\Pi = (V, [1 \ [2 \ [3 \]_3 \]_2 \]_1, \{S\$\}, \emptyset, \emptyset, R_1, R_2, R_3, h, \varphi),$

where $V = T \cup N \cup F \cup \overline{F} \cup \{\#, \overline{\#}, \$\}, F = \{F_A, | A \in N\}, \overline{F} = \{\overline{F}_A, | A \in N\}.$

We include into R_1 the following rules:

$$(AB, \#C, \alpha, here), \text{ if } AB \longrightarrow AC \in P;$$

$$(A, \#C, B\alpha, here), \text{ if } AB \longrightarrow CB \in P;$$

$$(A, C, \alpha, here), \text{ if } A \longrightarrow AC \in P;$$

$$(\varepsilon, C, A\alpha, here), \text{ if } A \longrightarrow CA \in P;$$

$$(A, \#\delta, \alpha, here), \text{ if } A \longrightarrow \delta \in P;$$

$$(\$, \varepsilon, \varepsilon, out),$$

where $\alpha \in V \setminus \{\#\}$, $A, B, C \in N$, and $\delta \in T \cup \{\varepsilon\}$. It may happen that the pair of letters AB subjected to be rewritten by a production $AB \longrightarrow AC$ or $AB \longrightarrow CB \in P$ is separated by letters that have been marked. We use two additional membranes to transfer a letter over marked ones. In order to transfer $A \in N$ we add for each $\alpha \in V \setminus \{\#\}$

$$r.1.1: (A, \#F_A, \alpha, in_2),$$

to the skin membrane. Then we add to the second membrane

$r.2.1:(F_A,\#A,\alpha',out),$	$r.2.2: (\overline{F_A}, \overline{\#}A, \alpha', out),$
$r.2.3:(F_AX,\#\overline{F_A},\#,in_3),$	$r.2.4:(F_A\overline{F_B},\overline{\#}F_A,\overline{\#},in_3),$
$r.2.5:(F_A\overline{\#},F_A,\alpha,in_3),$	$r.2.6:(\overline{F_A}\#,F_A,\alpha,in_3),$
$r.2.7:(F_A\overline{\#},F_A,\overline{\#},in_3),$	$r.2.8:(\overline{F_A}\#,F_A,\overline{\#},in_3),$
$r.2.9:(F_A\overline{\#},\overline{F_A},\#,in_3),$	$r.2.10: (\overline{F_A}\#, \overline{F_A}, \#, in_3),$

for every

$$X \in F \cup N, \overline{F_B} \in \overline{F}, \alpha \in V \setminus \{\#, \overline{\#}\}, \\ \alpha' \in \{ab \mid a \in N \cup T, b \in N \cup T \cup \{\$\}\} \cup \{\$\}$$

Finally, we add to the third membrane the rules

$$r.3.1: (F_A, \#, \alpha, out),$$
 $r.3.2: (\overline{F_A}, \overline{\#}, \alpha', out),$

for every $\alpha \in V \setminus \{\#\}, \alpha' \in V \setminus \{\overline{\#}\}.$

The morphism h and the weak coding φ are defined as

$$h(a) = \begin{cases} a, & \text{if } a \in V \setminus N, \\ a \# & \text{if } a \in N. \end{cases} \qquad \varphi(a) = \begin{cases} a, & \text{if } a \in T, \\ \varepsilon & \text{if } a \in V \setminus T. \end{cases}$$

We simulate productions in P by marking nonterminals from N and inserting corresponding right hand sides of the productions. This can be done with insertions in the skin membrane by rules of weight (2, 2, 2) since the grammar has such a form that production rewrites/adds at most one letter.

The simulation of the transfer is done in the second and third membranes. The idea of the simulation is (1) to mark the nonterminal we want to transfer, (2) jump over the marked letters with the help of one special letter, at the end (3) mark the special letter and insert the original nonterminal. Since we use two letter contexts, in one step we can jump only over a single letter. We also need to jump over the marking letter # as well as over the marked nonterminals, and the letters inserted previously. In order to jump over # we introduce one additional marking symbol $\overline{\#}$. We mark letters from \overline{F} by $\overline{\#}$, and all the other letters in $V \setminus \{\overline{\#}, \#\}$ by #, e.g., in a word $\overline{F_A}\#$, letter $\overline{F_A}$ is unmarked.

(1) The rule $r.1.1 : (A, \#F_A, \alpha, in_2)$, specifies that every unmarked letter from N may be subjected to the transfer.

(2) The rules r.2.3 - r.2.10 in the second membrane specify that F_A or $\overline{F_A}$ is copied to the right in such a way that inserted letters would not be marked. In order to do so, the appropriate rule chooses to insert either the overlined copy $\overline{F_A}$ or the simple copy F_A . The rules r.2.3, r.2.4 describe jumps over one letter not in $\{\#, \overline{\#}\}$, and r.2.5 - r.2.10 describe jumps over $\#, \overline{\#}$. Every rule r.2.3 - r.2.10 sends the sentential form to the third membrane, and the rules r.3.1, r.3.2 in the third membrane send the sentential form back to the second membrane after marking one symbol $F_A \in F$ or $\overline{F_A} \in \overline{F}$.

(3) The rules r.2.1 and r.2.2 may terminate the transferring procedure and send the sentential form to the first membrane if letter \$ or two letters from $\{ab \mid a \in N \cup T, b \in N \cup T \cup \{\$\}\}$ appear in the right context.

For example, consider the transfer of A in the string AX # C\$ (here, we underline inserting substrings)

$$AX \#C\$ \xrightarrow{r.1.1} A \underline{\#F_A} X \#C\$ \xrightarrow{r.2.3} A \#F_A X \underline{\#F_A} \#C\$ \xrightarrow{r.3.1}$$
$$A \#F_A \underline{\#X} \#\overline{F_A} \#C\$ \xrightarrow{r.2.6} A \#F_A \#X \#\overline{F_A} \underline{\#F_A} C\$ \xrightarrow{r.3.2}$$
$$A \#F_A \#X \#\overline{F_A} \ \overline{\#} \#F_A C\$ \xrightarrow{r.2.1} A \#F_A \#X \#\overline{F_A} \ \overline{\#} \#F_A \underline{\#AC} \$.$$

The sentential form preserves the following property: (i) The first membrane does not contain unmarked letters from $F \cup \overline{F}$; there is exactly one unmarked letter from $F \cup \overline{F}$ in the second membrane; and there are always two unmarked letters from $F \cup \overline{F}$ in the third membrane.

We mention that property (i) is preserved by every derivation. Indeed, we start derivation from the axiom S^{\$} that satisfies the property, then one unmarked symbol is inserted by r.1.1. Rules r.2.3 - r.2.12 always add one more unmarked letter, whereas rules r.2.1, r.2.2, r.3.1, r.3.2always mark one letter from $F \cup \overline{F}$.

In order to verify that Π generates the same language as G we note that every reachable sentential form in G will be reachable also in Π by simulating the same production.

We also note that the derivation in Π may terminate by the rule $(\$, \varepsilon, \varepsilon, out)$ only in the first membrane. Hence, every transfer will be completed. It follows from property (i) that the simulation of the transfer is deterministic in the second membrane. Also note, that there is a nondeterministic choice in the third membrane, where the rules r.3.1, r.3.2 may mark one of the two unmarked letters. In the case the rule marks the rightmost letter, the derivation has to "jump" again over the inserted letter.

In a special case if r.1.1 starts the transfer of a letter adjacently left from an unmarked one then the rules r.1.1, r.2.1 produce two marked symbols which do not affect the result of the simulation.

The output string w is in the language, iff $w' = \varphi(h^{-1}(w))$ is defined. Hence, the resulting output of Π does not contain unmarked nonterminals. On the other hand every final derivation in Π has its counterpart in G. By applying the inverse morphism h^{-1} we filter out every sentential form with unmarked nonterminals from N. Hence, the corresponding derivation in G is completed. Finally, the weak coding φ filter away all supplementary letters. Hence, we have $L(G) = L(\Pi)$.

4 Conclusions

This article investigates the generative power of insertion P systems with encodings. The length of insertion rules and number of membranes are used as parameters of the descriptional complexity of the system. In the article we exploit the fact that a morphism and a weak coding are incorporated into insertion P systems. The obtained family $MorLS[t]P_*(ins_*^{1,1})$ characterizes the matrix languages. When no membranes are used, the class $MorINS_*^{1,1}$ is equal to the family of context-free languages. We proved also the universality result regarding the family $MorLSP_*(ins_*^{2,2})$. Also, for the family $LSP_*(ins_*^{2,2})$ one can get an analogous computational completeness result by applying right(left) quotient with respect to a regular language.

We recall the open problem posed in [4], namely, whether $MorINS_*^{2,2}$ is computationally complete. Our work gives a partial solution of the problem by using the membrane computing framework. One may see that in order to solve the problem completely (by the technique

promoted in the article), it is enough to find a concise way to transfer a letter over a marked context. In our case this can be reduced to the question whether it is possible to compute the membrane regulation in the skin membrane.

One may mention that there is a trade-off between the number of membranes and the maximal length of productions. By introducing additional nonterminals and fitting the grammars into the normal forms we decrease the amount of membranes used. On the other hand by raising the number of membranes we can simulate larger production rules. Moreover, the descriptional complexity used in the paper may be extended by such internal system parameters as, e.g., the size of the alphabet, the number of rules per membrane, etc. It may be promising to continue the research of the minimal systems regarding these parameters. We are also interested in the computational power of the insertion P systems having only one sided (right or left) contexts.

Acknowledgments

The author acknowledges the support PIF program from University Rovira i Virgili, and project no. MTM2007-63422 from the Ministry of Science and Education of Spain. The author sincerely thanks the help of Yurii Rogozhin and Seghey Verlan without whose supervision the work would not been done. Also the author warmly thanks Gheorghe Păun and Erzsèbet Csuhaj-Varjú who draw attention, during the meeting BWMC-09, to the point of using different normal forms of grammars to simulate P insertion-deletion systems.

Bibliography

- M. Daley, L. Kari, G. Gloor, R. Siromoney, Circular Contextual Insertions/Deletions with Applications to Biomolecular Computation, In: *Proc. SPIRE'99*, pp. 47-54, 1999.
- [2] L. Kari, From Micro-soft to Bio-soft: Computing with DNA, In: Proc. of Biocomputing and Emergent Computations, 146-164, 1997.
- [3] L. Kari, Gh. Păun, G. Thierrin, S. Yu, At the Crossroads of DNA Computing and Formal Languages: Characterizing RE Using Insertion-Deletion Systems. In: Proc. of 3rd DIMACS, pp. 318-333, 1997.
- [4] L. Kari, P. Sosík, On the Weight of Universal Insertion Grammars, TCS, 396(1-3), pp. 264-270, 2008.
- [5] L. Kari, G. Thierrin, Contextual Insertion/Deletion and Computability, Information and Computation, 131, 1, pp. 47-61, 1996.
- [6] A. Krassovitskiy, Yu. Rogozhin, S. Verlan, Further Results on Insertion-Deletion Systems with One-Sided Contexts, In Proc. of 2nd LATA08, LNCS, 5196, pp. 333-344, 2008.
- [7] A. Krassovitskiy, Yu. Rogozhin, S. Verlan, One-Sided Insertion and Deletion: Traditional and P Systems Case, *CBM08*, pp. 53-64, 2008.
- [8] A. Krassovitskiy, Yu. Rogozhin, S. Verlan, Computational Power of P Systems with Small Size Insertion and Deletion Rules, In: Proc of CSP08, pp. 137-148, 2008.
- [9] S. Marcus, Contextual Grammars, Rev. Roum. Math. Pures Appl., 14, pp 1525-1534, 1969.
- [10] M. Margenstern, Gh. Păun, Yu. Rogozhin, S. Verlan, Context-Free Insertion-Deletion Systems, TCS, 330, pp. 339-348, 2005.

- [11] C. Martin-Vide, Gh. Păun, A. Salomaa, Characterizations of Recursively Enumerable Languages by Means of Insertion Grammars, TCS, 205, 1–2, pp. 195-205, 1998.
- [12] A. Matveevici, Yu. Rogozhin, S. Verlan, Insertion-Deletion Systems with One-Sided Contexts, LNCS, 4664, pp. 205-217, 2007.
- [13] M. Mutyam, K. Krithivasan, A. S. Reddy, On Characterizing Recursively Enumerable Languages by Insertion Grammars, *Fundamenta Informaticae*, 64(1-4), pp. 317-324, 2005.
- [14] K. Onodera, New Morphic Characterization of Languages in Chomsky Hierarchy Using Insertion and Locality, In Proc. of 3rd LATA09, LNCS, 5457, pp. 648–659, 2009.
- [15] Gh. Păun, Computing with Membranes, Journal of Computer and System Sciences, Vol. 61, 1, pp. 108-143, 2000.
- [16] Gh. Păun, Marcus Contextual Grammars, Kluwer, Dordrecht, 1997.
- [17] Gh. Păun, Membrane Computing. An Introduction, Springer-Verlag, Berlin, 163, pp. 226-230, 2002.
- [18] Gh. Păun, M. J. Pérez-Jiménez, T. Yokomori, Representations and Characterizations of Languages in Chomsky Hierarchy by Means of Insertion-Deletion Systems, *International Journal of Foundations of Computer Science*, 19, 4, pp. 859-871, 2008.
- [19] Gh. Păun, G. Rozenberg, A. Salomaa, DNA Computing. New Computing Paradigms, Springer-Verlag, Berlin, 1998.
- [20] G. Rozenberg, A. Salomaa, eds., Handbook of Formal Languages, Springer-Verlag, Berlin, 1997.
- [21] A. Takahara, T. Yokomori, On the Computational Power of Insertion-Deletion Systems, DNA 2002, Sapporo, LNCS, 2568, pp. 269–280, 2003.
- [22] S. Verlan, On Minimal Context-Free Insertion-Deletion Systems, Journal of Automata, Languages and Combinatorics, 12, 1/2, pp. 317–328, 2007.

Modelling, Implementation and Application of a Flexible Manufacturing Cell

F. Leighton, R. Osorio, G. Lefranc

Felipe Leigthon

Universidad de Las Américas Vina del Mar - Chile E-mail: felipe.leigthon@gmail.com

Román Osorio

Universidad Nacional Autónoma de México Instituto IIMAS E-mail: roman@servidor.unam.mx

Gastón Lefranc

Pontificia Universidad Católica de Valparaíso Escuela de Ingeniería Eléctrica, Valparaiso - Chile E-mail: glefranc@ucv.cl

Abstract: This paper presents Petri Nets Model, the implementation and application of a Assembly Flexible Cell. The Cell is composed by a robotic manipulator, a computer vision system and a conveyor. The system is applied to assembly several products, showing only two of them.

Keywords: Flexible Manufacturing Cells, Assembly, CIM, Petri Nets.

1 Introduction

Within the area of the manufacturing, CIM model (Computer Integrated Manufacturing) permits the automation all the activities in a manufacturing company: the organization, storage (ASRS), manufacturing (CAM), engineering (CAE) and Quality Control. These activities are supported by local area networks and related database. In CAM includes fabrication, assembly and quality control of products, utilizing robotics manipulators, computer vision, etc., to facilitate automation [1]. These industrial systems, allow flexible their production and to automate their processes, without human intervention.

Assembly systems are widely used in several industries. These systems are composed by robotics manipulators, computer vision and sensors to complete the production task accurately. However, since these systems are composed of many components, it is difficult to deal with unexpected situations. For example, an undetected error may propagate and end up as a detectable failure which may cause the whole line to stop its operation. In this case, it may take considerable time to diagnose the system and identify the main reasons for the failure. There exist approaches to detect and to predict failures, but the usage of these approaches is limited. [2], [3]

The ability to produce a variety of products through the combination of modular components is a meaningful benefit of product modularity. There are five different ways that modular products are developed in industry. Component-swapping modularity can be achieved when two or more alternative types of components are paired with the same basic product body to create different product variants. This modularity is often associated with the creation of product variety as perceived by the customer. [6]

The Assembled Flexible Cell is part of a developed Flexible System in the Laboratory of Robotics, Artificial intelligence and Advanced Automation Lab, in the Escuela de Ingeniería Eléctrica de la Pontificia Universidad Católica de Valparaíso, Chile. FMS system includes: an ASRS cell (Automatic Storage Retrieval System) [4], with an executive system and database [5], Servoing Systems [7], [9] and Stereo computer vision [8], [10], [11] and Flexible Manufacturing cell. [12]

In this paper is presented an assembly cell modeled by Petri nets, the implementation and the application to the assembly of some products, in a manufacturing flexible cell. The Assembly Flexible Cell is composed by a robotic manipulator, a computer vision system and a conveyor.

2 Structure of the Assembly Cell

The assembled flexible cell of is constituted by a robotic arm, a system of vision and an system of automated transport. The system robotic manipulator used is IBM 7547 type Scara; the computer vision system is an stereo vision (two webcams); and the automated system of transport, is conformed by four conveyor belts that form a "circuit" around the robotic manipulator and the assembled area. It function is to interconnect the different cells from the system, like the ASRS (storage), the assembly flexible cell, manufacturing flexible cell and the flexible cell of quality control. This conforms a Flexible system of Manufacture, as it is seen in the (Fig.1). Through this system of transport are transferred pallets with different kits to assembly. In (Fig.2), is the transfer.



Figure 1: The Flexible Manufacturing System



Figure 2: Pallets with kit in the conveyor belt

The stereo vision system, (Fig.3), determines what pieces are in pallet, and the pallet position of the in the work area to make the assemly. In order to prove the cell a wood figure is assembled (geometric figure), and then making assembly of a electrical interrupt and electrical plug, common in houses and offices. The robotics manipulator takes the pieces inside the pallet and make the assembly programmed, The vision system contributes to the determination of the position and direction of the parts to be assembled.



Figure 3: The system of Vision

3 Aplications of Flexible Assembly Cell.

The applications of the cell are to assemble geometric figure (Fig.4a), to test if the system works. Then, making assembly of a electrical interrupt and electrical plug with one, two or three plug per unit (Fig.4b).



Figure 4: a) Kit of wood parts. b) Kit of plugs.

4 Computer Vision System.

The vision system have to determine the space position of pallets and objects within the space of work of the robotic manipulator. The Stereo Vision has to be able of visualizing objects of different size, send the order to manipulator to pick up the pallets and put in the working area, recognize the oarts inside the pallets and then send the order to manipulator to assemble the product in other pallet. Finally, the vision system verify is the product is well/assembled, giving the orders to manipulator to transport the pallet with assembled product ro the conveyor.

The images are captured by two webcam cameras. These images pass through a process of digitalization during which is made a discretización of the space coordinates of the image (x,y) denominated sampling of the image, and a discretización of the amplitude of the intensity of light in each point or píxel of the image, which denominates quantification of the gray level (Fig.5). The system make image segmentation, separating the different objects and applying a threshold based on the histogram (Fig.6).

Then, the center of mass or centroide of an object of a binary image is computed (Fig.7).

The direction of an object determines by the minus existing distance between its center of mass and an element of his contour. Therefore, the first step to follow to calculate the direction of an object is to determine which are pixeles that they belong to his contour. For doing this, a gradient of the image is used. So that an object is taken, the effector of the manipulator must be oriented of such way that their direction is perpendicular to the direction of the straight line

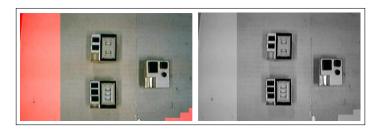


Figure 5: Digital image

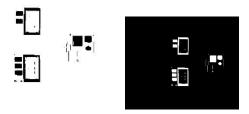


Figure 6: Segmentation and Thresholding

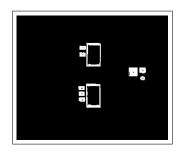


Figure 7: Centroide of the objects

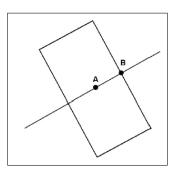


Figure 8: Direction of the object

that constitutes the minimum range between the center of mass of the object and its contour. This can be observed in following (Fig.8), where the point A corresponds to the center of mass of the rectangle and point B is determined finding the minimum range between the centroide and the contour of the object. The union of these points forms a perpendicular straight line to the direction that must have the effector of the manipulator to pick up the object.

5 Petri Nets Model of the Assembly Cell.

Assembly systems has many components and it is difficult to deal with unexpectes situations. Error non detected could propagate and to produce failure in the system, causing stop in the operation of the line. It takes time to diagnose, and to identify the reasons of the failure. A Petri Net Model (PN) of the Flexible Assembly Cell permits a simulation of the behaviour of the Assembly Cell, detecting failure nd an evaluation of the performance. The (Fig.9) shows the definitions of places and transitions and the (Fig.10), PN Model.

Places	Table 3: Place definition
	Arrive assembly order to cell
P ₁	Pallet in ASRS to CT
P ₂	Pallet in conveyor CT
P ₃	MR in initial position
P4	Pallet in place to be pick by MR over CT
P ₅	MR with pallet
P ₆	Pallet in workspace
P ₇	MR over pallet
P ₈	MR at 20 cm from workspace
Pg	Pallet in right position
P ₁₀	SV3D has digital image
P ₁₁	SV3D has position and orientation of objects
P ₁₂	MR has position and orientation of Obj1
P ₁₃	MR has obj 1
P ₁₄	Obj 1 in workspace
P15	MR has position and orientation obj 2
P16	MR has obj 2
P17	MR placed obj 2 over obj 1 in workspace
P18	Obj 2 assembled with obj 1
P19	MR at 20 cm from workspace
P20	MR has position and orientation obj 3
P21	MR has obj 3
P22	MR placed obj 3 over obj 1-2 in workspace
P23	Obj 3 assembled over obj 1-2
P24	MR has position and orientation obj 4
P25	MR has obj 4
P26	MR placed obj 4 over obj 1-2 in workspace
P27	MR has PE in the pallet
P28	MR has the Pallet
P29	Pallet in CT
P30	Pallet in ASRS

Tr	Table 4: Transitions definition
TO	MR moves to CT
T1	Move pallet from ASRS to CT
T2	CT transports pallet to assembly place
T3	MR pick pallet
T4	MR place pallet in workspace and open effector
T5	MR up 20 cm
T6	RM moves pallet to assembly zone
T7	RM moves out
T8	SV3D captures scene
T9	SV3D processes scene
T10	SV3D send to MR, position and orientation of Obj1
T11	MR pick Obj 1 from pallet that is in workspace
T12	MR moves obj 1 to assembly zone and return to 20 cm
T13	SV3D send to MR, position and orientation of Obj 2
T14	MR pick Obj 2
T15	MR moves obj 2 to assembly zone
T16	MR pushes for assembly
T17	MR places obj2 and return to 20 cm from workspace
T18	SV3D send to MR, position and orientation of Obj 3
T19	MR pick Obj 3
T20	MR moves obj 3 to assembly zone
T21	MR places obj 3 and return to 20 cm from workspace
T22	SV3D send to MR, position and orientation of Obj 4
T23	MR pick Obj 4
T24	MR moves obj 4 to assembly zone
T25	MR pushes for assembly and pick PE to pallet
T26	MR place PE and pick pallet
T27	MR moves pallet and place in CT
T28	CT transport pallet to ASRS
T29	Pallet is place in ASRS
T30	MR moves to initial position

Figure 9: Place and Transition definitions for PN Model

Definitions: MR - Robotics Manipulator; CT - Conveyor; SV3D - Stereo Vision Systems; Obj 2 - Internal Cover; Obj 4 - Plug 2;

ASRS: Automatic Storage and Retrieval System RM - Mobile Robot Obk 1 - Cover of the product Obj 3 - Plug 1 PE - Final assembly product

6 Evaluations of the Flexible Assembly Cell.

For the evaluations of the project different tests were made. One of them is to measure the time taken to make an assembly. Table 1 shows the times of joint, like for the evaluations, the iterations, with and without retiring pallet from the circuit of conveyor belts.

Table 1. Time to assembly a product

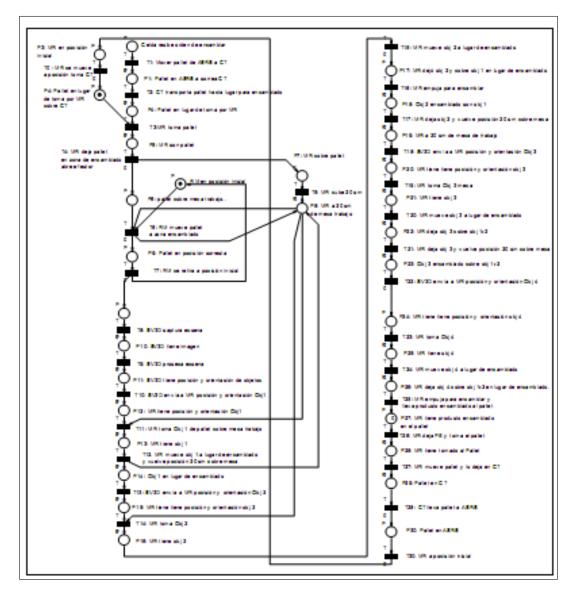


Figure 10: Flexible Assembly Cell PN Model

Geometric figures		Simple plug		Doble plug		Triple plug	
with	without	with	without	with	without	with	without
retiring	retiring	retiring	retiring	retiring	retiring	retiring	retiring
1'4"	1'37"	1'12"	1'34"	1'25"	1'59"	1'45"	2'22"

The time of assembled goes from 1 minute 12 seconds to 2'22", that is to say, the time when the manipulator puts the pallets in the work area, until the manipulator pick up the pallet with the assembled product.

A bad position of the pallets in the work area produces that the manipulator do not pick up the pallets or it pick up bad. One of the problem is the illumination affecting the quantification of the objects, like in the determination of the centroide and the direction of each one of the subparts. The Stero Vision System send to Manipulator wrong position. This problem it solves, in part, with the installation of a system of illumination superior to the existing one. A small mobile robots is used to move the pallet to the right place, using two walls to put the pallet in the desired position, the best place for stereo vision system (Fig.11). It can see the mobile robot, the walls and the pallet in the workspace. This tobots receives order from the computer cell.



Figure 11: Mobile Robot and the two walls.

Other evaluation is the measurement of the percentage of correct assembly. It is made ten measurement, for different assembled such as for the geometric figure, and each type of plug. Each one considers iterations with retiring pallet from the "circuit", and without retiring from them. The results given by Table 2, reflecting the percentage of error by amount of products, in other words, if the product to assemble has 3 pieces (geometric figure, simple plug) and fail in one, the result is 33% of error, or, if it has of 5 pieces, and it fails in two, have an error of the 40%, (triple plug).

	Geometric figures		Simple plug		Doble plug		Triple plug	
	with	without	with	without	with	without	with	without
	retiring	retiring	retiring	retiring	retiring	retiring	retiring	retiring
1	0%	0%	0%	0%	50%	0%	20%	20%
2	33%	0%	33%	0%	0%	25%	0%	20%
3	0%	33%	33%	0%	25%	0%	40%	0%
4	33%	0%	33%	0%	0%	25%	0%	0%
5	0%	0%	0%	0%	25%	0%	60%	0%
6	33%	0%	33%	33%	25%	0%	0%	40%
7	0%	0%	0%	33%	0%	0%	20%	40%
8	33%	33%	33%	0%	25%	50%	0%	60%
9	33%	0%	0%	33%	25%	50%	20%	20%
10	0%	0%	0%	0%	0%	0%	60%	20%

Table 2. Percentage of Faults to assembly a Product

7 Conclusions

In this paper has been presented a Petri Nets Model, the implementation and the of a Assembly Flexible Cell. The system assemblies several products, only is presented a wood assembly and assembly of different kind of electrical interrupt and electrical plugs. The Cell is composed by a robotic manipulator, a computer vision system and a conveyor. Different effectors for the robotics manipulator are designed for different kind of product to assembly, considering different factors of the object.

The Flexible Assembly Cell is modelled and simulated using Petri nets. This model and the simulation permit to evaluate the performance of the Cell, and correct it before starting a new assembly.

The real cell is tested assembled geometric objects and electrical interrupt and plugs. The time of assembled goes from 1 minute 12 seconds to 2'22", that is to say, the time when the manipulator puts the pallets in the work area, until the manipulator pick up the pallet with the assembled product.

The Stereo Vision System recognizes the objects whithin the working area, establishing the position of the pallets and the objects in the pallets. Problems affects the quatification determination of the centroid and direction of each object, producing error from 33% to 59%. The problem is solved with changing illumination and adding a small mobile robot to move the pallet to the desired position.

Bibliography

- [1] Lefranc, G., *La Manufactura Integrada por Computador: un Tutorial*, Magazine Automática e Innovación de la Asociación Chilena de Control Automático, vol 1, 2, 1993, pg 45.
- [2] T. C. Cao, A.C., Sanderson, 1992, Sensor-based error recovery for robotic task sequences using fuzzy petri-nets, Proceedings of the 1992 IEEE International Conference on Robotics and Automation, Vol.2, pp.1063-1069.
- [3] Q. Jing, W. Xisen, P. Zhihua, X. Youngcheng, 1996, A Research on Fault Diagnostic Expert System Basedon Fuzzy Petri Nets for FMS Machining Cell, Proceedings of IEEE International Conference on Industrial Technology, pp. 122-125.
- [4] Andrada I., Lefranc G. Sistema Experimental de un sistema de almacenamiento automatizado. I Parte. XV Congreso de la Asociación Chilena de Control Automático 2002.
- [5] Maizares M., Lefranc G., Base de datos y programa ejecutivo para un sistema AS/SR. IEEE Latin-American Conference on Robotics and Automation. Chile, 2003.
- [6] Lefranc, Gastón, Simulación de la generación de secuencias de ensamblado en una celda flexible de producción. Automática e Innovación de la Asociación Chilena de Control Automático. Ano 3. Volumen 2. N°7, 1996.
- [7] Lefranc G., Cano F., Sistema Servoing Pick and Place, Congreso Latinoamericano de Control Automático, 2003, Guadalajara, México.
- [8] Schleyer G., Lefranc G., *Experimental 3-D Visual servoing for FMS applications*. Third IFAC IEEE Conference on Management of Control of Production and Logistics MCPL'2004.
- [9] Opazo M., Lefranc G., Visual Servoing. Fourth IFAC IEEE Conference on Management of Control of Production and Logistics, MCPL 2007.
- [10] Schleyer G., Lefranc G., Tridimensional Visual Servoing, Studies In Informatics and Control, Volume 18. Issue 3. 2009.
- [11] Schleyer G., Lefranc G., Color Images Segmentation using Split&Merge and Region Growing Techniques in RGB and HSV Color Spaces, International Jornal of Computers, Communications & Control, Vol. V Issue 1, 2010.
- [12] Leighton F., Lefranc G., Flexible Assembly Cell using Scara Manipulator. Third IFAC IEEE Conference on Management of Control of Production and Logistics MCPL'2004.

Int. J. of Computers, Communications & Control, ISSN 1841-9836, E-ISSN 1841-9844 Vol. VI (2011), No. 2 (June), pp. 286-296

A Hybrid Artificial Bee Colony Algorithm for Flexible Job Shop Scheduling Problems

J. Li, Q. Pan, S. Xie, S. Wang

Jun-qing Li, Sheng-xian Xie School of Computer, Liaocheng University, Liaocheng Liaocheng, 252059, PR China E-mail: lijunqing.cn@gmail.com, xsx@lcu.edu.cn

Quan-ke Pan

 School of Computer, Liaocheng University, Liaocheng Liaocheng, 252059, PR China
 State Key Lab. of Digital Manufacturing Equipment and Technology, Huazhong University of Science and Technology
 Wuhan, 430074, PR China E-mail: qkpan@gmail.com

Song Wang

Department of Economic and Management Shandong University of Science and Technology Huangdao, 266510, PR China

> Abstract: In this paper, we propose a hybrid Pareto-based artificial bee colony (HABC) algorithm for solving the multi-objective flexible job shop scheduling problem. In the hybrid algorithm, each food sources is represented by two vectors, i.e., the machine assignment vector and the operation scheduling vector. The artificial bee is divided into three groups, namely, employed bees, onlookers, and scouts bees. Furthermore, an external Pareto archive set is introduced to record non-dominated solutions found so far. To balance the exploration and exploitation capability of the algorithm, the scout bees in the hybrid algorithm are divided into two parts. The scout bees in one part perform randomly search in the predefined region while each scout bee in another part randomly select one non-dominated solution from the Pareto archive set. Experimental results on the well-known benchmark instances and comparisons with other recently published algorithms show the efficiency and effectiveness of the proposed algorithm.

> **Keywords:** Flexible job shop scheduling problem, artificial bee colony, multiobjective optimization, hybrid algorithm

1 Introduction

The flexible job shop scheduling problem (FJSP), as a branch of the classical job shop scheduling problem (JSP), has been studied in very recent years. Brandimarte (1993) [1] is among the first author to solve the FJSP instances with tabu search (TS) algorithm. In very recent years, some meta-heuristic algorithms, such as TS algorithm [2] [3], particle swarm optimization (PSO) [4] [5], ant colony optimization (ACO) [6], and genetic algorithm (GA) [7] [8], have been used in solving the single-objective FJSPs. Although the single-objective FJSP has been widely investigated, the research on the multi-objective FJSP is still considered relative limited. Kacem et al. (2002a, 2002b) [9] [10] proposed an effective evolutionary algorithm. Xia and Wu (2005) [11] studied the problem with the hybrid algorithm of the PSO and the simulated annealing (SA). Zhang et al. (2009) [12] introduced a hybrid algorithm combining PSO algorithm with TS algorithm. Ho et al. (2008) [13]studied a hybrid evolution algorithm combined with a guided local search and an external Pareto archive set.

In this paper, we propose a hybrid algorithm combining an external Pareto archive set and the artificial bee colony (ABC) optimizer to solve the multi-objective FJSP. The rest of this paper is organized as follows: In Section 2, we briefly describe the problem formulation. Then, the artificial bee colony (ABC) algorithm is introduced in Section 3. The elements and framework of the hybrid algorithm are presented in Section 4 while Section 5 shows the experimental results and comparisons with other algorithms in the literature to demonstrate the superiority of the HABC performance. Finally, the last section presents conclusion of our work.

2 Problem Formulation

The FJSP considers n jobs to be processed on m machines. There are some assumptions and constrains in the FJSP considered in this study as follows: 1) each job has predefined number of operations and a known determined sequence among these operations; 2) each machine and each operation is ready at zero time; 3) each machine can only process one operation at a time, and each job must be processed on one machine at a given time; 4) each machine can process a new operation only after completing the predecessor operation; 5) each operation can be operated on a given candidate machine set instead of only one machine like in JSP; 6) given an operation O_{ij} and the selected machine M_k , the processing time p_{ijk} is also fixed.

Let C_i be the completion date of job J_i . W_k is the workload of machine M_k , which is the total processing time of operations that are operated on machine M_k . p_{ijk} be the processing time of $O_{i,j}$ on machine M_k . Three objectives are considered in this study, namely [13]:

1) minimization of maximum completion time (makespan):

$$F_1 = max\{C_i \mid i = 1, \dots, n\}$$
(1)

2) minimization of total workload

$$F_2 = \sum p_{i,j,k} \tag{2}$$

3) minimization of critical machine workload:

$$F_3 = max\{W_k \mid k = 1, \dots, m\}$$
(3)

3 Artificial Bee Colony Algorithm

Very recently, by simulating the behavior of honey bee swarm intelligence, an efficient bee colony (ABC) algorithm is proposed by Karaboga ([14] - [17]). Due to its simplicity and ease of implementation, the ABC algorithm has gained more and more attention and has been used to solve many practical engineering problems. In the basic ABC algorithm ([14] - [17]), there are two components: the foraging artificial bees and the food source. The position of a food source represents a possible solution to the optimization problem and the nectar amount of a food source corresponds to the quality or fitness of the associated solution. The artificial bee is divided into three groups, namely, employed bees, onlookers, and scouts bees. The employed bee is the one who is currently performing exploitation on a food source. A bee that is waiting in the hive for making decision to choose a food source is called an onlooker. The scout bee is

the one who perform exploration procedure and random exploitation search to find a new food source. The main steps of the algorithm are given as follows ([14] - [18]).

- **Step1**. Produce initial population;
- **Step2**. While stop criteria is not satisfied, perform steps 3 to steps 6.
- Step3. Send the employed bees onto their food sources.
- Step4. Send the onlooker bees onto the food sources depending on their nectar amounts.
- Step5. Send the scout bees to search possible new food sources.
- **Step6**. Memorize the best food source found so far.

4 The hybrid algorithm HABC

The basic ABC algorithm was originally designed for continuous function optimization. In order to make it applicable for solving the problem considered, a novel hybrid version of the ABC algorithm, named HABC, is proposed in this section.

4.1 Solution representation

The solution of the problem is represented with two vectors [19]: the machine assignment vector and the operation scheduling vector. The first part places the assigned machine number for each operation at the corresponding position, while the second part puts the same number symbol for each operation of a job and interpret them according to the occurrence in the operation scheduling vector.

4.2 Employed bee phase

The employed bee is to perform the local search around a given food source. Therefore, the employed bee takes the exploitation search of the algorithm. In order to generate good quality and diversity neighboring solutions, two types of local search operators are applied for the employed bees in this study, which are shown as follows.

(1)Local search operator in machine assignment component

The local search operator in machine assignment component is very simple and easy to be implemented. The perturbation is obtained by following steps.

Step1. Select a position in the machine assignment component, randomly or using some priority rules.

Step2. Assign a suitable machine different with the old one for the operation in the corresponding position.

Step3. Replace the machine number in the selected position and produce the new machine assignment component for the solution.

(2)Local search operator in operation scheduling component

The local search in the operation scheduling component is just like the perturbation in solving the JSP, where insert and swap operations are commonly used in the literature [20-22]. The insert operator is to remove a number symbol for an operation in the permutation π from its original position j and insert it into another position k such that $(k \neq j)$. The swap operator is to interchange two job symbols of π in the different positions.

After performing the above two local search approaches, the employed bee obtains a new neighboring food source around the old one. Then the new food source will be evaluated and compared with the old one. The better food source will be kept in the population as in the basic ABC algorithm which performs a greedy selection procedure.

4.3 Onlooker bee phase

In the classical ABC algorithm, each onlooker bee selects a food source based on the percent of the nectar amount of each food source among the total nectar amounts. However, the above approach consumes large computational time to compute the nectar amount of each food source. For this reason, we propose a tournament selection with the size of 3 in the HABC algorithm. In the tournament selection, three food sources are picked randomly from the population, and then the food source with highest nectar amount will be selected by the onlooker bee. After selecting the food source, each onlooker bee performs local search for the selected food source and produce a new neighboring food source. The better food source between the old one and the new neighboring one will be memorized in the population.

4.4 Scout bee phase

A scout bee performs randomly search in the basic ABC algorithm. This will increase the population diversity and avoid local minima, whereas this will also decrease the search efficacy. Since the food sources memorized in the Pareto archive set often carry better information than others, and the search space around these non-dominated solutions could be the most promising region. Therefore, in the HABC algorithm, the scout bees are first divided into two parts. One half of the scout bees randomly select a solution from the external Pareto archive set and perform several insert and swap operators to the selected solution, while the other half scout bees perform randomly search in the predefined search scope. In the hybrid algorithm, at least 5% - 10% of the population is scout bees.

4.5 Multi-objective optimizer

The Pareto archive set AS

To provide a set of solutions with good diversity, a Pareto archive set (AS) was introduced in this study, which is used to maintain a limited number of non-dominated solutions found so far. During the optimization process, the archive set is iteratively updated with adding some nondominated solutions and removing some dominated solutions to get closer to the Pareto-optimal front. Once a new non-dominated solution is found, it will be added to AS and any solution which is dominated by the added one will be removed from AS. In case AS becomes overfull, its member which is in the crowded domain is eliminated to maintain the diversity of the Pareto archive set.

The storage structure of AS

To reduce the computational time complexity consumed on the update process of the archive set, the members of the AS firstly sequence in an ascending order according to their first objective function value (Pan, 2009) [21].

Non-dominated sorting algorithm

For the population, we should sequence each solution according to a certain criteria. For multi-objective optimization problems, we can not use one objective function value to determine the solution quality. In this study, a non-dominated sort algorithm (Deb et al., 2002) [23]was

introduced to divide the population solutions into several levels according to their dominated solutions number.

4.6 The framework of HABC

The details steps of the proposed HABC algorithm are as follows:

Step1 Initialization phase;

Step 1.1 Set the system parameters;

Step 1.2 Produce the initial population.

- Step2 Apply the Pareto non-dominated sorting function on the population, and then update the external Pareto archive set by using the solutions in the first Pareto level front.
- **Step3** If the stopping criterion is satisfied, output the non-dominates solutions in the external Pareto archive set; otherwise, perform steps 4-7.
- Step4 Employed bee phase.

Step 4.1 Put each employed bee on each solution in the population.

- Step 4.2 For each employed bee, perform local search on the assigned solution and generate a new neighboring solution.
- Step 4.3 Evaluate the new neighboring solution and record the better solution among the new solution and the old one as the current solution and put it into the population. If the two solutions are non-dominated with each other, randomly select one as the current solution.
- **Step 4.4** If a solution has not been improved through limit cycles, then the corresponding employed bee becomes a scout bee and perform step 6.
- **Step 4.5** Evaluate each solution corresponding to each employed bee, apply the Pareto non-dominated sorting algorithm on the new population and update the external Pareto archive set using the solutions in the first Pareto level.
- Step5 Onlooker bee phase.
 - **Step 5.1** For each onlooker bee, randomly selects three solutions from the population and selects the best one as the food source. If the three solutions cannot dominate each other, then randomly select a non-dominate solution.
 - **Step 5.2** For each onlooker bee, performs local search for the selected food source and carries out greedy selection procedure to record the better solution in the population.
 - **Step 5.3** Evaluate each solution corresponding to each onlooker bee, apply the Pareto non-dominated sorting algorithm on the new population and update the external Pareto archive set using the solutions in the first Pareto level.
- Step6 Scout bee phase.
 - Step 6.1 Divide the scout bees into two parts with the same number of bees.
 - Step 6.2 The scout bees in the first part randomly select a food source and perform local search operator in the predefined region. After generating a new solution, performs greedy selection procedure.
 - Step 6.3 Each scout bee in the second part randomly select a non-dominate solution in the external Pareto archive set and perform several local search for the selected solution.. After generating a new solution, performs greedy selection procedure.

Step 6.4 Evaluate each solution corresponding to each scout bee, apply the Pareto nondominated sorting algorithm on the new population and update the external Pareto archive set using the solutions in the first Pareto level.

Step7 go to step 3.

5 Experiment Results

This section describes the computational experiments to evaluate the performance of the proposed algorithm. The test samples come from Kacem instances set [9]. The current instantiation was implemented in C++ on a Pentium IV 1.8GHz with 512M memory.

5.1 Setting parameters

Each instance can be characterized by the following parameters: number of jobs (n), number of machines (m), and the number of operations (op_num) . Followings are the detail parameters value:

The size of the population is equal to the number of employed bee and the number of onlooker, which is set to 5n; the maximum cycle of the algorithm is set to $10 \times n \times m$; the limit number of cycles through which no improvement occurs on the food source, then the employed bee becomes a scout bee; the limit number is set to $n \times \frac{m}{2}$; the percent of scout bee is set to a random number between 0.05 and 0.1.

5.2 Results comparisons

The five test instances come from Kacem [9] [10], which range from 4 jobs×5 machines to 15 jobs ×10 machines. Two tests are performed for comparison, i.e. the instances with single objective and the problems with three objectives. Several recently published algorithms are compared with the proposed HABC algorithm, such as the AL+CGA proposed by Kacem et al. (2002b) [10], the GENACE approach employed by Ho (2004) [24], the PSO+SA developed by Xia and Wu (2005) [11], the ant systems & local search optimization method (hereafter called ACO+LS) presented by Liouane et al. (2006) [6], and the PSO+TS introduced by Zhang (2009) [12].

The signle objective five Kacem instances

For solving the five instances with single objective to minimize the makespan criterion, the experimental results and comparisons are given in Table 1. It can be seen from Table 1 that the HABC algorithm can obtain the best results for all the Kacem instances. The proposed algorithm outperforms the AL+CGA in 4 out of 5 instances, while outperforms the GENACE method in 2 out of 4 cases. For comparison with the very recently published algorithms, the HABC algorithm obtained a superior result in solving the largest problem than the ACO+LS proposed by Liouane (2006). Particle swarm optimization (PSO) is an efficient swarm intelligent algorithm and the experimental results obtained by PSO+SA and PSO+TS are considered as the competitive results for the FJSP. Table 1 shows that the HABC algorithm outperforms the SO+SA algorithm in 2 out of 3 problems. The HABC algorithm can obtain the same experimental results with the PSO+TS in very short computational times. For example, in solving the largest problem 15×10 , our algorithm consumes just about 50 seconds to reach the best result so far.

Problem Set	AL+CGA	GENACE	ACO+LS	PSO+SA	PSO+TS	HABC
4×5	16	11	11	-	-	11
8×8	15	-	-	15	14	14
10×7	15	12	11	-	-	11
10×10	7	7	7	7	7	7
15×10	23	12	12	12	11	11

Table 1: Comparison of the five instances with single objective (makespan)

"-" denotes not given by the author

The multi objective five instances

Table 2 shows the comparison of the results on the five multi-objective FJSP instances. The three objectives are considered simultaneously, i.e. minimization of the makespan (denoted by f_1), the total workload (denoted by f_2), and the maximal workload (denoted by f_3). It can bee seen from Table 2 that the HABC algorithm is competitive to other algorithms. The experimental results of the proposed algorithm dominate the results of the AL+CGA for solving the four instances. For comparison with the very recently published algorithms, the HABC can either obtain more non-dominated solutions or obtain superior result than PSO+TS and PSO+SA algorithms. For example, our algorithm obtain three non-dominated solutions in solving the 8×8 instance, while the PSO+SA and PSO+TS can only obtain two results. In addition, our algorithm obtains all these results in a run while the other algorithms can obtain different results for an instance. Therefore, the external Pareto archive can enhance the population diversity of our algorithm. Fig. 1 shows the Gantt chart of the resulted solution of 15×10 instance, due to our proposed algorithm.

To make a further comparison with the ACO+LS proposed by Liouane (2006), we also test the problem listed in the paper [6]. The problem is given in Table 3. The comparison of the experimental results from ACO+LS and our algorithm are given in Table 4 and the two solutions obtained by the HABC algorithm are given in Table 5 and Table 6, respectively. It can be seen from Table 4 that the HABC algorithm obtained two non-dominated solutions for the example benchmark while the TS method in literature [6]can obtained only one solution. Furthermore, the resulted solutions obtained by our algorithm dominated all the results by the ACO method. In addition, our algorithm obtained the two non-dominated solutions consuming just only 0.01 seconds for the example benchmark. Therefore, Table 4 concludes that the proposed algorithm is efficient in solving the example problem especially when compared with the ACO method.

6 Conclusion

In this paper, we have proposed an efficient algorithm for solving multi-objective FJSPs. Instead of applying the basic ABC algorithm, we developed a hybrid ABC method. To memory the non-dominated solutions found so far and increase the population diversity, we presented an external Pareto archive set. A fast Pareto update function is also introduced in the algorithm to enhance the computational capability. In the hybrid algorithm, the balance of the capability of exploration and exploitation is considered. Experimental results on several well-known benchmarks show that our algorithm is competitive to other recently published algorithms for solving the FJSPs. The future work is to improve the neighborhood structure of the problem considered

		A	L+CGA	Р	SO+SA	Р	SO+TS		HABO	C
	f_1	16				11		11	12	13
4×5	f_2					32		32	32	33
1/(0	f_3	10				10		10	8	7
	f_1	15	16	15	16	14	15	14	15	16
8×8	f_2		75	75	73	77	75	77	75	73
	f_3	13	13	12	13	12	12	12	12	13
	f_1							11	12	
10×7	f_2							61	60	
	f_3							11	12	
	f_1	7		7		7		8	7	8
10×10	f_2			44		43		41	42	43
	f_3	5		6		6		7	6	5
	f_1		24	12		11		12	11	
15×10	f_2		91	91		93		91	93	
	f_3	11	11	11		11		11	11	

Table 2: Comparison of the five instances with three objectives

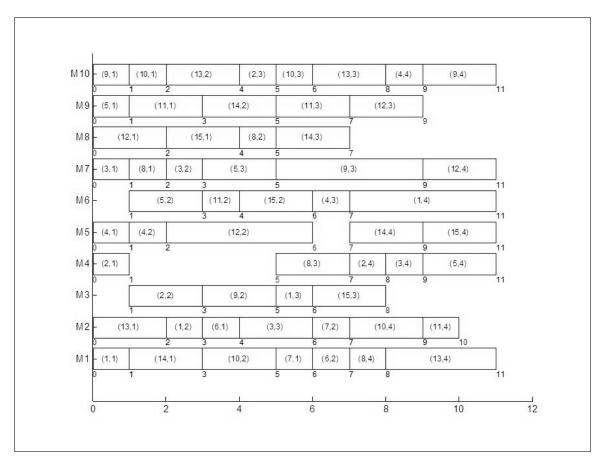


Figure 1: The Gantt chart of the solution 2 (f_1 =11, f_2 =93, f_3 =11) for the instance 15 × 10

		M_1	M_2	M_3	M_4	M_5	M_6
	O_{11}	10	7	6	13	5	1
T	O_{12}	4	5	8	12	7	11
J_1	O_{13}	9	5	6	12	6	17
	O_{14}	7	8	4	10	15	3
т	O_{21}	15	12	8	6	10	19
J_2	O_{22}	9	5	7	13	14	7
	O_{23}	14	13	14	20	8	17
т	O_{31}	7	16	5	11	17	9
J_2	O_{32}	9	16	8	11	6	3
	<i>O</i> ₃₃	6	14	8	18	21	14

Table 3: Example benchmark 3 jobs-6 machines from [6]

Table 4: Comparison of the example benchmark

	f_1	f_2	f_3	
Lower Bound Value	18	45	8	Avgtime(s)
	19	51	13	
	19	50	13	
ACO+LS	19	48	14	
	19	47	14	
TS	18	45	12	-
HABC	18	45	12	0.01
	19	46	10	

"-" denotes not given by the author

Table 5: Solution 1 for the example benchmark (f_1 =18, f_2 =45, f_3 =12)

	O_1	O_2	O_3	O_4
J_1	$M_6:[0,1]$	$M_1:[1,5]$	$M_2:[5,10]$	$M_6:[10,13]$
J_2	$M_4:[0,6]$	$M_5:[6,10]$	$M_5:[10,18]$	***
J_3	$M_3:[0,5]$	$M_6:[5,8]$	$M_1:[8,14]$	***

Table 6: Solution 2 for the example benchmark $(f_1=19, f_2=46, f_3=10)$

	O_1	O_2	O_3	O_4
J_1	$M_6:[0,1]$	$M_1:[1,5]$	$M_2:[11,16]$	$M_6:[16,19]$
J_2	$M_4:[0,6]$	$M_5:[6,11]$	$M_5:[11,19]$	***
J_3	$M_3:[0,5]$	$M_6:[5,8]$	$M_1:[8,14]$	***

and enhance the convergence ability of the algorithm.

7 Acknowledgments

This work supported by National Science Foundation of China under Grants 60874075, Science Research and Development of Provincial Department of Public Education of Shandong under Grant J08LJ20, J09LG29, J08LJ59, Soft Science Foundation of Shandong under Grant 2009RKB125.

Bibliography

- P. Brandimarte. Routing and scheduling in a flexible job shop by tabu search. Annals of Operations Research, Vol 22, pp. 158-183, 1993.
- [2] M. Mastrolilli and L. M. Gambardella. Effective neighborhood functions for the flexible job shop problem. *Journal of Scheduling*, Vol. 3(1), pp. 3-20, 2000.
- [3] J. Q. Li, Q. K. Pan, P. N. Suganthan and T. J. Chua. A hybrid tabu search algorithm with an efficient neighborhood structure for the flexible job shop scheduling problem. *International Journal of Advanced Manufacturing Technology*. doi: 10.1007/s00170-010-2743-y.
- [4] L. Gao, C. Y. Peng, C. Zhou and P. G. Li. Solving flexible job shop scheduling problem using general particle swarm optimization. In Proceedings of the 36th CIE Conference on Computers & Industrial Engineering, pp. 3018-3027, 2006.
- [5] J. Q. Li, Q. K. Pan and S. X. Xie. A hybrid variable neighborhood search algorithm for solving multi-objective flexible job shop problems. *ComSIS Computer Science of Software Information and System.* doi: 10.2298/CSIS090608017L.
- [6] N. Liouane, I. Saad, S. Hammadi and P. Borne. Ant Systems & Local Search Optimization for Flexible Job-Shop Scheduling Production. *International Journal of Computers, Commu*nications & Control, Vol. 2, pp. 174-184, 2007.
- [7] J. Gao, L. Sun and M. Gen. A hybrid genetic and variable neighborhood descent algorithm for flexible job shop scheduling problems. *Computers & Operations Research*, Vol. 35(9), pp. 2892-2907, 2008.
- [8] F. Pezzella, G. Morganti and G. Ciaschetti. A genetic algorithm for the Flexible Job-shop Scheduling Problem. Computers & Operations Research, Vol. 35, pp. 3202-3212, 2008.
- [9] I. Kacem, S. Hammadi and P. Borne. Pareto-optimality approach for flexible job-shop scheduling problems: hybridization of evolutionary algorithms and fuzzy logic. *Mathematics and Computers in Simulation*, Vol. 60, pp. 245-276, 2002a.
- [10] I. Kacem, S. Hammadi and P. Borne. Approach by localization and multi-objective evolutionary optimization for flexible job-shop scheduling problems, *IEEE Transactions on Systems* , Man and Cybernetics, Part C, Vol. 32(1), pp. 408-419, 2002b.
- [11] W. J. Xia and Z. M. Wu. An effective hybrid optimization approach for multi-objective flexible job-shop scheduling problems. *Computers & Industrial Engineering*, Vol. 48(2), pp. 409-425, 2005.

- [12] G. H. Zhang, X. Y. Shao, P. G. Li and L. Gao. An effective hybrid swarm optimization algorithm for multi-objective flexible job-shop scheduling problem. *Computers & Industrial Engineering*, Vol. 56(4), pp. 1309-1318, 2009.
- [13] N. B. Ho and J. C. Tay. Solving Multiple-Objective Flexible Job Shop Problems by Evolution and Local Search. IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews, Vol. 38(5), pp. 674-685, 2008.
- [14] D. Karaboga. An idea based on honey bee swarm for numerical optimization. Technical Report TR06. Computer Engineering Department. Erciyes University. Turkey. 2005.
- [15] D. Karaboga and B. Basturk. A powerful and efficient algorithm for numerical function optimization: Artificial Bee Colony (ABC) algorithm. *Journal of Global Optimization*, Vol. 39(3), pp. 459-171, 2007.
- [16] D. Karaboga and B. Basturk. On The performance of Artificial Bee Colony (ABC) algorithm. Applied Soft Computing, Vol. 8(1), pp. 687-697, 2008.
- [17] D. Karaboga and B. Akay. A comparative study of Artificial Bee Colony Algorithm. Applied Mathematics and Computation, Vol. 214, pp. 108-132, 2009.
- [18] Q. K. Pan, M. F. Tasgetiren, P. N. Suganthan and T. J. Chua. A discrete artificial bee colony algorithm for the lot-streaming flow shop scheduling problem. Information Sciences. doi: 10.1016/j.ins.2009.12.025, 2010.
- [19] J. Q. Li, Q. K. Pan and Y.C. Liang. An effective hybrid tabu search algorithm for multiobjective flexible job shop scheduling problems. *Computers & Industrial Engineering*, Vol. 59, pp. 647-662, 2010b.
- [20] F. Pezzella, G. Morganti and G. Ciaschetti. A genetic algorithm for the Flexible Job-shop Scheduling Problem. Computers & Operations Research, Vol. 35, pp. 3202-3212, 2008.
- [21] Q. K. Pan, L. Wang, B. Qian. A novel differential evolution algorithm for bi-criteria no-wait flow shop scheduling problems. *Computers & Operations Research*, Vol. 36(8), pp. 2498-2511, 2009.
- [22] L. Wang. Shop scheduling with genetic algorithms. *Tsinghua university press*, Beijing, China, 2003.
- [23] K. Deb, A. Paratap, S. Agarwal and T. Meyarivan. A fast and elitist multiobjective genetic algorithm: NSGA-II. *IEEE Transactions on Evolutinary Computation*, Vol. 6(2), pp. 182-197, 2002.
- [24] N. B. Ho and J. C. Tay. GENACE: An Efficient Cultural Algorithm for solving the Flexible Job-Shop Problem. In Proceedings of the IEEE Congress on Evolutionary Computation (CEC2004) Piscataway, pp.1759-1766, 2004.

The Communication in Distributed Client - Server Systems Used for Management of Flexible Manufacturing Systems

V. Lupu, D.E. Tiliute

Valeriu Lupu, Doru E. Tiliute "Stefan cel Mare" University of Suceava Department of Informatics, Faculty of Economics and Public Administration E-mail: valeriul@seap.usv.ro, dtiliute@seap.usv.ro

Abstract: Labour productivity growth is a necessary condition for social and economic progress, in general, and to overcome the economic crisis facing most of the world, in special.

Applying innovative solutions, based on the ITC, is one of the straight ways for achieving that objective, both important and necessary. This paper presents a software solution applicable to industrial production based on numerically controlled machines. It involves a distributed client - server communication system, combined with MLP neural networks for the recognition of the 2D industrial objects, viewed from any angle. The information on prismatic and rotational parts to be processed by numerically controlled machine, are stored on a database server together with the corresponding processing programs. The client applications run on the numerically controlled machines and on the robots serving groups of machines. While the machines are fixed, the robots are mobile and can move from a machine to another. As a novelty of the proposed solution, in some well defined situations, the clients are allowed to change messages among them, in order to avoid the server overload. The neural networks are used to help robots to recognize the parts before and during manipulation.

Keywords: client, server, web, Java, image.

1 Introduction in domain

The flexible production system consists in a number of fixed numerically controlled machines which process the parts, several mobile robots [1] whose main task is to load/unload the machines and storing parts in technologic stocks. The production system is driven by a complex distributed informatics system containing two servers, an application server and a database server, and several clients. The application server provides those applications required by the coordination of the robots and controlling the processing [2] stage of each part on each machine. It also has the communication protocols with the clients and with database server. Two video cams equip each robot and machine in order to ensure a 3D vision of parts to be processed or handled. The position of each robot is monitored continuously by the application server using an appropriate number of video cams fixed on the production hall ceiling.

An overview of the communication system is shown in figure 1.

For each part, the following geometrical features are stored and available in database server: aria, perimeter, scale factor, minimum radius, maximum radius, average radius, standard radius deflection, minimum embedded rectangle dimensions, maximum inertia momentum, minimum inertia momentum, elongation factor, average of the gray levels, Pentland ratio, $R_{min}(x)/R_{max}(x)$ ratio, $A_{elipsa}(x)/B_{elipsa}(x)$ ratio, form characterization using momentum, etc.). Functional features achieving an accurate enough description of the image objects using analytically developments easy to utilize. This category comprise Fourier coefficients of the inherent contour curve function, inertia momentum of different orders calculated on contour curve and the invariant computed using inertia momentum.

Topological features - conveying objects proprieties non depending on the distortions who aren't affecting the surface. These features are free from distances. We specify some of them:

C - the connected components number;

- H the cavities number;
- E = C-H Euler's number.

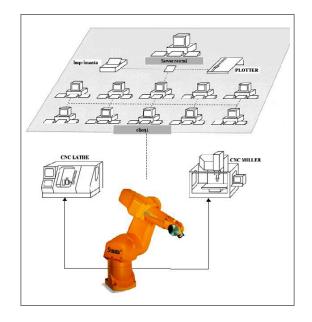


Figure 1: An example of distributed client-server system

Also, the relations among different areas are forming a good descriptor. These relations can be represented as a graph of connected components.

The features mentioned above are needed in order to recognize the parts about to be processed using numerically controlled machine. Images of the parts are captured by two video cameras placed on each robot in perpendicular planes, resulting in a 3D view [3]. For recognition it is used a neuronal architecture of multilayer perceptron type.

The solution presented in this paper employs neural networks for image processing and for the recognition of the model-based industrial objects [5]. In this order, a software application was created, giving to the robots the ability of recognizing the parts to be processed on the numericalcontrol machines. The software provides also a human interface that was used to verify correct functionality of neural network pattern recognition and communications, as described below.

The neural network [4] training is ensured by a software component that implements the wellknown backpropagation algorithm. For the classification of the vectors bearing the characteristics of the shapes in the input space, the model of the multilayer perceptron was used. The vectors of the training set present themselves to the network with the class they belong to. In other words, the training inputs are organized as association lists in which the associated vectors correspond to a classification of the key vectors. Consequently, an element of the list is a pair key vector - corresponding class. If sufficient objects have been memorized from each class, an internal representation can be made of each class by means of the connection weights of the network. The model of the multilayer perceptron used in this software solves the problem of the shape classification. The network inputs are the vectors that are to be classified, while its outputs are the corresponding classes.

2 The client-side program

The "client" program contains three classes: "client", "network" and "writer". Initially, the program imports three packages: "java.io.*" (used for input/output actions; io = Input / Output), "java.awt.*" (used for graphical objects; awt = Abstract Window Toolkit) and "sun.net.*" (used to establish a connection between client and server) [6], [7].

In order to run the client program, the user will type the following command: java client $\langle IPoraddress \rangle \langle ClientName \rangle$;

where:

 \rightarrow *java* - the name of the application that launch the classes;

 \rightarrow *client* - the name of the class that will be executed;

 \rightarrow *IP or address* - IP address or server name;

 \rightarrow ClientName - the name of the client who wants to connect to the server.

The following sections explain the three classes of the client application.

The "client" class.

The "client" class extends the Frame class, which has a method called add. Using this method we can insert in the middle of the frame an object of type network.

The class contains two methods: main and client.

In the main method it is verified if the number of arguments is smaller than two. The arguments meaning is:

- args[0] - the IP or the server name;

- args[1] - the name of the client who wants to connect to the server.

This function is the first executed when the client class starts.

With this method, an instance of the client class is created, having as parameters the server IP address and the client's name.

The client method receives as parameters the name or the IP address of the server (String host) and the user name (String username). Inside the frame we insert a new object, which is actually a new instance of the network class, having as parameters the same values received by the function. Then, the frame size is set to (500,500). The method show() displays the frame on the main screen.

The "network" class.

This class extends the Panel class. The meaning of the variables is the following:

• public static boolean afis - if it is true, then the messages that are received from server are shown in the green box from the client's main window;

• NetworkClient network_client - using this method, we can declare a client. We use the specifications from the sun.net.* package;

• DataInputStream net_input - using this variable, the client can receive messages from server;

• PrintStream net_output - using this variable, the client can send messages to the server;

• static int modif=0 - if it is 1, then the image received from server is shown in the client's main window. This image is called "imagine.gif". If it is 0, then this image isn't shown in the window.

• String username - user name;

• boolean connected=false - if it is true, this means that the client is connected to the server;

• writer w - using this variable, we can create a new instance of the writer class, which handles with the read of characters or strings received from server;

• String typed_line="" - the message that's typed by client.

Bellow are the methods of this class:

The class constructor. The constructor receives two parameters, which represents the name or the IP address of the server and the user name. Each time a client tries to connect to the server, the connect method is called. If the connection succeeded, the connected will be set on true and also a new instance of the writer class will be created. If the connection doesn't succeed, an error message is returned to the client (it may also be viewed on the human interface client) and the execution of the application stops.

The **connect** method. This method tries to connect to the server on a specific port (we chose port No. 1111 but it can be any unassociated port). The connection is achieved by creating a new instance of the class NetworkClient, which can be found in the sun.net.* package. The arguments of this constructor are the name or IP address of the server and the port number. After establishing the connection the net_input and net_output variables are initialized with instances of the classes that correspond to their definition. The serverOutput variable is declared inside the NetworkClient class and it is initialized when the constructor of the NetworkClient class is called.

The **read_net_input** method. This method reads a line from server and returns the string received or null in case that an input/output error occurs. It also returns null if there is nothing available from the server.

The write **net** output method. This method writes a string to the server.

The **close_server** method. This method tries to close the communication between client and server.

The **keyDown** method. This method is executed when a key is pressed in the client's window (in human interface client version). Also, it calls the repaint method in order to make a quick repaint of the client's window.

The **paint** method. This method paints the client's frame. If writing messages from server is allowed, a welcome message is displayed: "Hello, <UserName> !", then two colored boxes are drawn, one red and the other one green, under the welcome message. Each box displays different messages; the red one displays the messages sent by client and in the green one are displayed the messages received from server. If the modif variable is set, then the image "imagine.gif" will be shown. This image is received from the server.

The "writer" class.

This class handles with the reading of data from server. This method extends the Thread class, which is specific for repeated actions. If the keyword "RECV" is received from server, the variable primesc becomes true. The string specifies that a data transaction between server and client is about to begin. The transmission consists in the content of a ".gif" image, whose size is specified by "SIZE nnnn" string, which precedes the actual content. Data are received until the "STOP" string is detected. In this case the modif variable from network class will become 1 and the repaint method from the same class will be called.

3 The server-side program

The server application contains two classes [6], [7].

The server class extends the NetworkServer class, which is a part of "sun.net" package. In this class we consider that the maximum number of users is 1,000. The meaning of the variables from the server class is given bellow:

String user[] - this array will contains the name of the users that are currently connected to server;

String sir[] - in this array will be stored the last message from each user;

DataInputStream net_input[] - by using this array, the server can receive data from clients; PrintStream net_output[] - by using this array the server can send information and data to cliens;

reader r[] - in this vector will be stored the addresses of the "reader" class instances; static int client_counter = 0 - the current number of connected clients.

The main method creates a new instance of the server class.

The constructor of this class initializes the arrays presented above.

When idle, the server console displays Waiting for users...

The most important method in this class is serviceRequest, which is a function of the NetworkServer class. When an event occurs at the server, this function is automatically called and creates an instance of the reader class. This thing happens each time when a client connects to the server.

The method read_net_input reads data from clients.

The method write_net_output writes data to clients.

Then the server is started, by calling the "startServer(1111)" instruction. This instruction starts the server, specifying that the port on which the server will receive message is 1111. This port number can be changed, but if is changed in the server class it also has to be changed in client class.

If the server was started successfully, a message is printed on the server console:

Waiting for users ...

If it's not possible to open the server on the port 1111, then an error message will be shown and the execution of this application will be canceled.

The most important method in this class is serviceRequest, which is a function of the NetworkServer class. When an event occurs at the server, this function is automatically called. In this function it is created an instance of the reader class. This thing happens each time when a client is connecting to the server.

The method read_net_input reads data from clients.

The method write net output writes data to clients.

The *reader* class

This class extends the Thread class. The most important method of this class is run. This method is called each time when an event occurs (connection of a client, the send of a message from a client, etc.). When the end of sentence is detected (point, !, ?), it checks if the message represents the standard syntax for sending data (images). This syntax is:

SEND < ImageName > TO < user > where ImageName is the name of the image stored on the server and the user is the name of the user that has to receive the image. Examples:

SEND image1 TO machine1.

Send image2 TO machine2.

The syntax is "case-insensitive".

The application server has the possibility to search the image in database. Using the methods of the package "java.sql", queries may be sent to database and the results are fetched. Snapshots of server and clients with human interface are depicted in figure 2. Samples of pictures of parts to be processed are given in figure 3.

4 Case study - MLP type neural network for recognizing objects from any angle

The network has three layers. At the input, each neuron is connected to the outputs of all the neurons from the previous layer, its output being connected to the inputs of all the neurons of the following layer. Within the same layer, the neurons are not connected between them. The dimension of the input layer is given by the dimension of the five vectors describing each shape (sphericity, Pentland's ratio, radii ratio, ratio between the ellipses on the two axes, shape characterisation by means of moments). The dimension of the output layer is equal to the number of classes into which the input data will be separated, that is four (corresponding to the classes: disk, square, triangle, cross).

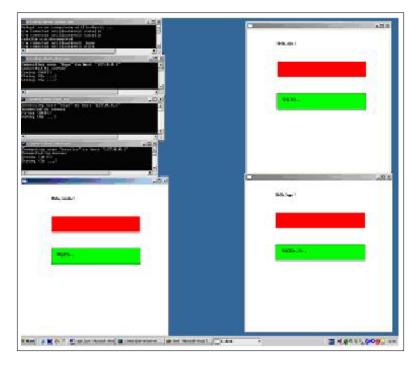


Figure 2: a) windows screen

The activation function used for the output of the neurons is:

$$f(a) = 1/(1 + exp(-a)).$$

As a result, the outputs of the network (y) will be in the domain $[0, 1]^4$. The representative of the class where the data set presented at the input is placed will have the output 1, while the others, 0.

For other types of pieces, the theoretical values cannot be determined directly, requiring measurements and experiments specific to each class.

The values obtained, following the implementation of the parameters calculation are given in table 2 (through medium values):

Table 2 Medium Values of the Shapes' Parameters

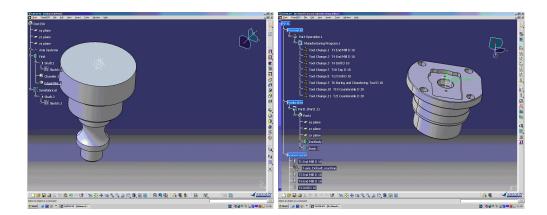


Figure 3: b) picture 1 c) picture 2

	C(X)	P(X)	$R_{min}({ m X})/$	$A_{eclipse}({ m X})/$	F(X)
			$R_{max}(\mathbf{X})$	$B_{eclipse}(\mathbf{X})$	
••.	1.014	1.120	0.921	1.042	0.040
-	0.841	0.531	0.673	2.512	0.125
	0.691	0.512	0.481	3.824	0.230
× , ,					
+×	0.265	0.401	0.213	14.412	0.612

Where: 1. sphericity $C(X) = 4 \pi A(X) / (P(X))^2$.

- **2.** Pentland's ratio $P(X) = 4 A(X) / (\pi FMD(X))^2$.
- 3. radii ratio $R_{max}(X) = maxdi / i = 1..N;$
- $R_{min}(X) = \min di / i = 1..N.$

4. ratio between the ellipses on the two axes.

5. shape characterisation by means of moments

$$F(X) = \frac{\sqrt{1/N \sum_{1}^{N} (d_i - M_1)^2}}{1/N \sum_{1}^{N} d_i}$$

Structure of the file data.txt

N patterns - number of input vectors for the network training

 $X_1, X_2, X_3, X_4, \ldots, X_{np}Y_i$ (i=1, nc, where nc = number of classes) the components of the input vector (shapes' parameters), where np = number of shapes' parameters.

The parameters were introduced as follows (according to the table corresponding to each shape):

class			
1-square	$0.867 \ 0.660 \ 2.145 \ 0.124$	$0.583 \ 0.583 \ 2.302 \ 0.164$	$0.738 \ 0.738 \ 1.968 \ 0.102$
2-disk	$1.018 \ 0.867 \ 1.000 \ 1.000$	$1.021 \ 0.857 \ 1.000 \ 0.039$	$0.994 \ 0.867 \ 1.168 \ 0.051$
3-cross	$0.160 \ 0.063 \ 22.924 \ 0.732$	$0.292 \ 0.250 \ 11.615 \ 0.462$	$0.169\ 0.143\ 21.596\ 0.647$
4-triangle	$0.670\ 0.462\ 3.696\ 0.229$	0.712 0.423 3.313 0.231	$0.708 \ 0.514 \ 3.347 \ 0.211$

5 Conclusions

The functionality of the proposed system was proved by a dedicated application software that responds to the following main requirements:

- provides the recognition of the 2D and 3D industrial objects viewed from any angle;
- provides border touching, partial overlapping and clipping of the objects in the image;
- it is noise tolerant.

The comparative study of the existent solutions, in correlation with the above-mentioned requirements, leads to the conclusion that 3D object-centered models are suitable in order to ensure the deduction of the model views from any angle.

Bibliography

- C.V. Kifor, C. Oprean, D.D.M. Banciu Intelligent Systems for Assisting Decisions in Advanced Product and Process Planning and Design, *Studies in Informatics and Control*, SIC Vol.18, No.3, 2009
- Mircea Ivanescu, Mihaela Cecilia Florescu, Nirvana Popescu The control of the Hyper - redundant Manipulators by Frequency Criteria, Studies in Informatics and Control, SIC Vol.18, No.3, 2009
- [3] Lupu Valeriu Contributions to the management of the flexible manufacturing systems (PhD Thesis), 2004
- [4] Mohammad Reza Soltanpour, Seyed Ehsan Shafie Design and Stability Analysis of a Robust Impendance Control System for a Robot Manipulator, *Studies in Informatics and Control*, SIC Vol.19, No.1, 2010
- [5] O.Khatib, K.Yokoi, K.Chang, D.Ruspini, R.Holmberg, A.Casal Coordination and Descentralized Cooperation of Multiple Mobile Manipulators, *Journal of Robotic Systems*, Volume 13, No.11, 1996, ISSN: 0741 - 2223
- [6] D.Logofatu Algoritmi fundamentali in Java. Aplicatii, ISBN: 978 973 46 0815 7, Editura Polirom, 2007
- [7] Java on-line documentation (http://java.sun.com)

Fuzzy Based Packet Dropping Scheme in Wireless Cellular Networks

J.D. Mallapur, S.S. Manvi, D.H. Rao

Jayashree D. Mallapur

Department of Electronics and Communication Engineering Basaveshwar Engineering College Bagalkot, India E-mail: bdmallapur@yahoo.co.in

Sunilkumar S. Manvi

Department of Electronics and Communication Engineering Reva Institute of Technology and Management Bangalore, India E-mail: sunil.manvi@revainstitution.org

D.H. Rao

Jain College of Engineering Belgaum, India E-mail: dr.raodh@gmail.com

> **Abstract:** Wireless multimedia networks are becoming very popular owing to the user demands for multimedia services. Packet dropping in the event of buffer congestion is one of the important issue in wireless multimedia networks. A packet dropping scheme has to be flexible and adaptive such that acceptable quality of an application is maintained. The paper presents a fuzzy based packet dropping scheme for wireless multimedia networks. A buffer manager placed at the base station performs packet dropping depending upon the traffic conditions and type of an application. Packet dropping is performed by computing dropping factor by considering packet priority, queue length and adaptive queue length threshold. The adaptive queue length threshold is used to dynamically adjust the dropping factor. The queue length threshold is varied by using two fuzzy input parameters, channel condition and rate of flow of an application. The scheme has been extensively simulated to test the performance in terms of acceptance and dropping probability of real-time handoff and new calls.

Keywords: Buffer, Multimedia, Fuzzy Logic, Packet Drop.

1 Introduction

The use of wireless connectivity has been steadily growing over the recent years. New generation of mobile networks are opting for multimedia applications, such as video conferencing, e-commerce, e-learning, e-gaming, etc. These applications need good end-to-end quality of service (QoS) control even under the conditions of network overload. The QoS parameters are bandwidth, buffers, delays, jitters, packet loss, etc. The bandwidth and buffer management have become critical issues because of their scarcity. Whenever sufficient bandwidth is not available, buffers are employed to avoid packet losses. Due to excess traffic, buffers may become full and cause buffer congestion. The buffer congestion leads to dropping of packets randomly irrespective of application importance (real time or non real time). There is a need for employing flexible and adaptable scheme for packet dropping, which may not deteriorate QoS of a real time application in wireless multimedia networks.

As per the literature survey, it is observed that most of the buffer management schemes in wireless networks, employ queue length as a threshold for dropping. Parameters such as round trip time and application priority are also used to make packet dropping decisions. Apart from these parameters, channel conditions and flow rates may also affect the dropping decisions. The rigid packet dropping schemes probably may not give flexibility to the buffer manager, hence soft packet dropping schemes by using fuzzy logic may be employed for wireless multimedia networks. Fuzzy system's characteristics are based on the concept of fuzzy partitioning of the information. The decision-making ability of the fuzzy model depends on the existence of rule base and fuzzy reasoning mechanism. The fuzzy logic based solutions are better for uncertain inputs which minimize the random outputs.

Some of the works done in the area of packet dropping in wireless networks are as follows. An algorithm called RED (random early detection), which drops the packets of all flows with the same probability depending on queue length is given in [1]. Flows with reduced transmission rates and shorter packet length are given unfair share in RED. The work presented in [2] proposes a packet dropping scheme that drops the packet based on the sharing index computed using differentiated services. In [3], an adaptive fuzzy based control algorithm which computes the packet dropping probability according to pre-configured fuzzy logic using only queue length as input variable is presented.

A fair packet dropping algorithm is presented in [4] depending on channel condition and fairness so as to achieve trade off between throughput and fair service. An adaptive link layer retransmission and active drop mechanism based on fuzzy logic considering queue length as fuzzy parameter is discussed in [5]. A scheme for handoff calls is presented in [6] by using user population, used bandwidth and received power level as fuzzy parameters for handoff decision at the base station. The scheme given in [7] is based on the fuzzy timer which changes the transmission rate to avoid congestion. A fuzzy green controller is presented in [8] which is expected to act as congestion controller in the routers, for both conditions such as with background load and without back ground load. A congestion index has been proposed for intelligent packet dropping scheme by using fuzzy logic in [9].

In [10], average queue length is replaced by the product of delay and retransmission times of RTS (request to send) as congestion indicator to calculate dropping probability. The work presented in [11] is a new scheduling algorithm, named as Delay-sensitive Dynamic Fair Queuing (DSDFQ) algorithm which is designed to meet delay requirements of multimedia applications as well as maintain high network efficiency to adapt to load fluctuations of different traffic classes and varying wireless channel conditions caused by user mobility, fading and shadowing. In [12], estimation of channel capacity in rayleigh fading environment and the gaussian noise environment is presented. Development of a systematic methodology of fuzzy logic modeling with three distinct features such as reasoning formation, fuzzy clustering and membership assignment for input and output is presented in [13]. A simple channel predictor is presented in [14] that responds to impairments in the channel at the packet transmission time scale and makes a decision on packet transmission.

The work given in [15] presents multihop relaying technology utilization for mitigating unfairness in QoS, which comes about due to the location-dependent signal quality and multihop system. The work presented in [16] depicts that cooperative relaying can improve both capacity and fairness in cellular network. The work proposed in [17] is an adaptive multirate auto rate fallback (AMARF) algorithm, in which each data rate is assigned a unique success threshold. This criterion is used to judge when to switch a rate to the next higher one, and the success thresholds can be adjusted dynamically in an adaptive manner according to the running conditions such as packet length and channel condition parameters. In wireless networks, the available bandwidth undergoes fast time-scale variations due to channel fading and error from physical obstacles [18]. In [20], an adaptive QoS handoff priority scheme which reduces the probability of handoff call failures in a mobile multimedia network with a micro/picocellular architecture is presented. The scheme given in [21] supports voice and data calls with some grade of service guarantee in mobile multiservice networks.

The variations in the channel condition and data rate have motivated us to use fuzzy controller for developing a packet dropping scheme. The proposed scheme provides soft packet dropping, which minimizes the bursty losses and also provides better QoS for wireless multimedia network applications. Rest of the paper is organized as follows. Section 2 presents the proposed fuzzy based packet dropping scheme. Sections 3 and 4 present simulation model and results, respectively. Finally, section 5 concludes the work with some remarks.

2 Proposed work

The proposed scheme is located at the base station. The components of the scheme are as shown in figure 1. It consists of a knowledge base, buffer manager, and a fuzzy controller. Knowledge base comprises of status information of base station and the running applications. The status information of the base station include queue length, queue length threshold, maximum bandwidth, bandwidth available, maximum buffers, buffers available, channel condition, aggregated flow rate, aggregated departure rate, correction factor, etc. The information details about the running applications include application ID, priority, type (handoff/new), rate of flow at previous and current time instants, packet service time limit, packet departure time, buffer required, sustainable dropping threshold, dropping factor, etc.

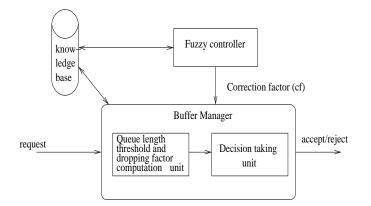


Figure 1: Fuzzy based packet dropping scheme

Buffer manager comprises of computational and decision taking units. Computational unit computes queue length, queue length threshold and dropping factor for an application packet request by using the data such as correction factor, priority, etc., available in the database. The dropping factor is made available to decision taking unit to decide either to drop or accept the packet. The buffer manager employs fuzzy controller to tune the queue length threshold. Buffer manager updates the knowledge base as and when required, either periodically or on arrival of application request. Fuzzy controller uses channel condition and packet flow rate of an application as inputs and provides correction factor as output, and updates the knowledge base. The correction factor is used to dynamically change the value of queue length threshold. The knowledge base is used for providing data for fuzzy controller to find out the correction factor. The computation of queue length, queue length threshold and dropping factor are discussed in the following sections.

2.1 Computation of queue length threshold and dropping factor

The following assumptions have been made to compute the queue length threshold and dropping factor. (1) Handoff calls and new calls both arrive at base station with different priorities, (2) packet dropping is considered in case of non-availability of buffers, and (3) priority is assigned depending upon the request time and current time within the different type of calls, i.e., high priority is assigned for handoff real-time calls. The computation of queue length, queue length threshold and dropping factor is divided into three phases: (1) queue length computation, (2) correction factor computation, and (3) priority an dropping factor computation.

The notations used in each of the phase are as follows. $q(t_i)$ - queue length in the buffer at time t_i , $r(t_i)$ - aggregated flow rate (sum of flow rate of all applications) at time t_i , $x(t_i)$ aggregated departure (sum of departure rate of all applications) rate of packets at time t_i , $q_{th}(t_i)$ - queue length threshold in buffer at time t_i , cf - correction factor, $chc(t_i)$ - channel condition at time t_i , df - dropping factor for a packet of an application, Bw_{max} - maximum bandwidth in the base station area, Bu_{max} - maximum buffer at the base station, and P - priority of a packet.

Computed queue length in Phase-I will be used in Phase-II and Phase-III for computing queue length threshold and dropping factor. Computation of queue length in the buffer is performed using two factors: aggregated flow rate and departure rate from the buffer. The queue length is given by formula 1.

$$q(t_i) = q(t_{i-1}) + (\delta r(t_i) * (t_i - t_{i-1})) - (x(t_i) * (t_i - t_{i-1}))$$
(1)

$$\delta r(t) = r(t_i) - r(t_{i-1}) \tag{2}$$

Where $\delta r(t_i)$ is the change of aggregated flow rate at time instant t_i .

In phase-II, correction factor (cf) is computed to tune the queue length threshold. The new queue length threshold computed in phase II is used by phase III for calculating the dropping factor. The correction factor is a value that depends upon channel conditions as given in formula 3. The value of cf is varied between 0 to 1.0 depending on fuzzy parameter chc. The variables 'h' and 'z' lie in the range 0 to 1 as defined by the system administrator.

$$cf = \begin{cases} 0 & \text{if } 0 < chc < h; \\ \mu(rf) \times \mu(chc) & \text{if } h \le chc < z; \\ 1.0 & \text{if } z \le chc < 1.0 \end{cases}$$
(3)

The correction factor depends upon the two fuzzy parameters, rate of flow $(\mu(rf))$ and channel condition $(\mu(chc))$ for one of cases. Where $\mu(rf)$ and $\mu(chc)$ are membership function for flow rate and channel condition, respectively. Before providing computation of new queue length threshold, let us first describe initial queue length threshold computation. Initial queue length threshold value is computed as given in formula 4 by using the Bu_{max} and threshold value setter f (varied in the range 0.7 to 0.9). The new queue length threshold is computed periodically as given in formula 5.

$$q_{th}(t_{initial}) = f * Bu_{max}.$$
(4)

$$q_{th}(t_i) = (q_{th}(t_{i-1})) - (cf * (q_{th}(t_{i-1}) - q_{th}(t_{i-2}))).$$
(5)

In phase-III, priority of the arriving packet and dropping factor are computed. The priority of the packet is used for computing dropping factor. The priority (p) of a packet from a connection of a call is computed (see formula 6) based on arrival time of a packet during a refreshing interval w, where, ref_{init} is the starting time of interval w and pkt_{arr} is packet arrival time. To provide higher priority to handoff packets, p' is recomputed using formula 7, where x is a value chosen between the range 1 to 3 (3 means 30 percent) or as desired by the administrator. The dropping factor of a packet is represented by formula 8 by considering n applications existing at the base station, where Bw denotes bandwidth used by an application. The computed value for dropping factor is communicated to buffer manager for taking the final decision to either drop or accept the packet. The accepted packet is placed into buffer. The fuzzy controller used in computation of correction factor is discussed.

$$p = 1.0 - (pkt_{arr} - ref_{init})/w.$$
(6)

$$p' = \begin{cases} p & for - handoff - packets \\ p/x & for - new - packets \end{cases}$$
(7)

$$df = \begin{cases} 1 & \text{if } q(t_i) > q_{th}(t_i); \\ 1 - max(q(t_i)/q_{th}(t_i), p') & \text{if } \sum_{k=1}^n Bw[k] > Bw_{max}; \\ 0 & \sum_{k=1}^n Bw[k] < Bw_{max} \end{cases}$$
(8)

2.2 Fuzzy controller

The fuzzy controller used in computation of correction factor in phase II is as shown in figure 2. The fuzzy inputs are flow rate and channel condition. The output parameter is a correction factor. Fuzzy controller has four components: fuzzification, inference, defuzzification and rule base. In the fuzzification step, fuzzy input parameter values (called crisp input) are converted into linguistic values (such as low, high or medium), each of which is represented by a fuzzy set. In the inference step, a set of rules called rule-base, which emulates the decision-making process of a human expert is applied to the linguistic values of the inputs to infer the output sets which represents the actual control signal for the process. We refer the reader to [22] for more complete background information on the fuzzy control.

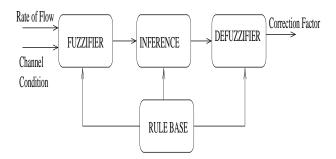


Figure 2: Fuzzy controller

Fuzzification

Fuzzy based packet dropping scheme considers two parameters for fuzzification: channel condition (chc) and the rate of flow of each application (rf). The output of linguistic parameter is the correction factor which is used for tuning of the queue length threshold. Membership

functions G(rf), G(chc) and G(cf) for each of the considered fuzzy parameters are depicted in figure 3, along with the linguistic values. Each of the fuzzy parameter is represented by triangular membership function since it represents minimum and maximum boundary conditions. The membership to each of fuzzy variables is assigned using intuition method.

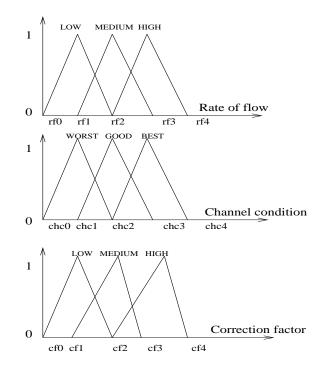


Figure 3: Membership function for input and output fuzzy parameters

- Rate of flow each application has different flow rate based on the requirements. The rate of flow is represented by linguistic values low (rf0 to rf2), medium (rf1 to rf3) and high (rf2 to rf4).
- Channel condition it is different at a given instant of time. Three linguistic values of channel conditions are represented by worst (chc0 to chc2), good (chc1 to chc3) and best (chc2 to chc4).
- Correction factor the correction factor is used as tuning parameter for queue length threshold. The linguistic values for correction factor are low (cf0 to cf2), medium (cf1 to cf3) and high (cf2 to cf4).

The fuzzy packet dropping scheme forms a fuzzy set of dimension G(rf) * G(chc). The values of fuzzy variables assigned depends on the network administrator, i.e., he/she can assign the different values at different instant of time depending upon the network conditions.

Inference and defuzzification

Mamdani controller is used for reasoning. The two fuzzy inputs are rf and chc with each input having three linguistic values. Total number of rules will be 9. The fuzzy rule base is as shown in figure 4. If the condition is true, we call the rule as being active. Each rule *i* is written as follows. For example, Rule i: IF rf_i is *low* and *chc_i* is *good*, THEN $cf_i = low$. To decide an appropriate output membership function, the strength of each rule is considered. For this reason, the output membership function is a complicated function and center of area method [22] is used for defuzzification. This method finds the center point of the fuzzy output membership function, which is used as output value. The defuzzified output parameter provides flexibility to the network administrator to perform soft packet dropping.

Rate of flow	Channel condition	Correction factor			
Low	Worst	High			
Medium	Worst	High			
High	Worst	High			
Low	Good	High			
Medium	Good	Medium			
High	Good	Medium			
Low	Best	Medium			
Medium	Best	Low			
High	Best	Low			

Figure 4: Fuzzy rule base table for correction factor

2.3 Algorithms

This section presents algorithmic description of the proposed work.

Algorithm 1: Packet dropping decision in buffer manager

Nomenclature: n = number of running applications, $Bu_{max} =$ Maximum buffer size at the base station, i = ith running application, $Bu_{req} =$ Buffer required, rf =Rate of flow, chc = Channel condition, $qt_i =$ queue length in the buffer at *i*th time instant, $qth_i =$ queue length threshold in the buffer at *i*th time instant, cf = correction factor, and df = dropping factor.

- 1. Receive buffer request Bu_{req} from new/handoff call to buffer manager;
- 2. If (buffers are available) then allocate and go to step 8; Else execute the following steps;
- 3. Initialize the computational unit with queue length and queue length threshold;
- 4. Check for channel conditions and set the cf according to formula 3; call algorithm 2 in case of chc lying between h and z.
- 5. Compute queue length threshold as given in the phase II (formula 5);
- 6. Compute df as given in the phase III (formula 8);
- 7. Decision unit in buffer manager compares the dropping factor with predefined dropping threshold; If the dropping factor is above the threshold, then the packet is dropped, else accepted;
- 8. Stop.

Algorithm 2: Computation of correction factor

- 1. Initialize fuzzy controller with aggregated rate of flow and channel condition;
- 2. Find the membership function of rate of flow and channel condition;
- 3. Compute cf by using rule base;
- 4. Return defuzzified cf;
- 5. Stop.

3 Simulation

The simulation of proposed work is carried out using C programming language in different scenarios on a Pentium IV machine. Simulation uses the network model, channel model, fuzzy model and traffic model. Models are presented as follows. Network model- A cell covers an area of X1 * Y1 sq. kms. Base station is located at some median point of the area. There are n number of mobile nodes arrive at the base station, which may include n1 handoff nodes and n2 new nodes that are generated with probability p and (1 - p), respectively. Mobile nodes move in random directions. Maximum bandwidth required in a cell is assumed to be Bw_{max} Mbps. Maximum buffer size at the base station of the cell is assumed to be Bu_{max} Gbytes.

Channel model - The wireless fading due to physical environment are often characterized to envelop fading of the carrier signal. Rayleigh fading has been used to model channel characteristics, where signal strength is varied by considering the mean and variance of signal strength. Each node is modeled to receive data if signal strength of the data is within the acceptable value. Fuzzy model - Fuzzy parameters considered are: rate of flow rf0, rf1, rf2, rf3 and rf4; channel condition chc0, chc1, chc2, chc3 and chc4; correction factor cf0, cf1, cf2, cf3 and cf4. The triangular membership is assigned to the input and output parameters because it is easy to compute and strictly specify the lower and upper limits of the fuzzy parameters. The intuition method of fuzzification is used to fuzzify the input parameters such as flow rate and channel condition. To decide an appropriate output crisp value for cf, we consider the center of gravity method. Traffic model - Poisson distribution is used to model the generation of data packets with arrival rate λ in interval t.

3.1 Simulation inputs

Following inputs are considered for simulation. X1 = 5 km., Y1 = 10 km., n=100, p = random number between 0 and 1, $n1=n^*p$, n2=(n-n1). Mobile nodes can move in any of eight directions: N, S, E, W, NE, NW, SE, SW. $Bw_{max}=20$ Mbps. $Bu_{max}=100$ Mbytes. Queue length threshold initialized with f=0.9. Queue length (qt_i) ranges from 10-90 Mbytes. The channel condition parameters such as μ set between 0.2 to 1.0, σ set between 0.3 to 0.5. The fuzzy input rate of flow data is randomly distributed between [rf0, rf2]=[500kbps,1mbps], [rf1, f3]=[750kbps,1.5mbps], [rf2, rf4]=[1mbps-2mbps]. Channel condition limits are z=0.9, h =0.1. Channel condition fuzzy values are: [chc0, chc2]=[0.1, 0.4], [chc1, chc3]=[0.2, 0.6], [chc2, chc4]=[0.4, 0.9]. The fuzzy output parameter correction factor is distributed between [cf0, cf2]=[0.0, 0.4], [cf1, cf3]=[0.2, 0.6], [cf2, cf4]=[0.4, 0.8]. The value of λ is positive real number between 0 to 10 and t is time interval set to 10.

Simulation procedure is as follows: Generate a wireless network and traffic across the network, apply the proposed model and find the dropping factor, and compute the performance of the system. The performance parameters measured in simulation are as follows. Acceptance of new

calls: It is defined as the ratio of acceptance of new calls to total number of new calls arrived. Acceptance of handoff calls : It is defined as the ratio of acceptance of handoff calls to total number of handoff calls arrived. Dropping of new calls : It is defined as the ratio of new calls dropped to total number of new calls arrived. Dropping of handoff calls: It is defined as the ratio of handoff calls dropped to total number of handoff calls arrived. Variations in queue length threshold: It is defined as change in the queue length threshold. Number handoff calls corrected: It is defined as number of handoff calls got corrected to total number of handoff calls arrived. Number of new calls corrected: It is defined as number of new calls got corrected to total number of new calls arrived.

4 Results

Figure 5 depicts the rise in acceptance of new calls upto certain limit with the variations in channel conditions. As the variance (sigma) of channel reduces, acceptance reduces. In a similar way, as the mean (mu) of channel condition reduces, acceptance reduces. This is because, the scheme attempts to get broader range of channel conditions for computing dropping factor. The packet drops for new calls change with respective change in the mean and variance of the channel conditions as observed in figure 6.

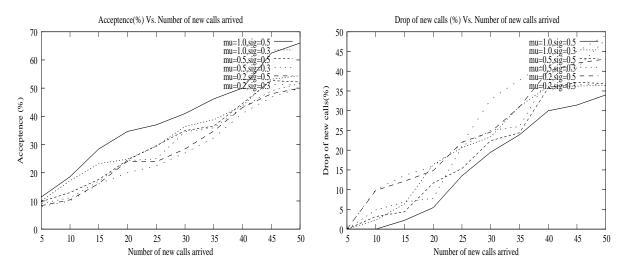


Figure 5: Acceptance of new calls (%) .Vs. New Figure 6: Dropping of new calls (%) .Vs. New call arrival calls arrival

More number of new calls get corrected as the number of calls increase as depicted in figure 7. This happens because of scarcity of buffers at the base station. The variations in acceptance for handoff calls is shown in figure 8. The handoff calls acceptance is more than new calls for the same mean and variance of the channel condition (see figures 5 and 8). This indicates that scheme gives higher priority to handoff and real-time calls.

Figure 9 shows the dropping percentage of handoff calls with respect to the handoff calls arrived. It is seen that less number of packets from handoff calls are dropped compared to new call packets (see figures 6 and 9). The corrected handoff calls increase with increase in the handoff call arrivals as observed in figure 10. The queue length threshold tuning with respect to channel conditions is shown in figure 11. It shows that as the channel condition becomes better, the tuning value increases to avoid buffer overflows. The queue length threshold variations with respect to time is shown in figure 12. This indicates the tuned value of queue length threshold at different instants of time.

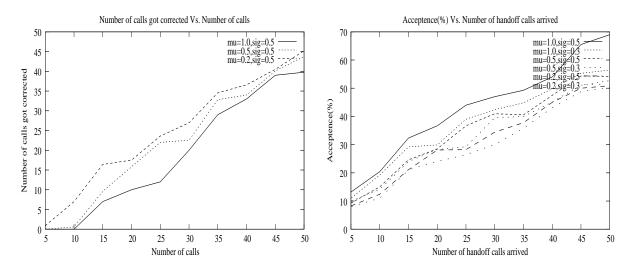


Figure 7: Corrected new calls (%) .Vs. New calls Figure 8: Acceptance of handoff calls (%) .Vs. arrival Handoff calls arrival

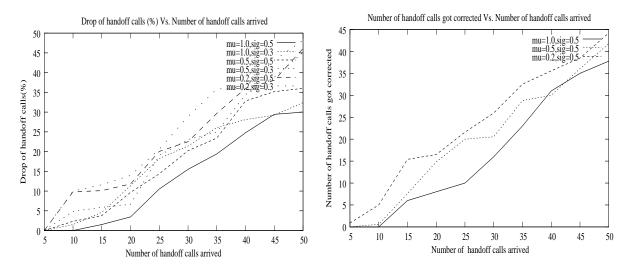


Figure 9: Dropping of handoff calls (%) .Vs. Figure 10: Number of handoff calls got corrected Handoff calls arrival (%) .Vs. Total handoff calls arrived

5 Conclusions

The main objective of the proposed work is to use the base station buffers efficiently and decrease the packet dropping especially for handoff and real-time calls during good channel conditions and high flow rate causing the buffer overflow. Extensive simulation results reveal that proposed scheme features very low call dropping, better handoff and new calls acceptance with different channel conditions. This is achieved because of tuning the queue length threshold by correction factor. The correction factor is derived by two fuzzy parameters such as channel condition and flow rate. The scheme can be extended for different fuzzy parameters such as device type, user priority, sustainable jitters, etc.

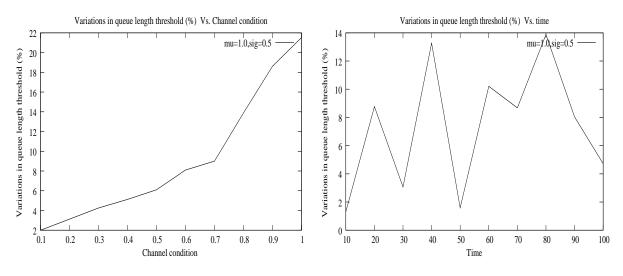


Figure 11: Queue threshold (%) .Vs. Channel Figure 12: Queue length threshold (%) .Vs. condition Time

Bibliography

- S. Floyd and V. Jacobson, Random Early Detection gateways for congestion avoidance. IEEE/ACM Transactions on Networking, vol. 1, pp. 397-413, Aug., 1993.
- [2] Yuan Chen and Lemin Li, A Fuzzy Fair Packet Dropping Algorithm Supporting Differentiated Services. *IEEE Proc. 5th International Conference on Computer and Information Technology* (CIT'05), Shangai, China, 2005.
- [3] Chonggang Wang, Bo Li, Kazem Sohraby and Yong Peng, An Adaptive Fuzzy-based Control Algorithm for Active Queue Management. Proc. IEEE International Conf. on Local Computer Networks, Bonn, Germany, pp.643-649, 2003.
- [4] Yuan Chen and Lemin Li, A wireless packet dropping algorithm considering fairness and channel condition. Proc. IEEE International Conf. on Communication, Circuits and Systems (ICCAS), Chengdu, China, vol. 1, pp. 369-373, June 2004.
- [5] Changvuan Luo and Chongsen Ran, An adaptive retransmission and active drop mechanism based on fuzzy logic. Proc. IEEE International Conf. on Radio Science, pp. 162-165, Aug. 2004.
- [6] Manpreet Dang, Amol Prakash, Manika and Rajeev, Fuzzy logic based handoff in indoor wireless networks. Proc. IEEE Vehicular Technology Conference, Tokyo, Japan, pp. 2375-2379, May 2000.
- [7] Huai-jen Liu, Chih-hsun Chou and Kuan-hu Ho, Fuzzy logic based solution for the congestion collapse problem. *Journal of Information Technology and Applications*, vol. 1, no.2, pp. 89-94 Sept., 2006.
- [8] Saman Taghavi Zargar and Mohammad Hossein, Fuzzy Green: A modified TCP equationbased on active queue management using fuzzy logic approach. *International Journal of Computer Science and Network Security*, vol. 6, no.5, pp. 50-58, May 2006.

- [9] Fan Yanfei, Ren Fengyuan and Lin Chuang, Active queue management based Fuzzy logic decision. Proc. International Conference Communication Technology, Beijing, China, pp. 286-289, 2003.
- [10] Dong Y., Makrakis D and Sullivan T, Network Congestion control in ad-hoc IEEE 802.11 Wireless LAN. Proc. IEEE Canadian Conference on Electrical and Computer Engineering, Montreal, Canada, vol. 3, pp. 1667-1670, 2003.
- [11] Huai-Rong Shao, Chia Shen, Daqing Gu, Jinyun Zhang and Philip Orlik, Dynamic Resource Control for High-Speed Downlink Packet Access Wireless. *Proc. ICDCS Workshop*, pp. 838-843, 2003.
- [12] William C.Y. and Lee, Estimation of channel capacity in rayleigh fading environment. *IEEE Transactions on Vehicular Technology*, vol. 39, no. 3, pp. 187-189, Aug. 1990.
- [13] Mohammad R., Emami I., Burhan T. and Andrew A, Development of a systematic methodology of fuzzy logic modeling. *IEEE Transactions on Fuzzy Systems*, vol. 6, no. 3, pp. 346-360, Aug. 1998.
- [14] Javier Gomez and Andrew T. Campbell, A Channel Predictor for Wireless Packet Networks. Proc. IEEE International Conference on Multimedia and Expo (ICME), New York, USA, July-August, 2000.
- [15] Jaeweon Cho and Zygmunt J. Haas, On the Throughput Enhancement of the Downstream Channel in Cellular Radio Networks Through Multihop Relaying. *IEEE Journal on Selected Areas in Communications*, vol. 22, no. 7, Sept. 2004.
- [16] Seungho Song, Kyuho Son, Hyang-Won Lee and Song Chong, Opportunistic Relaying in Cellular Network for Capacity and Fairness Improvement. Proc. IEEE Globecom, Washington, USA, Nov. 26-30, 2007.
- [17] Xi Yong, Huang Qingyan, Wei Jibo and Zhao Haitao, Rate adaptive protocol for multirate IEEE 802.11 networks. *Journal of Electronics China*, vol. 24, no.3, pp. 289-295, April 2007.
- [18] Samarth H. Shah, Kai Chen and Klara Nahrstedt, Available Bandwidth Estimation in IEEE 802.11-based Wireless Networks. Proc. of 1st ISMA/CAIDA Workshop on Bandwidth Estimation (BEst), San Diego, CA, Dec. 2003.
- [19] Wei Zhuang, Brahim Bensaou and Kee Chaing Chua, Adaptive quality of service handoff priority scheme for mobile multimedia networks. *IEEE Transactions on Vehicular Technology*, vol. 49, no. 2, pp. 494-505, March 2000.
- [20] Wei Zhuang, Brahim Bensaou and Kee Chaing Chua, Handoff priority scheme with preemptive, finite queueing and reneging in mobile multiservice networks. *Journal of Telecommuni*cation Systems, vol. 15, no. 1-2, pp. 37-51, Nov. 2000.
- [21] K. M. Passino and S.Yurkovich, *Fuzzy Control*, Addission Wesley, 1998.

Classification Performance Using Principal Component Analysis and Different Value of the Ratio R

J. Novakovic, S. Rankov

Jasmina Novakovic

"Faculty of Computer Science" Megatrend University Belgrade Serbia, 11000 Belgrade, Bulevar Umetnosti 29 E-mail: jnovakovic@megatrend.edu.rs

Sinisa Rankov

Megatrend University Belgrade Bulevar Umetnosti 29 E-mail: rankovs@megatrend.edu.rs

> Abstract: A comparison between several classification algorithms with feature extraction on real dataset is presented. Principal Component Analysis (PCA) has been used for feature extraction with different values of the ratio R, evaluated and compared using four different types of classifiers on two real benchmark data sets. Accuracy of the classifiers is influenced by the choice of different values of the ratio R. There is no best value of the ratio R, for different datasets and different classifiers accuracy curves as a function of the number of features used may significantly differ. In our cases feature extraction is especially effective for classification algorithms that do not have any inherent feature selections or feature extraction build in, such as the nearest neighbour methods or some types of neural networks.

> **Keywords:** feature extraction, linear feature extraction methods, principal component analysis, classification algorithms, classification accuracy.

1 Introduction

Data dimensionality reduction is an active field in computer science. It is a fundamental problem in many different areas, especially in forecasting, document classification, bioinformatics, and object recognition or in modelling of complex technological processes. In such applications datasets with thousands of features are not uncommon. All features may be important for some problems, but for some target concept only a small subset of features is usually relevant.

To overcome the curse of dimensionality problem, dimensionality of the feature space should be reduced. This may be done by selecting only the subset of relevant features, or creating new features that contain maximum information about the class label from the original ones. The former methodology is named feature selection, while the latter is called feature extraction, and it includes linear (PCA, Independent Component Analysis (ICA) etc.) and non-linear feature extraction methods. Finding new features subset are usually intractable and many problem related to feature extraction have been shown to be NP-hard ([1]).

Feature extraction brings the immediate effects for application: speeding up a data mining algorithm, improving the data quality and thereof the performance of data mining, and increasing the comprehensibility of the mining results. It has been a fertile field of research and development since 1970's in statistical pattern recognition ([2] and [3]), machine learning and data mining.

Some classification algorithms have inherited ability to focus on relevant features and ignore irrelevant ones. Decision trees are primary example of a class of such algorithms ([4] and [5]), but also multi-layer perceptron (MLP) neural networks with strong regularization of the input

layer may exclude the irrelevant features in an automatic way ([6]). Such methods may also benefit from independent feature selection or extraction. On the other hand, some algorithms have no provisions for feature selection or extraction. The k-nearest neighbour algorithm (k-NN) is one family of such methods that classify novel examples by retrieving the nearest training example, strongly relaying on feature selection or extraction methods to remove noisy features.

2 PCA

PCA is a standard statistical technique that can be used to reduce the dimensionality of a data set. It ([7]) is known as Karhunen-Loeve transform, has proven to be an exceedingly useful tool for dimensionality reduction of multivariate data with many application areas in image analysis, pattern recognition and appearance-based visual recognition, data compression, time series prediction, and analysis of biological data - to mention a few.

The strength of PCA for data analysis comes from its efficient computational mechanism, the fact that it is well understood, and from its general applicability. For example, a sample of applications in computer vision includes the representation and recognition of faces ([8], [9], [10] and [11]), recognition of 3D objects under varying pose ([12]), tracking of deformable objects ([13]), and for representations of 3D range data of heads([14]).

PCA is a method of transforming the initial data set represented by vector samples into a new set of vector samples with derived dimensions. The basic idea can be described as follows: a set of n-dimensional vector samples $X = \{x_1, x_2, x_3, ..., x_m\}$ should be transformed into another set $Y = \{y_1, y_2, ..., y_m\}$ of the same dimensionality, but y-s have the properties that most of their information content is stored in the first few dimensions. So, we can reduce the data set to a smaller number of dimensions with low information loss.

The transformation is based on the assumption that high information corresponds to high variance. If we want to reduce a set of input dimensions X to a single dimension Y, we should transform X into Y as a matrix computation

$$Y = A \cdot X \tag{1}$$

choosing A such that Y has the largest variance possible for a given data set. The single dimension Y obtained in this transformation is called the first principal component. This component is an axis in the direction of maximum variance. The first principal component minimizes the distance of the sum of squares between data points and their projections on the component axis. In practice, it is not possible to determine matrix A directly, and therefore we compute the covariance matrix S as a first step in features transformation. Matrix S ([15]) is defined as

$$S_{n \times n} = \frac{1}{n-1} \sum_{j=1}^{n} (x_j - x')^T \cdot (x_j - x')$$
(2)

where

$$x' = \frac{1}{n} \sum_{j=1}^{n} x_j \tag{3}$$

In the next step, the eigenvalues of the covariance matrix S for the given data should be calculated. Finally, the m eigenvectors corresponding to the m largest eigenvalues of S define a linear transformation from the n-dimensional space to an m-dimensional space in which the features are uncorrelated. To specify the principal components we need the following additional explanations about the notation in matrix S: 1) The eigenvalues of $S_{n \times n} \lambda_1, \lambda_2, ..., \lambda_n$, where $\lambda_1 \geq$

 $\lambda_2 \geq \dots \geq \lambda_n \geq 0$ and 2) The eigenvectors e_1, e_2, \dots, e_n correspond to eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$, and they are called the principal axes.

Principal axes are new, transformed axes of *n*-dimensional space, where the new variables are uncorrelated, and variance for the *i*-th component is equal to the *i*-th eigenvalue. Most of the information about the data set is concentrated in a few first principal components. In this paper we research how many of the principal components are needed to get a good representation of the data. In other words, we try to find the effective dimensionality of the data set. For this purpose we analyze the proportion of variance. Dividing the sum of the first m eigenvalues by the sum of all the variances (all eigenvalues), we will get the measure for the quality of representation based on the first m principal components. The result is expressed as a percentage. The criterion for features selection is based on the ratio of the sum of the m largest eigenvalues of S to the trace of S. That is a fraction of the variance retained in the m-dimensional space. If the eigenvalues are labeled so that $\lambda_1 \geq \lambda_2 \geq , ..., \geq \lambda_n$, then the ratio can be written as

$$R = \frac{\sum_{i=1}^{m} \lambda_i}{\sum_{i=1}^{n} \lambda_i} \tag{4}$$

All analyses of the subset of m features represent a good initial estimate of the *n*-dimensionality space if the ratio R is sufficiently large, it means greater than the threshold value. This method is computationally inexpensive, but it requires characterizing data with the covariance matrix S.

In implementation, the transformation from the original attributes to principal components is carried out through a process by first computing the covariance matrix of the original attributes and then, by extracting its eigenvectors to act as the principal components. The eigenvectors specify a linear mapping from the original attribute space of dimensionality N to a new space of size M in which attributes are uncorrelated.

The resulting eigenvectors can be ranked according to the amount of variation in the original data that they account for. Typically, the first few transformed attributes account for most of the variation in the data set and are retained, while the remainders are discarded.

PCA is an unsupervised method, which makes no use of information embodied within the class variable. Because, the PCA returns linear combinations of the original features, the meaning of the original features is not preserved. Over the years there have been many extensions to conventional PCA. For example, Independent Component Analysis (ICA) ([16] and [17]) is the attempt to extend PCA to go beyond decorrelation and to perform a dimension reduction onto a feature space with statistically independent variables. Other extensions address the situation where the sample data live in a low-dimensional (non-linear) manifold in an effort to retain a greater proportion of the variance using fewer components ([18], [19], [20], [21], [22] and [23]) and yet other (related) extensions derive PCA from the perspective of density estimation (which facilitate modeling non-linearities in the sample data) and the use of Bayesian formulation for modeling the complexity of the sample data manifold ([24]).

3 Classification Algorithms

Four supervised learning algorithms are adopted here to build models, namely, IB1, Naive Bayes, C4.5 decision tree and the radial basis function (RBF) network. This section gives a brief overview of these algorithms.

3.1 IB1

IB1 is nearest neighbour classifier. It uses normalized Euclidean distance to find the training instance closest to the given test instance, and predicts the same class as this training instance. If multiple instances have the same (smallest) distance to the test instance, the first one found is used. Nearest neighbour is one of the simplest learning/classification algorithms, and has been successfully applied to a broad range of problems ([25]).

To classify an unclassified vector X, this algorithm ranks the neighbours of X amongst a given set of N data $(X_i, c_i), i = 1, 2, ..., N$, and uses the class labels c_j (j = 1, 2, ..., K) of the K most similar neighbours to predict the class of the new vector X. In particular, the classes of these neighbours are weighted using the similarity between X and each of its neighbours, where similarity is measured by the Euclidean distance metric. Then, X is assigned the class label with the greatest number of votes among the K nearest class labels.

The nearest neighbour classifier works based on the intuition that the classification of an instance is likely to be most similar to the classification of other instances that are nearby to it within the vector space. Compared to other classification methods such as Naive Bayes', nearest neighbour classifier does not rely on prior probabilities, and it is computationally efficient if the data set concerned is not very large. If, however, the data sets are large (with a high dimensionality), each distance calculation may become quite expensive. This reinforces the need for employing PCA and information gain-based feature ranking to reduce data dimensionality, in order to reduce the computation cost.

3.2 Naive Bayes

This classifier is based on the elementary Bayes' Theorem. It can achieve relatively good performance on classification tasks ([26]). Naive Bayes classifier greatly simplifies learning by assuming that features are independent given the class variable. More formally, this classifier is defined by discriminant functions:

$$f_i(X) = \prod_{j=1}^{N} P(x_j \mid c_i) P(c_i)$$
(5)

where $X = (x_1, x_2, ..., x_N)$ denotes a feature vector and $c_j, j = 1, 2, ..., N$, denote possible class labels.

The training phase for learning a classifier consists in estimating conditional probabilities $P(x_j | c_i)$ and prior probabilities $P(c_i)$. Here, $P(c_i)$ are estimated by counting the training examples that fall into class c_i and then dividing the resulting count by the size of the training set. Similarly, conditional probabilities are estimated by simply observing the frequency distribution of feature x_j within the training subset that is labeled as class c_i . To classify a class-unknown test vector, the posterior probability of each class is calculated, given the feature values present in the test vector; and the test vector is assigned to the class that is of the highest probability.

3.3 C4.5 Decision Tree

Different methods exist to build decision trees, but all of them summarize given training data in a tree structure, with each branch representing an association between feature values and a class label. One of the most famous and representative amongst these is the C4.5 tree ([27]). The C4.5 tree works by recursively partitioning the training data set according to tests on the potential of feature values in separating the classes. The decision tree is learned from a set of training examples through an iterative process, of choosing a feature and splitting the given example set according to the values of that feature. The most important question is which of the features is the most influential in determining the classification and hence should be chosen first. Entropy measures or equivalently, information gains are used to select the most influential, which is intuitively deemed to be the feature of the lowest entropy (or of the highest information gain). This learning algorithm works by: a) computing the entropy measure for each feature, b) partitioning the set of examples according to the possible values of the feature that has the lowest entropy, and c) for each are used to estimate probabilities, in a way exactly the same as with the Naive Bayes approach. Note that although feature tests are chosen one at a time in a greedy manner, they are dependent on results of previous tests.

3.4 RBF Networks

A popular type of feed forward network is RBF network. Usually, the RBF network consists of three layers, i.e., the input layer, the hidden layer with Gaussian activation functions, and the output layer. Each hidden unit essentially represents a particular point in input space, and its output, or activation, for a given instance depends on the distance between its point and the instance-which is just another point. Intuitively, the closer these two points, the stronger the activation. This is achieved by using a nonlinear transformation function to convert the distance into a similarity measure. A bell-shaped Gaussian activation function, whose width may be different for each hidden unit, is commonly used for this purpose. The hidden units are called RBFs because the points in instance space for which a given hidden unit produces the same activation form a hypersphere or hyperellipsoid.

The output layer of an RBF network takes a linear combination of the outputs of the hidden units and in classification problems-pipes it through the sigmoid function. The parameters that such a network learns are (a) the centers and widths of the RBFs and (b) the weights used to form the linear combination of the outputs obtained from the hidden layer.

One way to determine the first set of parameters is to use clustering, without looking at the class labels of the training instances at all. The simple k-means clustering algorithm can be applied, clustering each class independently to obtain k basis functions for each class. Intuitively, the resulting RBFs represent prototype instances. Then the second set of parameters can be learned, keeping the first parameters fixed. This involves learning a linear model using one of the techniques such as, linear or logistic regression. If there are far fewer hidden units than training instances, this can be done very quickly.

RBF networks ([28]) are a special class of neural networks in which the distance between the input vector and a prototype vector determines the activation of a hidden neuron. Prototype vectors refer to centers of clusters obtained during RBF training. Usually, three kinds of distance metrics can be used in this network, such as Euclidean, Manhattan, and Mahalanobis distances. The RBF network provides a function $Y : \mathbb{R}^n \to \mathbb{R}^M$, which maps *n*-dimensional input patterns to *M*-dimensional outputs ($\{(X_i, Y_i) \in \mathbb{R}^n \times \mathbb{R}^M, i = 1, 2, ..., N\}$). Assume that there are *M* classes in the data set. The *m*-th output of the network is as follows ([29]):

$$y_m(X) = \sum_{j=1}^{K} w_{m_j} \theta_j(X) + w_{m_0} b_m$$
(6)

In this case X is the n-dimensional input pattern vector, m = 1, 2, ..., M, and K is the number of hidden units. M is the number of classes (outputs), w_{m_j} is the weight connecting the j-th hidden unit to the *m*-th output node, b_m is the bias, and w_{m_0} is the weight connecting the bias and the *m*-th output node.

The radial basis activation function $\theta(x)$ of the RBF network distinguishes it from other types of neural networks. Several forms of activation functions have been used in applications ([29]):

•
$$\theta(x) = e^{\frac{-x^2}{2\sigma^2}}$$
 (7)

•
$$\theta(x) = (x^2 + \sigma^2)^{-\beta}, \beta > 0$$
 (8)

•
$$\theta(x) = (x^2 + \sigma^2)^{\rho}, \beta > 0$$
 (9)

•
$$\theta(x) = x^2 \ln(x) \tag{10}$$

here σ is a parameter that determines the smoothness properties of the interpolating function. A disadvantage of RBF networks is that they give every feature the same weight because all are treated equally in the distance computation. Hence they cannot deal effectively with irrelevant features.

4 Experiments and results

Real datasets called "Statlog (Australian credit approval)" and "Statlog (German credit data)" for tests were used, taken from the UCI repository of machine learning databases. These datasets were used to compare the classification performance using IB1, Naive Bayes, RBF networks and C4.5 decision tree classifiers, in conjunction with the use of PCA and different value of the ratio R. The classification performance is measured using ten-foldcross-validation.

German Credit Data

This dataset classifies people described by a set of features as good or bad credit risks. Data set characteristics is multivariate, feature characteristics are categorical and integer. Number of instances is 1000, number of features is 20, and there are no missing values.

Australian Credit Approval

This file concerns credit card applications. Data set characteristics is multivariate; feature characteristics are categorical, integer and real. Number of instances is 690, number of features is 14, and there are missing values. This dataset is interesting because there is a good mix of features continuous, nominal with small numbers of values, and nominal with larger numbers of

values. There are also a few missing values.

Australian credit approval data	Classification accuracy without PCA				
Naive Bayes	77,7				
C4.5 decision tree	86,1				
IB1 classifier	81,2				
RBF network	79,7				

Table 1: Classification results for Australian credit approval data without using PCA.

The number of input compo-	20	22	25	29	29	30	32	33
nents produced - Australian								
credit approval data								
R	0,8	0,85	0,9	0,95	0,96	$0,\!97$	0,98	0,99

Table 2: The number of input components produced using PCA at various ratio R values for Australian credit approval data.

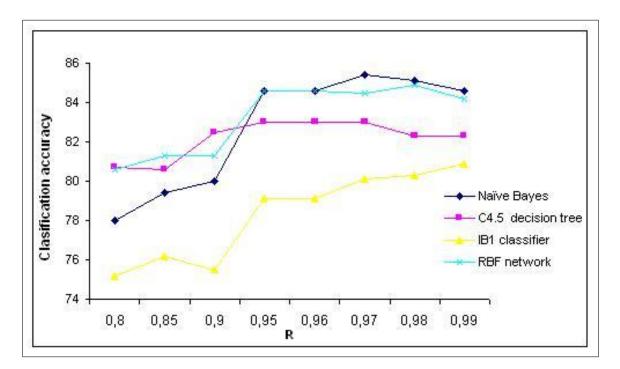


Figure 1: The classification performance using Naive Bayes, C4.5 decision tree, IB1 and RBF network classifiers, in conjunction with the use of PCA and different value of the ratio R. Statlog (Australian credit approval) data set

Classification results without using PCA as a standard statistical technique that can be used to reduce the dimensionality of a data set, for Australian credit approval are presented in Table 1, and for German credit data in Table 3.

Table 2 and 4 show the number of input components produced for each ratio R value investigated. It can be observe from the tables that the number of input components reduces with decreasing values of the ratio R used.

In Figure 1 classification performance for Naive Bayes and RBF network are significantly better with PCA. For IB1 and C4.5 decision tree classifiers the results are better without PCA.

For greater value of the ratio R classification accuracy of IB1 classifier is better. The others types of classifier in our experiment have the better results with greater value, but when value reached some boundaries performance of classifier are the same or worst.

German credit data	Classification accuracy without PCA					
Naive Bayes	75,4					
C4.5 decision tree	70,5					
IB1 classifier	72					
RBF network	74					

Table 3: Classification results for German credit data without using PCA.

The performance of some classifiers depends on its generalization capability, which in turn is dependent upon the data representation. One important characteristic of data representation is uncorrelated. This is because correlated data reduce the distinctiveness of data representation and thus, introduce confusion to the classifier model during the learning process and hence, producing one that has low generalization capability to resolve unseen data. The results demonstrated that the elimination of correlated information in the sample data by way of the PCA method improved Naive Bayes and RBF network classification performance (Figure 1).

At point 0.95% of the ratio R value with 29 input components, all classifiers significantly improved the classification accuracy. After that the ratio R value, the classification accuracy is about the same with little variations between classifiers. It suggests that this number of inputs is sufficiently optimal for the all classifiers to learn distinct features in the data and perform better input/output mapping.

The number of input com-	31	34	38	42	43	44	45	46
ponents produced - German								
credit data								
R	0,8	0,85	0,9	0,95	0,96	0,97	0,98	0,99

Table 4: The number of input components produced using PCA at various ratio R values for German credit data.

German credit data doesn't consist of correlated information caused by overlapping input instances. Without correlation in sampled data there is not confusion in classifiers during the learning process (Figure 2) and thus, no degrades their generalization capability. In this case all classifiers' classification performance doesn't improved by PCA. For greater value of the ratio R classification accuracy of IB1 classifier is better. Naive Bayes classifier has the worst results with greater values of the ratio R. Classification accuracy of RBF network have oscillation values. Classification accuracy for C4.5 decision tree doesn't change too much with different values of the ratio R.

This final part of the comparative study is set to investigate the differences between different classifiers, in terms of their classification ability. It is clear from Figures 1 and 2 that on average, Naive Bayes and RBF network classifiers tend to significantly outperform the decision tree and IB1 classifiers.

5 Conclusions

Feature extraction leading to reduced dimensionality of the feature space. PCA is one of the most popular techniques for dimensionality reduction of multivariate data points with application

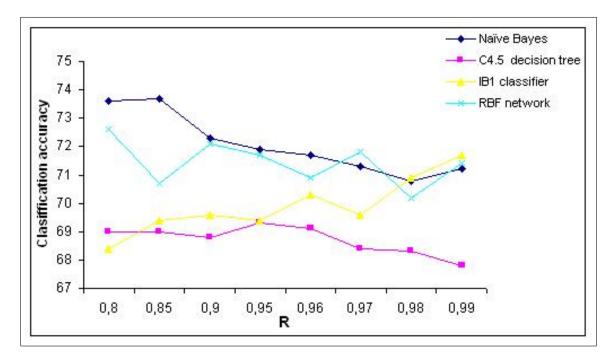


Figure 2: The classification performance using Naive Bayes, C4.5 decision tree, IB1 and RBF network classifiers, in conjunction with the use of PCA and different value of the ratio R. Statlog (German credit data) data set

areas covering many branches of science. This is especially effective for classification algorithms that do not have any inherent feature selections or feature extraction builds in, such as the nearest neighbour methods or some types of neural networks. PCA has been used for feature extraction with different values of the ratio R, evaluated and compared using four different types of classifiers on two real benchmark data sets. Accuracy of the classifiers is influenced by the choice different values of ratio R (Figure 1 and Figure 2).

There is no best value of the ratio R, for different datasets and different classifiers accuracy curves as a function of the number of features used may significantly differ. But, in more cases, the value of 0.95 gave the best results.

Several improvements of the feature extraction method presented here are possible:

- The algorithms and datasets will be selected according to precise criteria: different algorithms with PCA as linear feature extraction method, and several datasets, either real or artificial, with nominal, binary and continuous features.
- ICA and others linear feature extraction methods may be included.
- Problem of data dimensionality reduction may be analysed with non-linear feature extraction methods.

These conclusions and recommendations will be tested on larger datasets using various classification algorithms in the near future.

Bibliography

- A.L. Blum, R.L. Rivest, Training a 3-node neural networks is NP-complete, Neural Networks, 5:117 - 127, 1992.
- [2] N. Wyse, R. Dubes, A.K. Jain, A critical evaluation of intrinsic dimensionality algorithms. In E.S. Gelsema and L.N. Kanal, editors, *Pattern Recognition in Practice*, pp 415–425. Morgan Kaufmann Publishers, Inc., 1980.
- [3] M. Ben-Bassat, Pattern recognition and reduction of dimensionality, In P. R. Krishnaiah and L. N. Kanal, editors, *Handbook of statistics-II*, pp 773-791. North Holland, 1982.
- [4] L.Breiman, J.H. Friedman, R.H. Olshen, Stone C.J., Classification and Regression Trees, Wadsworth and Brooks, Monterey, CA, 1984.
- [5] J.R. Quinlan, C4.5: Programs for machine learning, San Mateo, Morgan Kaufman, 1993.
- [6] W. Duch, R. Adamczak, K. Grabczewski, A new methodology of extraction, optimization and application of crisp and fuzzy logical rules, IEEE Transactions on Neural Networks, vol. 12, pp. 277-306, 2001.
- [7] I.T. Jolliffe, Principal Component Analysis, Springer-Verlag, New York, 1986.
- [8] L. Sirovich, M. Kirby, Low dimensional procedure for the characterization of human faces, Journal of the Optical Society of America, 4(3) 519-524, 1987.
- [9] M. Turk, A. Pentland, Eigen faces for recognition, J. of Cognitive Neuroscience 3(1), 1991.
- [10] B. Moghaddam, A. Pentland, B. Starner, View-based and modular eigenspaces for face recognition, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recogni*tion (CVPR), pp 84-91, 1994.
- [11] P.N. Belhumeur, J.P. Hespanha, D.J. Kriegman, Eigenfaces vs. Fisherfaces: Recognition using class specific linear projection, *Proceedings of the European Conference on Computer* Vision, 1996.
- [12] H. Murase, S.K. Nayar, Learning and recognition of 3D objects from appearance, *IEEE 2nd Qualitative Vision Workshop*, pp 39-50, New York, NY, June 1993.
- [13] M. J. Black, D. Jepson, Eigen-tracking: Robust matching and tracking of articulated objects using a view-based representation, *Proceedings of the European Conference on Computer Vision (ECCV)*, pp 329-342, Cambridge, England, 1996.
- [14] J.J. Atick, P.A. Griffin, N.A. Redlich, Statistical approach to shape-from-shading: deriving 3d face surfaces from single 2d images, *Neural Computation*, 1997.
- [15] M. Kantardzic, Data Mining: Concepts, Models, Methods, and Algorithms, John Wiley & Sons, 2003.
- [16] P. Comon, Independent component analysis, a new concept? Signal processing pages 36(3), pp 11-20, 1994.
- [17] A.J. Bell, T.J. Sejnowski, An information maximization approach to blind separation and blind deconvolution, *Neural Computation*, pp 1129-1159, 1995.

- [18] C. Bregler, S.M. Omohundro, Nonlinear manifold learning for visual speech recognition, iccv, Boston, Jun 1995.
- [19] T. Heap, D. Hogg, Wormholes in shape space: Tracking through discontinuous changes in shape, *iccv*, 1998.
- [20] T. Hastie, W. Stuetzle, Principal curves, Journal of American Statistical Association 84, pp 502-516, 1989.
- [21] M.A. Kramer, Non linear principal component analysis using autoassociative neural networks, AI Journal 37(2), pp 233-243, 1991.
- [22] A.R. Webb, An approach to nonlinear principal components-analysis using radially symmetrical kernel functions, *Statistics and computing* 6(2), pp 159-168, 1996.
- [23] V. Silva, J.B. Tenenbaum, J.C. Langford, A global geometric framework for nonlinear dimensionality reduction, *Science*, 290, December 2000.
- [24] J.M. Winn, C.M. Bishop, Non-linear bayesian image modelling, Proceedings of the European Conference on Computer Vision, Dublin, Ireland, June 2000.
- [25] M. Kuramochi, G. Karypis. Gene classification using expression profiles: a feasibility study, International Journal on Artificial Intelligence Tools, 14(4):641-660, 2005.
- [26] P. Domingos, M. Pazzani, Feature selection and transduction for prediction of molecular bioactivity for drug design, *Machine Learning*, 29:103-130, 1997.
- [27] E. P. Xing, M. L. Jordan, R. M. Karp Feature selection for high-dimensional genomic microarray data, Proceedings of the 18th International Conference on Machine Learning, 601-608, 2001.
- [28] C.M. Bishop, Neural Network for Pattern Recognition, Oxford University Press Inc., New York, 1995.
- [29] L. Wang, X. Fu, Data Mining with Computational Intelligence, Springer-Verlag Berlin Heidelberg, Germany, pages 9-14, 2005.

Modeling Uncertainty in a Decision Problem by Externalizing Information

I. Parpucea, B. Pârv, T. Socaciu

Ilie Parpucea

Babes-Bolyai University Faculty of Economic Sciences and Business Administration Romania, 400591 Cluj-Napoca, 58-60 Teodor Mihali E-mail: ilie.parpucea@econ.ubbcluj.ro

Bazil Pârv

Babes-Bolyai University Faculty of Mathematics and Computer Science Romania, 400084 Cluj-Napoca, 1 M. Kogălniceanu E-mail: bparv@cs.ubbcluj.ro

Tiberiu Socaciu

Stefan cel Mare University Faculty of Economic Sciences and Public Administration Romania, 720229 Suceava, 13 Universității E-mail: socaciu@seap.usv.ro

Abstract: This paper deals with decision problems under uncertainty. The solution of a decision problem involves observation, processing, and modeling of statistical data in order to quantify the uncertainty. Better data measurement and estimation of uncertainty add more consistency to the solution of a decision problem. The paper proposes a new way of predicting the Bayesian-Nash equilibrium which uses information sources to measure new information received by information consumers. Thus, the estimation of uncertainty is based on a more solid mathematical foundation, needed (as in the case of artificial intelligence) to produce logical inferences. From another perspective, the externalization of information helps the software designers to produce better software architectures for decision support systems. An theoretical example illustrates a market situation with a small number of firms, each firm's output being likely to have a large impact on the market price.

Keywords: Bayesian-Nash equilibrium, information source, conditional probability distribution.

1 Introduction

Game theory, founded by von Neumann and Morgenstern (1947), studies situations in which multiple agents or players interact in order to each maximize an objective (payoff) function. The payoff function of a player is determined not only by its own actions, but also by the actions of other players. In a game with incomplete information, the payoffs also depend on information that is private to the individual agents. This information is known as an agent's *type*.

Bayesian decision theory is concerned with the question of how a player (decision maker) should choose a particular action from a set of possible choices if the outcome of the choice also depends on some unknown state (from the states of the world). In our approach, the decision maker is modeling the information received by the system (i.e. new information) as an *information source* ([2]). A decision problem involves one or several information sources. We

assume that each person is able to represent his beliefs, as the likelihood of the different n states of the information source, by a subjective discrete probability distribution ([5]).

The structure of this paper is as follows. After this introductory section, the next one introduces the teoretical background related to games with incomplete information and information sources. The third section presents the Bayes-Nash equilibrium in the presence of information sources, and the fourth discusses a market-related example. The last section compares the proposed approach with the classical one, and outlines future work.

2 Theoretical background

This section introduces the main concepts discussed in this paper. It starts with some definitions stating the context, and then defines the concept of information source, both taken from [1] and [3]. The last sub-section introduces the new concept of Bayesian-Nash equilibrium based on information sources.

In what follows, *information* means a message about an event that has occurred, will occur, or is likely to occur. The received information regarding a possible realization of an event is extremely important. Information is a particular case of reflection, as an interaction between two processes; one's properties (the process that generates or *produces* information) will be reproduced in another process or several other processes (that *consume* information). Interaction between two or more processes involves an exchange of information.

2.1 Games with incomplete information

Definition 1. A game with incomplete information ([1]), is denoted by:

$$\Gamma_t = (I, (F_i)_{i \in I}, (p_t^i(f, \theta))_{i \in I}, (\Theta_i)_{i \in I}, \mu_t),$$
(1)

where:

- I is the set of players, |I| = m,
- F_i is the strategy set for player $i, i = \overline{1, m}$, and $F = F_1 \times F_2 \times \cdots \times F_m$ is the set of all possible strategy profiles;
- $f = (f_1, f_2, \cdots, f_m) \in F$ is a joint strategy or strategy profile;
- Θ_i is the set of types for the player *i*, and $\Theta = \Theta_1 \times \Theta_2 \times \cdots \times \Theta_m$ is the joint type space;
- $\theta = (\theta_1, \theta_2, \cdots, \theta_m) \in \Theta$ is the joint type of all players;
- $p_t^i(f,\theta)$ is the payoff function for player *i* at the moment *t* if the strategy *f* and the type combination θ are chosen. Note that the payoff for the player *i* may depend not only on its type θ_i , but also on the other players' type, denoted by θ_{-i} .
- μ_t the probability distribution on the set Θ at the moment t.

In our exposition, we assume that type sets Θ_i are finite; consequently, Θ is a finite set also. $\mu_t(\theta), \theta \in \Theta$ denotes the probability of chosing type combination θ at the moment t. As in [4], we assume, without loss of generality, that players have incomplete information about their opponents' payoffs but have complete information about the strategies of all other players.

The following definition, taken from [1,3], introduces the classical Bayes-Nash equilibrium of a game with incomplete information.

Definition 2. A strategy profile $f(\theta) = (f_1(\theta_1), f_2(\theta_2), \dots, f_m(\theta_m))$ constitutes a **Bayes-Nash** equilibrium of a game Γ_t with incomplete information if the following inequality:

$$\sum_{\theta_{-i}\in\Theta_{-i}} p_t^i(f_i^*(\theta_i), f_{-i}^*(\theta_{-i}), \theta_i, \theta_{-i}) \mu_t(\theta_{-i}|\theta_i) \ge \sum_{\theta_{-i}\in\Theta_{-i}} p_t^i(f_i(\theta_i), f_{-i}^*(\theta_{-i}), \theta_i, \theta_{-i}) \cdot \mu_t(\theta_{-i}|\theta_i)$$
(2)

holds for all possible players $i \in I$ and all types $\theta_i \in \Theta_i$ and all strategies $f_i \in F_i$.

2.2 Information sources

A process (information producer or consumer) is specified by a set of n variables, denoted by $V = \{V_1, V_2, \dots, V_n\}$, where V_j is the set of values for the j^{th} variable: $V_j = \{v_j^1, v_j^2, \dots, v_j^{n_j}\}$. The state of a process at a certain moment is given by the vector $v = \{v_1, v_2, \dots, v_n\}, v_j \in V_j, j = \overline{1, n}$. Future values of all variables are random, and the realization of the states depends on received information. This information is in message form, decreasing or increasing the uncertainty of the realization of an event.

An information source is a way of specifying the states of a process, regarding one or several variables. The information source assigned to the j^{th} variable is denoted by S^j , and the set of distinct values $v_j^k \in V_j, k = \overline{1, n_j}$, represents a complete space of events. The simultaneous realization of two events is impossible, and the union of the events represents a certain event. A state s_j^k of the information source S^j is assigned to each event $v_j = v_j^k$. For a bounded interval of time, only information sources with a finite number of states will be taken into account.

Consider the following assumptions:

- each player is able to represent his beliefs, as to the likelihood of the different n_j states of the information source S^j , by a subjective discrete probability distribution.
- the information source S^j has discrete states and the individual is supposed to be able to assign to each state s_j^k a degree of belief, in the form of (normalized) numerical weights p_j^k , between zero and one and whose sum is one:

$$\forall j = \overline{1, n} : 0 \le p_j^k \le 1, \forall k = \overline{1, n_j}; \sum_{k=1}^{n_j} p_j^k = 1.$$

 $(p_j^k$ is the probability that the state s_j^k occurs).

If the information source S^{j} has n_{j} states, the set of states and probabilities defined at a moment t, forms a discrete random variable denoted by:

$$S_t^j: \begin{pmatrix} s_j^k \\ p_j^k(t) \end{pmatrix}_{k=\overline{1,n_j}}, j=\overline{1,n}.$$

A simple information source is an information source defined with respect to a single variable V_j . A complex information source is an information source defined with respect to two or more variables, which can be independent or dependent. In the second case (i.e. the variables are related to one another), the mathematical model of the complex information source needs to contain this dependency.

In order to illustrate how a complex information source is constructed, let's consider for the beginning the simplest case of two independent variables, $V_1 \in V$ and $V_2 \in V$ and assign to each variable a simple information source, S_t^1 and S_t^2 respectively:

$$S_t^1 : \begin{pmatrix} s_1^k \\ p_1^k(t) \end{pmatrix}_{k=\overline{1,n_1}}, S_t^2 : \begin{pmatrix} s_2^l \\ p_2^l(t) \end{pmatrix}_{l=\overline{1,n_2}}$$

where $s_1^k, k = \overline{1, n_1}$ are the states of information source S_t^1 and $s_2^l, l = \overline{1, n_2}$ are the states of information source S_t^2 .

The complex information source $SC_t^{1,2}$, built with respect to the variables V_1 and V_2 at the moment t, has the following mathematical model:

$$SC_t^{1,2}: \begin{pmatrix} s_1^k s_2^j \\ p_1^k(t) p_2^j(t) \end{pmatrix}_{k=\overline{1,n_1,j}=\overline{1,n_2}}.$$

Now let us consider the case of two

Now let us consider the case of two dependent variables, with the simple information source S_t^1 assigned to the independent variable V_1 , and the simple information source S_t^2 assigned to the dependent variable V_2 . The discrete random variable S_t^1 is:

$$S_t^1 : \left(\begin{array}{c} s_1^{\kappa} \\ p_1^k(t) \end{array}\right)_{k=\overline{1,n_1}}.$$

For a state of the source S_t^1 , denoted by s_1^k , the information source S_t^2 conditioned by the state s_1^k is defined as follows:

$$S_t^{2/1}(s_1^k) : \left(\begin{array}{c} s_2^l \\ p_2^{k,l}(t) \end{array}\right)_{l=\overline{1,n_2}},$$

where $p_2^{k,l}(t) = P(S_t^2 = s_2^l | S_t^1 = s_1^l)$ is the probability of occurence of state s_2^l conditioned by state s_1^k .

The complex information source $SC_t^{2/1}$, constructed by considering the two variables, has the following form:

$$SC_t^{2/1} : \left(\begin{array}{c} s_1^k s_2^l \\ p_1^k(t) \cdot p_2^{k,l}(t) \end{array}\right)_{k=\overline{1,n_1}, l=\overline{1,n_2}}$$

where the probability of occurence of the state $s_1^k s_2^l$ is:

 $P(S_t^1 = s_1^k; S_t^2 = s_2^l) = P(S_t^1 = s_1^k) \cdot P(S_t^2 = s_2^l | S_t^1 = s_1^k) = p_1^k(t) \cdot p_2^{k,l}(t).$ If S_t^2 is a discrete probability distribution and S_t^1 is an information source, then, according

to the above discussion, we can say that S_t^2 is a distribution conditioned by information source S_t^1 . Therefore we have a probability distribution updated by an information source.

3 Bayes-Nash equilibrium in the presence of information sources

Let us consider now the probability distribution μ_t defined on the discrete set Θ and a information source S_t (simple or complex), common to all players. According to the probability distribution conditioned by an information source discussed above, we have:

 $P(\mu_t = \theta; S_t = s^j) = P(S_t = s^j) \cdot P(\mu_t = \theta | S_t = s^j)$, where $\theta \in \Theta$ and $s^j \in S_t$.

3.1Notations

Considering the following notations:

- $P(\mu_t, S_t)$ the joint probability of μ_t and S_t occurring simultaneously, referred to as the historic probability distribution,
- $P(\mu_t/S_t)$ the probability of μ_t occurring conditional on S_t having occurred (i.e. the conditional probability of μ_t given S_t), also known as the probability distribution of information source, and
- $P(S_t)$ the marginal probability of S_t , referred to as the posterior probability distribution,

the following equation holds (as in [7]): $P(\mu_t/S_t) = \frac{P(\mu_t,S_t)}{P(S_t)}$. According to the above, posterior means historical updated with information, i.e. the probability distribution μ_t conditioned by the information source S_t , denoted by μc_t , is the probability distribution μ_t updated by the information source S_t .

The information source S_t is a possible probability distribution of states at the moment t+1. Denoting the game with incomplete information at the moment t+1 and based on S_t with Γ_{t+1} , it can be defined recursively as follows:

 $\Gamma_{t+1} = \Gamma_t(S_t)$, where $\Gamma_{t+1} = (I, (F_i)_{i \in I}, (p_t^i(f, \theta))_{i \in I}, (\Theta_i)_{i \in I}, \mu c_t).$

The game Γ_{t+1} is the updated game Γ_t based upon S_t . This information source updates probability distribution μ_t on Θ and thus the equilibrium of Γ_t is modified.

3.2 Decision functions and strategies

If a player receives information about his/her own type, then he/she can choose a particular strategy to maximize his/her expected payoff.

Definition 3. A decision function of player $i \in I$, denoted by $f_i(.)$, is a function that, for each type $\theta_i \in \Theta_i$, specifies the strategy $f_i(\theta_i) \in F_i$ this player will choose if his/her type turns out to be θ_i .

Let $\mu c_t^j(\theta_{-i}|\theta(i))$ be the updated probability obtained by using Bayesian updating rule, of a particular type combination for the opponents θ_{-i} , given that player *i* has type θ_i . For each type profile $\theta \in \Theta$, there are updated beliefs for each player, i.e. a list of conditional probability distributions $(\mu c_t^1(\theta_{-1}|\theta_1), \cdots, \mu c_t^m(\theta_{-m}|\theta_m))$. Players' beliefs after they have received information about their types, are no longer identical.

Definition 4. The strategy combination of all players except player *i*, that will be played according to the decision functions $f_{-i}(.)$, if type combination θ_{-i} occurs, is a list of decision functions $f_{-i}(.) = (f_1(.), \dots, f_{i-1}(.), f_{i+1}(.), \dots, f_I(.))$ for all players (other than player *i*) and $\theta_{-i}(.) = (\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_I)$, a type combination for the other players, $f_{-i}(\theta_{-i})$, that is: $f_{-i}(\theta_i) = (f_1(\theta_1), \dots, f_{i-1}(\theta_{i-1}), f_{i+1}(\theta_{i+1}), \dots, f_m(\theta_m))$.

3.3 Bayes-Nash equilibrium in the presence of information sources

The above definitions allow us to give the following:

Definition 5. The *Bayes* - *Nash equilibrium* of the game Γ_{t+1} is a list of decision functions $(f_1^*(.), \dots, f_I^*(.))$, such that for all possible players $i \in I$ and all types $\theta_i \in \Theta_i$:

$$\sum_{\theta_{-i}\in\Theta_{-i}} p_t^i \left(f_i^*(\theta_i), f_{-i}^*(\theta_{-i}), \theta_i, \theta_{-i} \right) \cdot \mu c_t(\theta_{-i}|\theta_i) \ge \sum_{\theta_{-i}\in\Theta_{-i}} p_t^i \left(f_i, f_{-i}^*(\theta_{-i}), \theta_i, \theta_{-i} \right) \cdot \mu c_t^i(\theta_{-i}|\theta_i)$$
(3)

holds for all strategies $f_i \in F_t$.

The equilibrium of Γ_{t+1} can differ from the equilibrium of Γ_t due to the information in S_t . For a given player $i \in I$ the updated equilibrium of Γ_{t+1} is:

$$f_i^*(.) = \sum_j f_{ij}^*(.) \cdot p^j(t),$$

where $f_{ij}^*(.)$ is the equilibrium of player *i* for the state s^j of the information source S_t and $p^j(t)$ is the probability that $S_t = s^j$. The above equation represents the updated equilibrium as a weighted average of all equilibria for the states of information source.

4 A market-related example

Consider a game with incomplete information, where the players are two firms supplying slightly different products (produced with zero production costs, as in [1]), with prices denoted

by p_1 and p_2 . As a result of new received information, the variation of each price around the average price can be modeled using two simple information sources, S_t^1 and S_t^2 , as it follows:

$$S_{t}^{1}: \begin{pmatrix} s_{1}^{1}: p_{1} \leq \overline{p_{1}} & s_{1}^{2}: p_{1} > \overline{p_{1}} \\ \pi_{1} & 1 - \pi_{1} \end{pmatrix}; S_{t}^{2}: \begin{pmatrix} s_{2}^{1}: p_{2} \leq \overline{p_{2}} & s_{2}^{2}: p_{2} > \overline{p_{2}} \\ \pi_{2} & 1 - \pi_{2} \end{pmatrix}$$

where $\overline{p_1}$ represents the average price of the first firm, and $\overline{p_2}$ the average price of the second firm. Both information sources have two states, (s_1^1, s_1^2) and (s_2^1, s_2^2) , with the probabilities $(\pi_1, 1 - \pi_1)$ and $(\pi_2, 1 - \pi_2)$, respectively.

4.1 Notations

For our example, consider the following problem-specific notations:

• The demand functions for the goods of the two firms:

 $d_1(p_1, p_2) = a \cdot \Delta p_1 + b \cdot \Delta p_2, d_2(p_1, p_2) = c \cdot \Delta p_1 + d \cdot \Delta p_2,$

where $\Delta p_i = p_i - \overline{p_i}$, $i = \overline{1, 2}$, are the deviation of price p_i from the average price $\overline{p_i}$. Firm one does not know parameters c and d; firm two does not know parameters a and b.

• The sets of possible types of the two players:

$$\Theta_1 = \{(a_i, b_j), i = 1, 2; j = 1, 2\}, \Theta_2 = \{(c_i, d_j), i = 1, 2; j = 1, 2\}$$

• The payoff functions of the two players:

$$\Pi_1 (p_1, p_2, (a_i, b_j)) = (a_i \cdot \Delta p_1 + b_j \cdot \Delta p_2) \cdot p_1, \Pi_2 (p_1, p_2, (c_i, d_j)) = (c_i \cdot \Delta p_1 + d_j \cdot \Delta p_2) \cdot p_2.$$

4.2 Building the complex information source

With the above notations, the complex information source is rewritten as:

 $SC_{t+1}: \begin{pmatrix} p_1 \leq \overline{p_1} \land p_2 \leq \overline{p_2} & p_1 \leq \overline{p_1} \land p_2 > \overline{p_2} & p_1 > \overline{p_1} \land p_2 \leq \overline{p_2} & p_1 > \overline{p_1} \land p_2 > \overline{p_2} \\ \pi_1 \cdot \pi_2 & \pi_1 \cdot (1 - \pi_2) & (1 - \pi_1) \cdot \pi_2 & (1 - \pi_1) \cdot (1 - \pi_2) \end{pmatrix},$ no matter what dependency relation is between the two prices considered. The information source SC_{t+1} describes the behavior of the market of the both products considering the variation of their prices in the next period, as a result of the received information.

In a Bayesian-Nash equilibrium, each firm is supposed to choose a type contingent strategy, that is decision functions $p_1(.)$ and $p_2(.)$ respectively, which is the best response to the opponent's decision function. In this example, μc_t is a probability distribution defined on $\Theta_1 \times \Theta_2$ and conditioned by the information source SC_t . For a state $s_1^k s_2^l$ of source SC_t we build two conditional distributions, μc_t^1 and μc_t^2 as can be seen next.

4.3 Computing the best response for the firm one

Consider $p_2(.) = (p_2(c_1, d_1), p_2(c_1, d_2), p_2(c_2, d_1), p_2(c_2, d_2))$ as given (fixed) and suppose that firm one has just learned that it has the demand parameters (a_1, b_1) . Firm one's expected payoff can be rewritten as:

$$\Pi_{1} (p_{1}(a_{1}, b_{1}), p_{2}(c_{1}, d_{1})) \cdot \mu c_{t}^{1} ((c_{1}, d_{1})|(a_{1}, b_{1})) + \\ +\Pi_{1} (p_{1}(a_{1}, b_{2}), p_{2}(c_{1}, d_{2})) \cdot \mu c_{t}^{1} ((c_{1}, d_{2})|(a_{1}, b_{1})) + \\ +\Pi_{1} (p_{1}(a_{2}, b_{1}), p_{2}(c_{2}, d_{1})) \cdot \mu c_{t}^{1} ((c_{2}, d_{1})|(a_{1}, b_{1})) + \\ +\Pi_{1} (p_{1}(a_{2}, b_{2}), p_{2}(c_{2}, d_{2})) \cdot \mu c_{t}^{1} ((c_{2}, d_{2})|(a_{1}, b_{1})) = \\ = a_{1} \cdot p_{1}^{2}(a_{1}, b_{1}) + p_{1}(a_{1}, b_{1}) (b_{1} \cdot \overline{p_{2}}(c_{i}d_{j}|a_{1}b_{1}) - a_{1}\overline{p_{1}} - b_{1}\overline{p_{2}}).$$

$$(4)$$

In the above equation we used the following notation for the average price $\overline{p_2}$ conditioned by state (a_1, b_1) :

$$\overline{p_2}(c_i, d_j)|a_1b_1) = p_2(c_1, d_1) \cdot \mu c_t^1((c_1, d_1)|(a_1, b_1)) + p_2(c_1, d_2) \cdot \mu c_t^1((c_1, d_2)|(a_1, b_1)) + p_2(c_2, d_1) \cdot \mu c_t^1((c_2, d_1)|(a_1, b_1)) + p_2(c_2, d_2) \cdot \mu c_t^1((c_2, d_2)|(a_1, b_1)).$$

The payoff function (4) is continuously differentiable in firm one's strategy $p_1(.)$. Therefore, any p_1 satisfying the first - order condition for a maximum will be the best response to the type-contingent strategy $p_2(.)$, previously considered. Solving the first - order condition for the maximum of the expected payoff function (4), one obtains the following best response $p_1^*(a_1, b_1)$ for firm one of type (a_1, b_1) to the decision functions $(p_2(c_1, d_1), p_2(c_2, d_2), p_2(c_2, d_1), p_2(c_2, d_2))$ of firm two:

$$p_1^*(a_1, b_1) = \frac{1}{2} (\overline{p_1} + \frac{b_1}{a_1} (\overline{p_2} - \overline{p_2} (c_i d_j | a_1 b_1))).$$
(5)

In a similar way, one obtains for all other types of firm one:

$$p_1^*(a_1, b_2) = \frac{1}{2} (\overline{p_1} + \frac{b_2}{a_1} (\overline{p_2} - \overline{p_2} (c_i d_j | a_1 b_2))), \tag{6}$$

$$p_1^*(a_2, b_1) = \frac{1}{2} (\overline{p_1} + \frac{b_1}{a_2} (\overline{p_2} - \overline{p_2} (c_i d_j | a_2 b_1))), \tag{7}$$

$$p_1^*(a_2, b_2) = \frac{1}{2} (\overline{p_1} + \frac{b_2}{a_2} (\overline{p_2} - \overline{p_2} (c_i d_j | a_2 b_2)).$$
(8)

4.4 Computing the best response for the firm two

Now consider that firm two learns that its type is (c_1, d_1) . For a fixed type-contingent strategy of firm one $p_1(.) = (p_1(a_1, b_1), p_1(a_1, b_2), p_1(a_2, b_1), p_1(a_2, b_2))$, the expected payoff of firm two will be as follows:

$$\Pi_{2} (p_{1}(a_{1}, b_{1}), p_{2}(c_{1}, d_{1})) \cdot \mu c_{t}^{2} ((a_{1}, b_{1})|(c_{1}, d_{1})) + \\ +\Pi_{2} (p_{1}(a_{1}, b_{2}), p_{2}(c_{1}, d_{2})) \cdot \mu c_{t}^{2} ((a_{1}, b_{2})|(c_{1}, d_{1})) + \\ +\Pi_{2} (p_{1}(a_{2}, b_{1}), p_{2}(c_{2}, d_{1})) \cdot \mu c_{t}^{2} ((a_{2}, b_{1})|(c_{1}, d_{1})) + \\ +\Pi_{2} (p_{1}(a_{2}, b_{2}), p_{2}(c_{2}, d_{2})) \cdot \mu c_{t}^{2} ((a_{2}, b_{2})|(c_{1}, d_{1})) = \\ = d_{1} \cdot p_{2}^{2}(c_{1}, d_{1}) + p_{2}(c_{1}, d_{1}) (c_{1} \cdot \overline{p_{1}}(a_{i}b_{j}|c_{1}d_{1}) - c_{1}\overline{p_{1}} - d_{1}\overline{p_{2}}).$$
(9)

In the equation (9) of the payoff function for the second firm, the average price $\overline{p_1}$ conditioned by the state (c_1, d_1) , is given by:

$$\overline{p_1}(a_i, b_j)|c_1d_1) = p_1(a_1, b_1) \cdot \mu c_t^2((a_1, b_1)|(c_1, d_1)) + p_1(a_1, b_2) \cdot \mu c_t^2((a_1, b_2)|(c_1, d_1)) + p_1(a_2, b_1) \cdot \mu c_t^2((a_2, b_1)|(c_1, d_1)) + p_1(a_2, b_2) \cdot \mu c_t^2((a_2, b_2)|(c_1, d_1)).$$

First-order condition gives the best response function for a firm of type (c_1, d_1) :

$$p_2^*(c_1, d_1) = \frac{1}{2} (\overline{p_2} + \frac{c_1}{d_1} (\overline{p_1} - \overline{p_1}(a_i b_j | c_1 d_1))).$$
(10)

A similar calculation yields firm two's best response for all other types of the following typecontingent strategies:

$$p_2^*(c_1, d_2) = \frac{1}{2} (\overline{p_2} + \frac{c_2}{d_1} (\overline{p_1} - \overline{p_1}(a_i b_j | c_1 d_2))), \tag{11}$$

$$p_2^*(c_2, d_1) = \frac{1}{2} (\overline{p_2} + \frac{c_1}{d_2} (\overline{p_1} - \overline{p_1}(a_i b_j | c_2 d_1))), \tag{12}$$

$$p_2^*(c_2, d_2) = \frac{1}{2} (\overline{p_2} + \frac{c_2}{d_2} (\overline{p_1} - \overline{p_1}(a_i b_j | c_2 d_2))).$$
(13)

4.5 Conclusion

In order to find a Bayesian-Nash equilibrium, one has to solve the system of equations given by the best response functions. With two players and four types for each player, this leads to a system of eight equations. The solution, $(p_1^*(.), p_2^*(.))$, is the Bayesian -Nash equilibrium. The probability of realization of equilibrium prices $(p_1^*(.), p_2^*(.))$, as a result of information received, is equal to the probability of realization of state $s_1^k s_2^l$ of complex source SC_{t+1} .

5 Conclusions and Future Works

5.1 Our approach vs other approaches

The main points of our approach are as follows:

- The *complexity of uncertainty* is given by the great (huge) number of variables; when the complexity of a decision problem (and the number of components dominated by uncertainty) grows, it is recommended to use a Bayesian network ([6]); in our case, we use a simple Bayesian network, subject to a learning algorithm;
- The original idea is to *separate* the set of problem components into two disjoint subsets: (a) deterministic components, and (b) components dominated by uncertainty; the separation of game information into external and internal can be done for each decision problem dominated by uncertainty;
- This separation allows you to study the influence of each individual factor to the solution of the game in a more efficient way; also, it suggests some architectural patterns (styles) to be used when designing a decision support system. The paper [8] discusses this issue in more detail.

The essential difference between the classic approach and those proposed in this paper is given by the separation of the information external to the game from the game-specific information. This separation follows the *separation of responsibilities* principle. This way, both external and internal elements of the game are easier to model and understand.

The classical approach does not make any difference between these two categories of information; more precisely, the influence of external information on the uncertainty that dominates the game is not taken into account/quantified. By splitting the game information into external and internal, the former being modeled by information sources, the influence of external environment on the variation of the solution is better captured and quantified. This provides a better evaluation of the contribution of individual factors to the predicted equilibrium.

Another advantage of this separation is that it allows a better, easier calibration of the model, by comparing the computed equilibrium with real solution, taken from historical data.

5.2 Future work

Our future efforts are directed to apply this general algorithm to various games with incomplete information and to build decision support systems based on it.

Acknowledgements

This work was supported by the grant ID_2586, sponsored by NURC - Romanian National University Research Council (CNCSIS).

Bibliography

- [1] Eichberger, I., Game Theory for Economists, Academic Press, 1993.
- [2] Florea, I., Parpucea, I., Economic Cybernetics (Romanian), Babeş-Bolyai University, 1993.
- [3] Fudenberg, D., Tirole, J., Game Theory, MIT Press, 1991.
- [4] Harsanyi, J.C., Games with Incomplete Information Played by Bayesian Players, Parts I, II, III, Management Science, 14, 1967, pp. 159-182, 320-334, 486-502.
- [5] Hirshleifer ,J., Riley J. G., The Analytics of Uncertainty and Information, Cambridge University Press, 1995.
- [6] Jensen, F.V., Nielsen, T.D., Bayesian Networks and Decision Graphs, Springer, 2nd ed., 2007
- [7] Koop, G., Bayesian Econometrics, Wiley, 2003.
- [8] Pârv, B., Parpucea, I., Computing Bayes-Nash Equilibrium for Games with Incomplete Information by Using Information Sources, submitted to *IJCCC*.
- [9] Reiz, B., Csato, L. Tree-Like Bayesian Network Classifiers for Surgery Survival Chance Prediction, Int. J. of Computers, Communications & Control, III, 2008, pp. 470-474.

Impact of Poor Requirement Engineering in Software Outsourcing: A Study on Software Developers' Experience

I. Perera

Indika Perera

Department of Computer Science and Engineering, University of Moratuwa, Sri Lanka E-mail: indikaperera@uom.lk

> **Abstract:** The software Requirement Engineering (RE) is one of the most important and fundamental activities in the software life cycle. With the introduction of different software process paradigms, the Requirement Engineering appeared in different facets, yet remaining its significance without a doubt. The software development outsourcing is considered as a win-win situation for both developed and developing countries. High numbers of low paid, yet talented workforce in developing countries could be employed for software outsourcing projects with the demanding power of the outsourcer to decide the projects, their scope and priorities with the intention of profit maximization. This study was conducted to analyze the impact of poor Requirement Engineering in outsourced software projects from the developers' context (sample size n = 57). It was identified that the present outsourcing scenario has created to have frequent requirement changes, shrunk design and stretched development phases, and frequent deliverables, which have to be accommodated by the software developer with extra effort and commitment beyond the project norms. The results reveal important issues and open policy level discussions while questioning our insights on the outsourcing benefits as a whole.

> **Keywords:** Requirement Engineering, Software Outsourcing, Project Management, Software Developer Productivity

1 Introduction

With the development of Internet and other communication media, the entire world is becoming a single village with the reach of each other at arms length. Many explain this scenario with the combination of other concepts as the Globalization. Some views Globalization with a positive look, where as others have strong critiques on its impact to individual nations and ultimately to the global economy. The development of the software and computing technologies has triggered the Globalization and world's operations to a large scale. In fact, the software outsource process gains its significance in the global economy with more and more projects being outsourced [1]. Software applications are complex and intangible products, which are difficult to manage. Hence, software lifecycle management becomes one of the key research areas in Software Engineering. Due to the nature of the software, software researchers and practitioners are focused to improve the software processes, which are used to develop software [2]. The underline assumption is that there is a direct correlation between the quality of the process and the quality of the developed software [3]. A software process can be considered as a set of tools, methods, and practices, where we use to produce a software product [4]. Requirement Engineering is one of the key process phases in the software development process; in most cases, essential enough to determine the success or failure of the project. Therefore, having an error free RE activity and specification can be more significant for the project success than any other activity. Software requirement errors can cost up to 70-85% of additional cost of reworking in an average software project [5]. This shows the importance of managing the Requirement Engineering works productively especially in an outsourced project where the cost and profit concerns are high motives [18]. Furthermore, it is a known fact that many small to medium scale software companies practice flexible, lightweight, Agile like informal processes to suit their development; more or less for the sake of having a process practice to convince the outsourcers and win the projects [6]. However, in practice, these organizations tend to follow ad-hoc set of activities than any established process, which are camouflaged as the Agile practice. These informalities make the RE phase more complex, resulting the developers' programming work unnecessarily tedious and unproductive. This research has tried to study on the issues of poor RE practices with outsourced software projects and how it affects to junior level software developers' productivity and work life balance. The rest of the paper presents the research work and outcomes with the following organization. Section 2 will discuss the background literature related to this study. Section 3 discusses the research problem in brief. Thereafter, the section 4 is included with the experiment methodology and section 5, the analysis section explains the analysis done based on the collected results. The conclusion with the possible future directives is included thereafter along with the acknowledgements. The references will compile the paper.

2 Background

This section contains a comprehensive synopsis of the referred literature for this study in the areas of Software Requirement Engineering and Software Outsourcing. It is important to mention that, there have very few research works been done to study how the current software outsourcing practices affect to the software developers work.

2.1 Software Requirements Engineering

Software Requirement Engineering is a mature and one of the key life cycle activities in any software process paradigm. As Nuseibeh and Easterbrook explained, software systems' Requirements Engineering (RE) is the process of discovering that purpose, by identifying stakeholders and their needs, and documenting these in a form that is amenable to analysis, communication, and subsequent implementation [7]. More precisely, Zave clearly defines the Requirement Engineering as; "Requirements Engineering is the branch of software Engineering concerned with the real world goals for, functions of, and constraints on software systems. It is also concerned with the relationship of these factors to precise specifications of software behavior, and to their evolution over time and across software families." [8]. Sometimes, it is impossible to have complete and finalized set of requirements at the beginning of a project. This allows requirement changes to happen during the latter stages of the project and create conflicts with the software process been practiced. To overcome such scenarios Agile like flexible processes have been introduced and well established in the software industry. Agile process welcomes frequent requirement changes even at late stages of the project. With frequent deliverables, Agile process measures its progress through the norm of working software [9]. However, the most important point to remember, even if someone practices the Agile process is, that still he needs to perform a certain set of Requirement Engineering activities to ensure smooth and correct process flow [17]. Therefore, whatever the process practice used in the development, following proper Requirement Engineering activities is a must, even though the practices may have different facets.

2.2 Software Outsourcing

There are number of definitions for the concept outsourcing in general and specifically in the context of software development. Rajkumar and Mani explain offshore outsourcing as: "When the supplier of software development is from another country than the firm that decides to outsource information systems" [10]. The software development outsourcing has, in the recent past, taken up global dimensions in which companies from the developed countries sub contract software lifecycle function to the developing countries [11]. Importantly, outsourcing has become an appealing option to organizations operating within the global economy for a variety of reasons [12]. One of the major and foremost reasons to outsource software development is the lower cost and relative high profit margins for outsourcer [13]. Another advantage of software development outsourcing is the opportunity for the outsourcer to focus more on their core businesses and strategic activities [1]. Sometimes, this makes the outsourcing company to act solely as an intermediary organization, except for legal and contractual works. Furthermore, there are some issues with outsourcing to countries with vast heterogeneous socio-economic and cultural factors. Having local representatives may be a key option to overcome cultural and linguistic barriers in offshore outsourcing [14], [15]. However, this can be an extra cost item and will not help developers if the representatives are not technically sound. Nevertheless, the outsourcing attracts both developed and developing countries due to their benefits from different perspectives. Developing countries trying to secure more and more outsourcing projects, aiming the monetary gains, without conducting grass root level research on how it impacts to their workforce, industry and society at large.

3 Research Problem

The research problem which this study has focused mainly based on the software developers' experience in an outsourcing environment. There have been number of business and economic forums to attract outsourcing projects for developing countries, but not visible effort been done for assessing the impact of developing an outsourced industry on various contexts, except the monetary gains and employment opportunities. Since the majority of studies done to examine how to manage the outsource projects form the view point of the outsourcer, there is a significant vacuum of knowledge to understand the issues faced by the end developers due to the poor Requirement Engineering. As the outsourcing process consists with two parties mainly; the outsourcer and the service/solution provider, there are different issues for both parties depending upon their context. In many cases, the software developers are employed by the solution providers to engage in outsourced projects. This research primarily focuses on the problems that software developers encounter during their project activities, due to the poor software requirement practices of the projects. It is a known fact that many software outsourced projects get schedule overrun and extended developer time. One research found out that 60-80% of software projects get schedule or effort overrun and the average overrun varies from 30-40% of the initial project scope [16]. Why this level of scope variation during the project happens even with mature, experienced staff is a concern need to be examined with extensive studies. There can be other factors, which cause the over time work. Furthermore, the software developers' attitude towards such over work situation need to be assessed. This was a major objective of this study. Another questioning area is whether the process of passing requirements to the developer either directly from the main client or from the outsourcer, is effective and efficient enough to facilitate the development process. With the technological advancement and the studies to alleviate geographical barriers, have come up with various methods of communication for outsourcing [15]. In fact there are number of informal communication approaches used in the software outsourcing and this study has analyzed the developer perception towards each approach and their relative significance to hinder the development activities. Identifying Key Success Factors for software process practices may be helpful for the improvement of current processes [19]. This study will also help to provide necessary guidance for possible improvements to the Requirement Engineering activities in outsourced software projects.

4 Research Methodology

The research methodology was based on a comprehensive survey conducted using the Sri Lankan outsourcing companies and their employees. Since the main objectives of the study and the research problem are generic to any outsourcing company and country as a whole the survey and its findings could easily be generalized with minor alterations. For this study, 57 software developers from 8 Small to Medium software developing enterprises were used. All these developers are engaged in outsourced projects for more than one year. To avoid any extreme cases in the study and to reduce skill variances, the developers who have industrial experience for less than 3 years were selected. The study was on voluntary basis and to strengthen the trust on the confidentiality of the participants' responses the study was conducted on individual basis without the intervention of the participants' employing organization. The study had two major phases for data collection. First all the participants were individually contacted through e-mail and telephone to explain the objective of the study and the confidentiality of their responses. The study was designed flexible enough to conduct by the participants as their own experiment. For the first phase, a worksheet was given to each participant to fill every day for 3 weeks. The worksheet was simple to fill and it has only four fields to be entered, based on the following measurement parameters.

- N Number of requirement changes/emerge for a given day
- $T_{\mathcal{S}}$ Start time of the work
- $T_{\cal F}$ Finish time of the work
- T_M -Time spent on meetings/discussions to clarify/understand requirements for a given day.

These are basic parameters for this study and those were used to analyze the average daily commitment of the developer and as a validating factor for the second phase data. The second phase was based on a questionnaire which was distributed electronically through e-mail and the participants had to fill their one and send back. Strong assurances were given on the confidentiality of their responses in order to have more accurate information. The following section describes the collected data and their analysis

5 Analysis

During the first phase of the study, the above mentioned 4 parameters were recorded and the distribution of average time spent by each developer on meetings/ discussions to clarify requirements was calculated using the following model (1). μ_{M_j} is the mean time spent for the j^{th} developer, and T_{M_i} is the time spent on the i^{th} day up to n, i.e. the number of days he worked during the experiment.

$$\mu_{M_j} = \frac{\sum_{i=1}^n T_{M_i}}{n} \tag{1}$$

The obtained values are sampled for displaying purpose and shown in the Table 1.

<u>S the development</u>	
Time Spent on Meetings for Requirements	Number
less than 0.5 hours	5
0.5 - 1 hours	9
1 - 1.5 hours	13
1.5 - 2 hours	14
2 - 2.5 hours	11
2.5 - 3 hours	3
3 - 3.5 hours	1
more than 3.5 hours	1

Table 1: The distribution of average daily times spent on meetings/discussions for clarifying requirements during the development

To obtain the population mean value μ_M of T_M , the following model (2) was used and m represents the number of participants in the sample.

$$\mu_M = \frac{\sum_{j=1}^{m} \frac{\sum_{i=1}^{n} T_{M_{ij}}}{n}}{m}$$
(2)

The obtained μ_M value for the entire population is 1.571 hours. The value shows that on average a developer has to spend more than one and half hours of their development time, daily to clarify requirements that they are developing. The main problem with this is that these times are not included in the project schedules and the developers have to get these times from their own expense as the deliverable deadlines are fixed.

Similarly with the same meaning of the notations as above, the population mean value for N, i.e. μ_N was calculated with the following model in (3).

$$\mu_{M} = \frac{\sum_{j=1}^{m} \frac{\sum_{i=1}^{n} N_{ij}}{n}}{m}$$
(3)

The obtained value for μ_N was 2.84. This indicates that nearly three requirements alterations happen per day on average during the development phase. There are many reasons for such significant change to occur. The first fact is that everybody wants to have the working software so quickly and success of getting future work depends on the delivery of current work. Because of this rapid nature of project development, the outsourcer cuts down the time on Requirement Engineering, design, and testing phases of the life cycle. However, when the frequent deliverable is done the outsourcer frequently improves his expected functionalities and gives frequent alterations to the planed requirements. Moreover, due to the high demanding power, the outsourcer gets what he requests over multiple times of changing the same requirements.

Finally the work pattern of the participants were analyzed with the model described in (4).

$$\mu_{W_j} = \frac{\sum_{i=1}^{n} \{T_{F_i} - T_{S_i} - K\}}{n} \tag{4}$$

In this model T_{F_i} and T_{S_i} are the finish time and the start time of a developer on i^{th} day, respectively. K is an average constant to represent the summation of none working times during a given day. It was considered to be 1.5 hours based on preliminary studies. With that model the following sample of work patterns were derived for a given week. In fact there is a reason to analyze this for a week specially. All the developers have in their employment contract that they are supposed to work on weekdays i.e. Monday to Friday and they have the opportunity of holidays including Saturday, Sunday and any other declared holidays by the government. However, it was observed that most of the participants have worked in Saturdays, special holidays, and in some cases Sundays to meet their deadlines. The mean value (μ_W) for the population, was

Average working hours per week	Number
35 - 40 hours	2
40 - 45 hours	7
45 - 50 hours	16
50 - 55 hours	21
55 - 60 hours	10
60 - 65 hours	1

Table 2: The distribution of average working hours per week

calculated using the model described in equation (5); hence, $\mu_W = 51.76$ hours, which indicates an average developer spends nearly 12 hours of additional work than their contractual norm of 40 hours per week. That means they put effort on additional work load on average worth of 29.4% of their monthly salary. In fact when the results finding were informed to the people who participated, they were astonished, and one of their response is worth to mention here. "...I can't believe this. This means up to now I have worked nearly one million rupees work just for free ...I should have thought of this before ..." [Participant's response]

$$\mu_W = \frac{\sum_{j=1}^m \sum_{i=1}^n \{T_{F_i} - T_{S_i} - K\}}{m}$$
(5)

5.1 Questionnaire Analysis

The questionnaire was consisted with 16 questions. Some of the questions were simple Yes/No type and some have to be answered with further details. Out those 16 questions, 6 were used to condition the participant's mind to answer correctly, and also used to validate the accuracy of their responses for the main questions. For this analysis, only the relevant questions, i.e. the 10 main questions were used.

Q1. Do you get a sufficient Requirement Specification for the project?

For this question 41 participants said 'No', 12 said 'Yes' and 4 said Don't Know. Therefore as a percentage, 71.92% believes that they do not get a sufficient requirement specification for their development.

Q2. What are the most frequent methods you use to communication on Requirements? Please rank them

The Table 3 summarizes the responses for the Question 2. Only 13 developers are working with proper Requirements Specification in which only 5 have the opportunity to consider it as the

Table 5. The distribution of average w	Orking	s noui	s per	WEEK	
Method of Communication	1^{st}	2^{nd}	3^{rd}	4^{th}	Total
Formatted Req. Spec. as E-mail attachment	5	6	2	0	13
E-mail messages	24	21	8	4	57
Instant Messaging (IM) / Chat	19	18	11	9	57
Voice Conference / Telephone	6	10	24	11	51
Video Conference	3	2	12	19	36
Face to Face communication	0	0	0	4	4
Travel to Client places for Discussions	0	0	0	1	1

Table 3: The distribution of average working hours per week

main form of requirement communication. Lack of proper Requirement Specification makes it enormously vulnerable to requirement alternations and scope creep, during the later stages of the project. Also only 1 participant has the opportunity to go to client premises and experience their expectations where as only four have the opportunity to meet the clients for discussions. Email messages considered as the main form of communication followed by instant messaging or chat. Also everybody (all 57) use E-mail and IM/chatting with different significance. Voice conferencing and Video conferencing are used as secondary and tertiary forms of communication when extra clarification is needed on requirements, in generally.

Q3. What are the reasons you think for requirement issues in your projects?

Reason	Number
Poor Technical Knowledge at Outsourcer	37
Optimistic Scheduling	29
Scope Creep and Stretching Requirements	42
Poor Communication	17
Informal Ways of Practice	31
Management Issues	24
Practice of "Do What They Say"	44

Table 4: The distribution of average working hours per week

The Table 4 indicates developers' responses on what they believe as the causes of requirement issues. There were multiple reasons given and most of the participants have concluded that the current practices in the outsourced industry are mainly responsible for these RE issues. As mentioned in the Background section, the Agile practices has significant impact on this issue as when it is practiced without a proper knowledge on how to use, the Agile process becomes a worse practice and creates a chaotic environment. According to other findings of this study it shows that developers have to accommodate such chaos within their activities and meet the deadlines and outcomes as expected. However, these results open a huge necessity for a discussion to improve the outsourcing industry in the developing countries.

Q4. Do you experience spending more time on assigned tasks than they were scheduled? If yes what is the reason for that?

As the Figure 1 shows, a significant impact to schedule overrun is due to the requirement issues. The other significant factors are technical difficulties, which may be due to lack of training and development opportunities for developers, and the Management issues specially having optimistic schedules to delight the outsourcer. The most important fact is that all of the participants experience to spend more time than they were scheduled, making them to work more time than the employment norms, without any compensation for the extra effort.

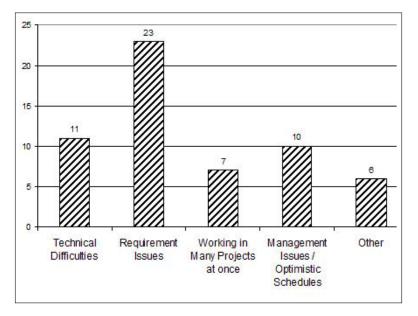


Figure 1: Reasons make the schedule overruns

Q5. Do you experience difficulties for understanding requirements from foreigners? Comment on this

Out of the 57 participants, 49 has responded and 37 had issues with verbal communication due to the accent of the foreigners. 42 of the participants complained about email and instant messages with insufficient information to take a decision; As a result, they have been spending more time and effort to clarify requirements. More importantly, the natural language based e-mails or chat massages, instead of a proper requirement specification make it further difficult to understand what the client is actually expecting.

Q6. If you are provided with complete and final requirement specification at the beginning of development, will it helps to your work?

23 people said it will help them lot and 9 were not sure about whether it helps them or not. 25 have indicated that it will not help as they will never getting a finalized requirement specification at any time of the project. This shows how much significant the issue of ever changing and poor Requirement Engineering process in the outsourced industry.

Q7. Are you the direct person to get requirements from client? If not how many to be passed for requirements to reach you?

8 participants have direct access with their clients to get requirements, 17 participants get through one person, 29 participants get through two persons, and 3 participants gets through three persons. The most obvious scenario was to have one person (non-technical) at the out-source end and one person (technical) at the development end to form 2 person layer to process requirements. Due to the non-technical managerial person at the client end, the requirements process deters its effectiveness, making developers to work more on requirements fine tuning and error fixing during the development phase. Q8. How good are you at work and life balancing?

The idea of this question is to assess the health and social impact from the outsourcing work to the developers. There was a 5 level lickert scale given to rank them considering their extra activities, exercise time, recreational activities and good life style practices. Results are shown in figure 2. 45.61% of the participants are not pay attention as required, for their social, health and other extra activities. When they spend more time on their work they get less time to attend those activities. The impact of this may not visible for short term but there are huge social, economical and health consequences in the long run which may not be easily compensated by their wee bit of extra salary they get compared to other disciplines. The good point is still the majority are happy about their work life balance, but there should be a proper mechanism to improve the situation.

Q9. Are you properly compensated for your effort?

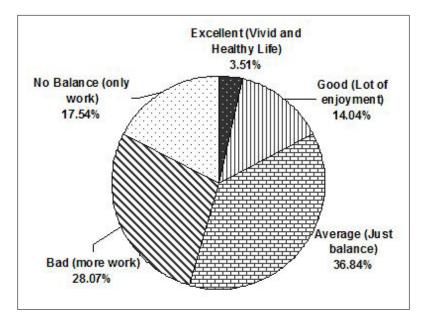


Figure 2: Work - Life balance of the developers

For this question the 21 participants (36.84%) said Yes and 26 participants (45.61%) have said No. 3 people said they are not sure about it. 7 participants did not respond to this question. In fact, out of the participants who responded, more than 50% is unhappy about their compensation and believe that they are under paid and over utilized.

Q10. What is the impact from time zone difference for your work times?

This was an open ended question where developers have to express their feeling about the impact from time zone difference. Idea of this question was to assess whether their work patterns have been affected from that. In fact, many have responded in similar fashion indicating there is a significant impact from the time zone difference. Most of the outsourcing projects are from United States of America and European Union member countries, where there is a significant time difference with the Sri Lankan time, around 11-13 hours and 4-6 hours, respectively. What the participants have indicated was that when the outsourcer wants to change a requirement or impose a new one, the developers have to communicate with them which will mostly the evening or late night in Sri Lankan Local time. In the following day they have to implement those from the morning. Therefore, invariably they have to spend at least 11-12 hours per day despite what was agreed in the contract as 8 hours. And these don't count for any overtime payment or additional benefit.

The above individual analyses of questions provide a clear view on how the poor Requirement Engineering practices affect to software developers in outsourced projects. The questions were focused into different aspects of significance relevant to the research problem. Even though the fact that, responds for these questions covers a larger spectrum of issues related to the problem, there should be further researches with more focus to these significant areas identified in this analysis.

5.2 Study Limitations

In this study, there were some research limitations which are worth to mention. Since this study was involved with human activities, this had the experimental limitation of different skill levels between the participants. However this was general scenario of the software industry and has no significant impact for this study's outcome.

The second and the most crucial limitation with this study was the difference between the requirement changes. Some requirements were very simple and some were not. Though it was really hard to eliminate this, in the general situation every project has both difficult requirement changes and simple changes. Therefore, in the large population conditions (n > 30), this behaves with normal distribution and nullifies the impact as the standard error is 0.

Another limitation with the study was the truncation errors of the collected data. Literally, what happened was, the developers were confident on expressing their values with integer figures of hours or in minutes over the decimal or fractional values. For an example they might have said actual time as 7.15 but the time may be 7.11 However, since there is no comparison done based on their given values, the impact for the study outcome was negligible. Furthermore, these errors are also normally distributed with the standard error 0 in the large population samples.

6 Conclusion

Despite the outcome of this research, there are some possible future studies relevant to this research, which can be considered as further extensions. This study mainly focused on the Requirement Engineering issues in the development phase of the outsourced software project life cycle. There are other important life cycle phases in the outsourced projects such as System Testing, System Design, System Implementation, etc. Then again, there can be similar future studies to examine the social, economic and health impact from malpractices of software outsource industry in the context to employees, countries, and to industry as in whole. Furthermore, impact on software outsourcing due to different cultural characteristics, social and demographic parameters, and geographical and ethnic factors can be another fruitful research area.

Domain specific examine of the impact of poor Requirement Engineering on developers is another possible study as a further extension to this study, which will indeed provide an in depth insight for the prevailing issues. Due to the limited resources this study was conducted in a limited environment. Though the results prove impressive observations, it is expected, and encouraged the other scholars to conduct further researches based on the findings of this study to formulate a global knowledgebase on outsourcing and impact to its stakeholders due to various parameters. For that level of generalization, essentially there should more similar researches have to be conducted within other major countries within the global software outsourcing industry.

This study was initiated to examine an unaddressed issue in the outsourcing industry. Though the study does not cover each and every issue experienced by the developers, it provides sufficient results and analyses to understand the gravity of the issue faced by developers engaged in outsourced industry. It is high time for policy makers to do further research on the identified issues and formulate standardization mechanisms to strengthen outsource industry in many developing countries while facilitating the developer community. In fact, the new computational paradigms such as Cloud Computing can cause significant impact towards reforming how outsourcing works at present. It is possible to have influential impact on present software outsourcing countries, when emerging solutions offering outsourcing both hardware and software [20]. It is therefore, a must to investigate possible improvements with new technologies to enhance the outsourcing industry, both individual country level and global level.

With all this regard it is fair to believe that, this study will create a significant paradigm shift towards rethink of the outsourced software development process. In conclusion, it is evident that this research is one of the significant achievements in the present Software Engineering field. This study's outcome will direct the policy implications to benefits the stakeholders of the software outsource industry, thus assist to improve the social, economical and health levels of software developers while avoiding crisis situation within the industry, in the long run.

Acknowledgement

Author would like to thank the people who helped for this study in various forms, especially, the developers who participated in this study for their effort and contribution to make this a success.

Bibliography

- R. E. Ahmed, Software Maintenance Outsourcing: Issues and Strategies. Computers and Electrical Engineering, Vol. 32(6), pp. 449-453, 2006.
- [2] G.I.U.S. Perera, M.S.D. Fernando, Enhanced Agile Software Development Hybrid Paradigm with LEAN Practice, in Proc. of IEEE 2nd ICIIS conference, Peradeniya, pp.239-244, 2007
- [3] A. Fuggetta, Software process: a roadmap, in Proc. of the conference on The future of Software Engineering, ICSE, Limerick, p.25-34, 2000
- [4] W.S. Humphrey, Managing the Software Process, SEI, Pearson Education, India, pp.03, 2006
- [5] K.E. Wiegers, Software Requirements, 2nd Ed. Redmond, Wash, Microsoft Press, 2003
- [6] G.I.U.S. Perera, M.S.D. Fernando, Bridging the gap Business and information systems: A roadmap, in Proc. of 4th ICBM conference, pp. 334-343, 2007
- [7] B. Nuseibeh, S. Easterbrook, "Requirement Engineering: A Roadmap", Proc. of the conference on The future of Software Engineering, 22nd ICSE, Limerick, p. 35-46, 2000
- [8] P. Zave, Classification of Research Efforts in Requirements Engineering. ACM Computing Surveys, Vol. 29(4), pp. 315-321, 1997
- [9] K. Beck, et. al., Manifesto for agile software development, 2001, available at http://agilemanifesto.org/, [accessed on 07th December 2008]
- [10] T.M. Rajkumar, R.V.S. Mani, Offshore Software Development. The View from Indian Suppliers, *Information Systems Management*, pp. 63-73, 2001

- [11] E. Carmel, Taxonomy of New Software Exporting Nations, Electronic Journal of Information Systems in Developing Countries, Vol. 13(2), p. 1-6, 2003
- [12] A. Yalaho, Plugging Into Offshore Outsourcing Of Software Development: A Multiple Case Study, *Issues in Information Systems*, Vol. 8(2), pp. 499-515, 2007
- [13] N. Levina, J.W. Ross, From the Vendor's Perspective: Exploring the Value Proposition in Information Technology Outsourcing, MIS Quarterly, Vol. 27(3), pp. 331-365, 2003
- [14] C. T. Coward, Looking Beyond India: Factors that Shape the Global Outsourcing Decisions of Small and Medium Sized Companies in America, *Electronic Journal of Information Systems* in Developing Countries, Vol. 13(11), pp. 1-12, 2003
- [15] E. Carmel, R. Agarwal, Tactical Approaches for Alleviating Distance in Global Software Development, *IEEE Software*, Vol. 18(2), pp. 22-29, 2001
- [16] K. Molřkken, M. Jřrgensen, A Review of Surveys on Software Effort Estimation. In Proc. of the International Symposium on Empirical Software Engineering, pp. 223-231, 2003
- [17] I. Perera, Impact of using agile practice for student software projects in computer science education, International Journal of Education and Development using Information and Communication Technology, Vol. 5(3), Online [http://ijedict.dec.uwi.edu/viewarticle.php?id=755], 2009
- [18] G.I.U.S. Perera, Fernando M.S.D., Rapid Decision Making for Post Architectural Changes in Agile Development - A Guide to Reduce Uncertainty, *International Journal of Information Technology and Knowledge Management*, Vol. 2(2), In Press, 2009
- [19] G.I.U.S. Perera, "Key Success Factors for e-Learning Acceptability: A Case Based Analysis on Blended Learning End-User Experience", In Proc. of IEEE International Advance Computing Conference, IACC'09, pp. 2379-2384, 2009
- [20] I. Perera, Reshaping the Computation with Clouds: an Analysis on Opportunities and Issues of Cloud Computing, Journal of Advances in Computational Sciences and Technology, 2(3), pp.305-324, 2009

Accessing Information Sources using Ontologies

D. Sun, H. Jung, C. Hwang, H. Kim, S. Park

Dongeun Sun, Hyosook Jung, Chunsik Hwang, Heejin Kim

Department of Computer Science Education Korea University, Anam-dong Seongbuk-gu, Seoul, 136-701, Korea E-mail: sunde41@korea.ac.kr, est0718@comedu.korea.ac.kr, vollfeed@korea.ac.kr, prin@korea.ac.kr

Seongbin Park

Department of Computer Science Education Korea University, Anam-dong Seongbuk-gu, Seoul, 136-701, Korea E-mail: hyperspace@korea.ac.kr Corresponding author

Abstract: In this paper, we present a system that helps users access various types of information sources using ontologies. An ontology consists of a set of concepts and their relationships in a domain of interests. The system analyzes an ontology provided by a user so that the user can search and browse Wikipedia [1], DBpedia [4], PubMed [5], and the Web by utilizing the information in the ontology. In particular, terms defined in the ontology are mapped to Wikipedia pages and the navigation history of a user is saved so that it can serve as a personalized ontology. In addition, users can create and edit ontologies using the proposed system. We show that the proposed system can be used in an educational environment.

Keywords: Wikipedia, ontology, education.

1 Introduction

As the amount of information available from various information sources such as the Web is increasing, it becomes more difficult to find relevant information in large information spaces. In this situation, an ontology that describes concepts and their relationships in a domain of itnerests can help since terms provided by ontologies can help novices within a specific domain or people who are not familiar with searching [22]. An ontology can be also utilized for effective navigation [28].

In this paper, we present an ontology-based system that provides users with appropriate keywords for searching so that users can easily access various types of information sources such as Wikipedia, DBpedia, PubMed, and the Web. The system uses an ontology to provide keywords for searching. If a user selects a concept defined in the ontology, our system shows the corresponding Wikipedia article using http://en.wikipedia.org/wiki/reference, where reference is the selected concept. In addition, the user can access DBpedia dataset and PubMed that contain the information associated with the selected concept. Users can also access the information related to the concept on the Web using Google search engine [21]. On top of these, users can share ontologies with others and the list of pages visited by users can be saved as navigational histories in RDF [6] or OWL [7]. The saved ontology can serve as a personalized ontology that can be utilized for web search personalization [29].

The left figure in Fig. 1 is the screenshot of the page at http://en.wikipedia.org/wiki/Gene which a user can see without using our system. The right figure shows the screenshot of the same web page that a user can see using the proposed system. The ontology viewer shows an ontology imported by a user and the class and individual viewer shows the hierarchy of classes

in the ontology. When a user clicks a class in the class and individual viewer, the class' detail viewer shows the subclasses and properties of the class. The navigational history viewer shows the list of the pages visited by the user.

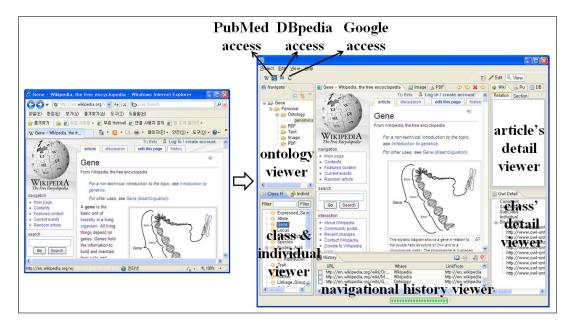


Figure 1: A snapshot of our system

An ontology is usually defined by domain experts and represents essential concepts and meaningful relationships between the concepts. If a user obtains a well-organized ontology defined by domain experts, the user can understand which concepts are suited to the searching goal and which concepts are related meaningfully. It is helpful for the user to establish search strategies. For example, because the user can know which concepts are related to one another, it is easy to determine which concepts the user should look for next. This is important especially in a Web-based educational environment where students who are not sure about how to proceed to find right information for leaning tend to get lost. This symptom is called lost in hyperspace. In such a case, teachers can provide their students with an appropriate navigational guidance by offering ontologies that define important concepts and their relationships so that the students can figure out where to proceed to achieve learning goals by referring to the contents of the ontologies. Ontologies can also help navigating on Web sites since users can click categories defined in ontologies that can be shown on a part of the Web sistes to expand categories.

There are two types of users who can potentially benefit from our system. One is a provider who can create ontologies for consumers. The other is a consumer who uses the ontologies created by the providers. In a learning environment, teachers are providers. Teachers create ontologies by themselves or modify ontologies defined by others using our system. They provide ontologies for their students. Students are consumers. Most probably, they do not know what the ontologies are or how the ontologies are created. They just use the ontologies to search or browse through learning materials. Our system targets on both novices and experts about ontologies. In a more general environment, some users can be both a provider and a consumer. For example, a user might want to define an ontology that captures background knowledge on the domain of interests and utilizes the ontology to improve the accuracy of searching. Depending on the type of a user, our system helps in creating and editing ontologies. For example, an advanced user who knows how to create an ontology can easily edit an existing ontology to personalize it. It is also possible that users can collaboratively edit ontologies. This paper is organized as follows. Section 2 describes related works. In section 3, the functionalities of the system are explained. Section 4 describes illustrative examples that show how the system can be used. Section 5 concludes the paper and describes future works.

2 Related works

Ontology is referred to as the explicit and formal specification of a shared conceptualization. As an engineering artifact, it consists of terms and relationships that describe a certain reality, plus a set of explicit assumptions regarding the intended meaning of the vocabulary [8]. An ontology can serve as background knowledge [24] and it can help users refine the search results from domains that they are not familiar with [25]. Constructing a formal ontology generally relies on an interactive process to elicit knowledge and formalize it. The construction of scientific ontologies needs tools required in ontology designers and domain experts collaborate and combine their knowledge [20].

Wikipedia is a freely available online encyclopedia developed by a community of users to which anyone can contribute [1]. Wikipedia supports information collection in e-Science by using Web 2.0 technology. The Gene Wikipedia is a knowledge collection about genes obtained by gathering information from multiple data sources based on a set of ontologies [2]. While Wikipedia is a valuable source of information in many areas, search capabilities are limited to full-text search [3]. [17] sketches a system that helps managing knowledge using Wikipedia. There are browsers that help users access Wikipedia. Pathway is a browser that searches for Wikipedia articles related to a subject and displays the related articles as a graph or network [11]. Gollum allows a user to access Wikipedia through the toolbars and menus and allows the user to add dedicated Wikipedia Bookmarks in a sidebar [12]. Indywiki [13] supports browsing through images by extracting all images on Wikipedia articles related to a search keyword. [18] describes a Wikipedia browsing system using Semantic Web technologies.

Sealife project supports context-based information integration through the use of semantic hyperlinks which map the ontology terms to the web/grid services. It allows users to add the relevant terms to a shopping cart while the users are browsing through various web pages. Then, it presents the contents of the shopping cart and then enables the users to link to the web/grid services according to semantically identified terms [24]. Our system uses a domain specific ontology as background knowledge and links ontology terms to relevant Wikipedia articles. It enables browsing semantically related articles based on the ontology. It also allows users to build a new ontology by using their browsing history that consists of terms automatically extracted from the articles visited by the users. The users can use the ontology in order to browse articles based on user's interest.

Magpie [10] uses an ontology to annotate web pages and provides meaningful information related to the web pages. When a user visits a web page, it annotates semantic entities on a web page based on the ontology. Our system uses an ontology to support the navigation of semantically related web pages within a specific domain. When a user selects a semantic entity, it shows a web page related to the entity and provides supplementary information related to the entity and the web page. In addition, our system allows the users to modify the ontology suited to their needs or interest. It means that it can support the creation and utilization of the personalized ontology.

[15] presents a system that can generate an ontology based on user's seed. A seed is composed of a single sample data instance embedded within the concepts of interest. In our system, a user selects a class in an ontology and the system regards the class as a concept which the user is interested in. The user continues navigating information of interest with following link anchors in the article or selecting other classes in the ontology. Our system can create a personalized ontology that consists of information about a history of user's navigation. For example, a class that is selected by a user at first can be considered as a seed and an ontology about the navigational history can be viewed as a generated ontology based on the seed.

[23] presents a system that helps students navigate in the web. It connects each page to one concept of the ontology and creates new links between the pages based on the relationships between the concepts. It can offer semantically related links by using the ontology although there are not real links between the pages. In its experiment, the students who navigated the web pages including semantic links experienced less disorientation, revisited web pages less, spent less time completing tasks, and have done better tasks than those who navigated the general web pages. It means that using ontology can help students quickly and easily find what they want on the web.

[14] is an ontology-based biomedical literature search engine. It retrieves relevant abstracts from PubMed which contain the biomedical concepts related to a search query. The biomedical concepts belong to the Gene Ontology (GO), Medical Subject Headings (MeSH) and Universal Protein Resource (UniProt). In GoPubMed, users can navigate relevant abstracts by selecting one of categorized ontological concepts. GoPubMed supports quick navigation through the abstracts by category and so users can jump to search results related for a certain concept.

Ontoverse [9] creates and links ontological concepts and the information extracted from scientific publications. Ontoverse enables annotating and interlinking knowledge sources and supports cooperative knowledge management within life science based on the ontologies. Our system uses the ontologies as knowledge basis for semantic information access and navigation. Our system obtains ontological concepts from user-defined or existing ontologies and also extracts conceptual information from Wikipedia articles from the history of user's navigation. It is possible to create ontologies collaboratively because it allows users to edit existing ontologies and build new ontologies according to their interest.

When a user navigate in a complex hyperspace, the user can feel the symptom of lost in hyperspace. Our system helps users navigate in a large hyperspace such as the Web using ontologies so that they can find right kind of information easily without being distracted. This is important in an area such as Web-based education, where one type of users is students who need to understand learning materials on the Web. They could have different backgrounds and ontology-based guidance can help them access necessary learning materials. Our system allows teachers to create ontologies suited to their students. It presents the students the ontological concepts related to information to be learned. The students navigate the web according to relationships between the ontological concepts. The system can prevent the students from getting lost in hyperspace or cognitive overload while browsing web pages. In addition, our system can provide adaptation by cooperating with an adaptive hypermedia system. It offers users additional information suited to user's characteristics such as knowledge, interest, etc. In the learning environment, our system can provide students the learning content suited to their learning ability.

Fig. 2 shows how our research is related to other works. Our system uses an ontology to search for semantic information such as Sealife and Magpie and to guide navigation in hyperspace semantically like ontology-based web navigation guidance. It also generates personalized ontology to retrieve information suited to user's need or interest like seed-based ontology generation. Our system enables users to search Wikipedia articles such as Pathway, Gollum and Indywiki and abstracts of biomedical literatures form PubMed like GoPubMed. It also supports searching information from Google and DBpedia. Like Ontoverse, our system supports cooperative ontology building.

The differences between our system and systems mentioned in this section are as follows. Sealife and Magpie are semantic web browsers to identify ontological terms from web pages

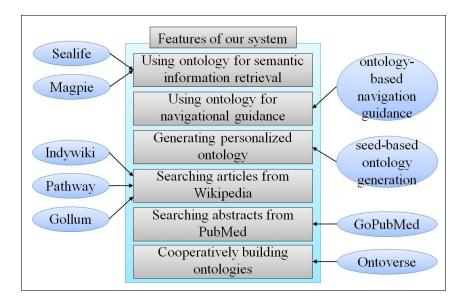


Figure 2: Related works to our research

and enable users navigate the web according to the relationships between the identified terms. The ontology-based navigation guidance lets users navigate web pages which are linked according to the relationships between concepts in ontologies. The seed-based ontology generation and Ontoverse are tools which generate personalized ontology consisting of the concepts of user's interest and gathers information based on the ontology. In order to use these tools, the uses should know ontologies to be used. However, in our system, even though the users like students do not know them, students can search for information by using ontologies defined by experts or teachers.

Pathway, Indywiki, and Gollum are web browsers to support navigating Wikipedia. Pathway keeps track of what users have read on Wikipedia and displays the track in a data map. Indywiki extracts images from related pages on Wikipedia and displays the page that the image links to. Gollum gives users special functions which support working with Wikipedia easily. When using these tools, the use should choose proper queries to find for information. However, our system helps the users to determine the queries as showing the concepts related to the information to be searched.

GoPubMed is a search engine to search for abstract of biomedical literatures from PubMed, GO, MeSH, and UniPort. While it is too professional for novices about biology to retrieve proper information, out system helps the users search for information about different domains as well as biology by using ontologies customized according to the user's features.

3 Functionalities of the system

In this section, we explain the functionalities provided by our system. The system parses an ontology and shows the classes and individuals of the ontology in a tree structure. A user can use one's own ontology or an ontology found from Swoogle which is a Semantic Web search engine [16]. Users can add new classes or individuals to the ontology or remove the classes or individuals from the ontology. If a user selects a class, it shows the subclasses, properties, or individual of the class in a tree structure. The system also displays the corresponding Wikipedia article. In addition, it provides concept categories which the article belongs to and incoming links to the article.

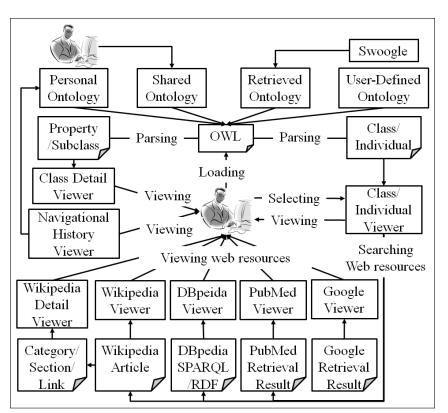


Fig. 3 depicts the functionalities of the proposed system.

Figure 3: Functionalities of the proposed system

A user can save text data, image, PDF files into project, which is the data collection including ontology, text, image, PDF in our system. Besides, for saved data in project, the system has its own viewing function of TEXT, image (jpg, gif, PNG format), PDF. The system supports OWL editing such as syntax coloring, formatting, highlighting, content outline, word wrap and so on. A user can export information about the Wikipedia articles visited by the user in an ontology. The representative concept of accessed articles is automatically created as individuals of the ontology.

3.1 Ontology-based browsing

A user can retrieve information from DBpedia by selecting a class and clicking a DBpedia button. Fig. 4 shows an RDF document corresponding to the selected class from DBpedia.

When double-clicking a class 'protein' in DBpedia viewer, the user can see the HTML page which the RDF document of 'protein' is transformed into. If clicking an icon on the bottom of the page, the user can visit a web page to use SPARQL query for searching DBpedia dataset. If the user types a query in the text box like "select distinct? Protein where [] a ?Protein" and clicks a button 'Run Query', the user can see information related to the query such as data related to 'protein' (Fig. 5)

It is also possible that the user can access information from PubMed which is an archive of biomedical and life sciences journals. It provides a list of the publications related to a selected class or individual and shows information about a selected publication such as authors, abstract, title, etc.

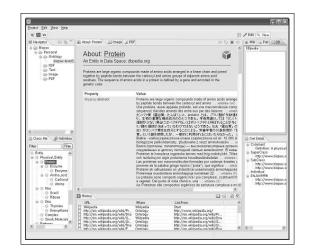


Figure 4: A result from DBpedia

Broject Edit View Help						
W 🛃 W.				E .	/ Edit 🔍 View	
	Mp://dbpedia.org/sparqi?default-g	raph-uri=http%34%	37%3Fdbpedia.org&should-spong 🖉 🔒 Im	iage 🔔 PDF	i 🚳 Wiki 🔒 Pu 🔛 DB	
B Biology B Dersonal				> % X 🗇	DBpedia	
B 😁 Ontology	Protein					
biopax-level:						
a 🍅 Text	http://www.w3.org/2004/02/skos/core#Concept http://dboedia.org/ontology/MusicalWork					
- 😓 Image - 🛵 POF	http://dopedia.org/ontology/MusicalWork					
	http://dopedia.org/ontology/Work	ice				
	http://dopedia.org/ontology/Wonc					
	http://dopedia.org/ontology/Musical http://unibel.org/unibel/sc/Product					
	http://umbel.org/umbel.ac/Artifact				a Owl Detail	
	http://dopedia.org/class/yago/1982Novels					
	http://dopedia.org/ontolog//Book					
< > >	http://dopedia.org/class/yago/EnglishAstronomers					
Class He de Individu	http://dopedia.org/class/yago/Engl					
	http://utipedia.org/unite/sc/Astronomer				- Comment	
Filter Filter	http://dopedia.org/class/vago/Livin				Definition: A phys	
Ently A Physical_Ently	http://dopedia.org/class/yago/AmericanCompetitiveEaters				 SuperClass http://www.biops SubClass http://www.biops 	
B O Protein	http://umbel.org/umbel/sc/Person					
G Enzyme G Enzyme	http://umbel.org/umbel/sc/Writing				htp://www.biop	
Amino_acic	http://unbel.org/unbel/sc TextualMaterial				 DisJointWith http://www.blopa http://www.blopa 	
Garboxs Garboxs	http://unibel.org/umbel/sc/Story					
- Amine	http://dopedia.org/ontology/Place			-		
	History					
	URL	Where	Linkfrom			
	Wikpedia http://en.wikipedia.org/wiki/Pr http://dbpedia.org/page/Protein http://dbpedia.org/page/Protein http://dbpedia.org/page/Protein http://dbpedia.org/page/Protein	Wikipedia Ontology DBpedia DBpedia DBpedia DBpedia	Start http://www.wkipadia.org/ http://en.wkipadia.org/wki/Pr http://dbpedia.org/sage/Protein http://dbpedia.org/spard http://dbpedia.org/spard			

Figure 5: SPARQL query result

3.2 Collaborative ontology building

Users can write an ontology collaboratively with other users while sharing their opinions and knowledge for their common goal. For example, an expert accesses a Wikipedia article about a class Protein. While reading the article, he wants to know information about Amino acid. He continues navigating the article about Amino acid. Then, he edits the ontology by adding a subclass Amino_acid of the class Protein. He also adds subclasses Carboxyl_group and Anime of the class Amino_acid. He shares the ontology with other experts. (Fig. 6)

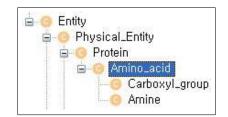


Figure 6: Editing a part of a shared ontology

Now, another expert interested in Enzyme downloads the edited ontology. She follows a link anchor Enzyme while reading the Wikipedia article about Protein and continues reading the article about Enzyme. Then, she edits the ontology by adding a subclass Enzyme of the class Protein and a subclass Enzyme_catalysis of the class Enzyme. (Fig. 7)

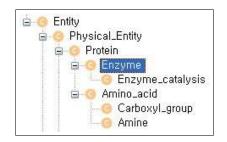


Figure 7: Editing parts of the shared ontology

3.3 Creation of a navigation history

Our system exports the navigation history that consists of the information about pages visited by a user in OWL and it can serve as a personalized ontology which can be used for semantic searching. For example, a personalized ontology can be created as follows. (Fig. 8)

- 1. A user is interested in water activities and loads an ontology about travel (travel.owl).
- 2. The user clicks individual kayaking of class Water_activity for requesting a Wikipedia article about kayaking.
- 3. To return the Wikipedia article corresponding to the class, the system creates a URL string including a class name as a query like http://en.wikipedia.org/wiki/Name, where Name is a concept or term to search for the article. Our system creates a URL string including the individual name like "http://en.wikipedia.org/wiki/Kayaking".
- 4. Wikipedia sends the article Kayaking and our system shows the article to the user.

- 5. While navigating in Wikipedia, the user is interested in canoeing which is similar kayaking. The user adds canoeing as a new individual of class Water_activity to the ontology.
- 6. The user can see the article Canoeing by clicking the individual canoeing.
- 7. The user can obtain RDF data of DBpedia about canoeing.
- 8. The user can also obtain searching results of PubMed about canoeing.
- 9. The user continues navigating other interesting Wikipedia articles such as Niagara_Falls, Horseshoe_Falls, Whitewater_canoeing, etc. Then, the user builds an ontology which consists of the navigational history of the user. In the ontology, the visited articles are represented as individuals of class Article.

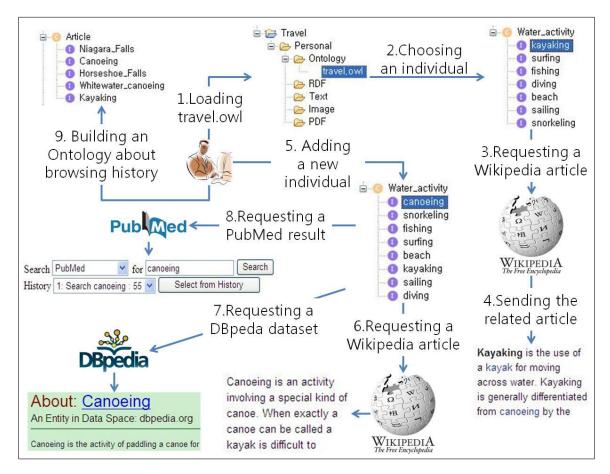


Figure 8: Creation of a personalized ontology

Many ontologies are defined by domain experts. While one important reason to use ontologies is to support knowledge sharing for different types of applications, ontologies can be also re-used and customized for personal use. Especially, since the number of ontologies on the Web is everincreasing and there are search engines for ontologies, average users can find ontologies for the domain of interests under consideration easily. If the users can edit ontologies found on the Web using either our system or other ontology editors, they can customize them so that the modified ontologies can be used for interesting applications. One area where this can help is Web-based education where teachers can prepare a set of concepts or terminologies that need to be accessed by students on the Web. In this situation, teachers serve as so-called trail blazers who provide helpful information for someone else. It is also possible that personalized ontologies can help learning in adaptive hypermedia area since different types of adaptations can be made possible using the information contained personalized ontologies. In other words, a base ontology exists and there also exist different ontologies that are modified slightly from the base ontology in order to support adaptive learning.

3.4 An illustrative example

In this section, we show navigational efficiency of our system by comparing navigation using our system and navigation using a general web browser.

We believe that users can access information from different websites faster by using our system than by visiting each web site manually. In our system, they can search information from different websites in an application. They do not need to load web pages in a web browser, type the URLs of each web site, and enter queries. Instead, they can navigate the web sites by clicking terms defined in ontologies. As an illustrative example, let's assume that a teacher wants students to find out information about human nutrition source in Wikipdia. Specifically, the students should find out Wikipedia articles about "Photosynthesis", "Herbivore", and "Carnivore" but the problem is that students do not know the terms exactly. So, a typical way of finding the information using a web browser would be that a student first visits Wikipedia and types query "nutrition" in the search box. The student reads the Nutrition article and clicks link anchor "food" which is shown in the left figure in figure 9. Then the student reads the Food article and clicks "plant" which is shown in the right figure in figure 9.



Figure 9: Search results using a general web browser

After that, the student reads the Plant article and clicks "photosynthesis" and eventually finds the Photosynthesis article that is shown in the left figure in figure 10. Then, the student goes back to the Food article and clicks "animal" which is shown in the right figure in figure 10.

Now, the student reads the animal article and clicks "carnivore" and "herbivore". Then, student eventually finds the Herbivore article and the Carnivore article as shown left figure in figure 11 and right figure in figure 11, respectively.

In our system, a student can open an ontology defined by the teacher. The student clicks "Photosynthesis" on the class view and then obtains the Photosynthesis article immediately which is shown in figure 12.

Then, the student also clicks "Herbivore", and "Carnivore" on the class view and then obtains the Herbivore article and the Carnivore article which is shown in left figure in figure 13 and right

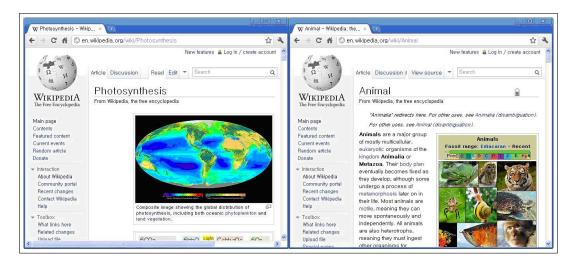


Figure 10: Search results using a general web browser



Figure 11: Search results using a general web browser

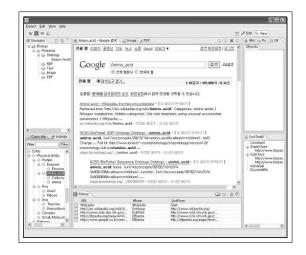


Figure 12: Screenshot of our system

figure in figure 13, respectively.

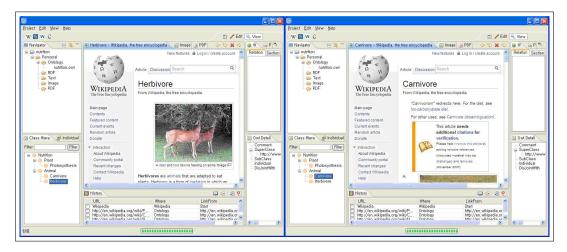


Figure 13: Search results using our system

4 Applications

In this section, we show how the proposed system can be used in an educational environment.

4.1 Management of studying materials on the Web

Assume that a student is studying the history of the Middle Ages by browsing resources on the web. The student wants to organize information by concepts related to events or developments of the middle ages. Our system can help the student manage the information conveniently.

Fig. 14 is the screenshot of the system, where the student loads an ontology related to the history of the Middle Ages and their relationships. For example, it defines classes such as war, health, religion, building, etc. It defines subclasses of each class like the subclasses of class building are castle, dungeon, walled town, moat, etc. If the student double clicks a class, the student can see a Wikipedia article corresponding to the class. The student reads the article and saves a part of the article inside the system.

If the student wants to search health information related to the Middle Ages such as the Bubonic plague and herbal cure, the student selects a class the bubonic plague and clicks an icon "PubMed". The student can see a list of the publications retrieved from PubMed. (Fig. 15)

The student can also search the RDF documents related to the Wikipedia articles from DBpedia. (Left figure in Fig. 16) In addition, the student can see a list of retrieval results from Google. (Right figure in Fig. 16)

The navigational history is recorded while the student browses web pages. Like bookmarks, the student can save the information about visited web pages in OWL or RDF which can serve as a personalized ontology.

4.2 Utilization of a navigation history

Assume that a teacher plans to teach about cell. In the class, the teacher is planning to make the students search useful resources related to the lesson in the web after teaching basic concepts. However, it is hard for them to search or define appropriate ontologies. So, the teacher



Figure 14: Initial screenshot

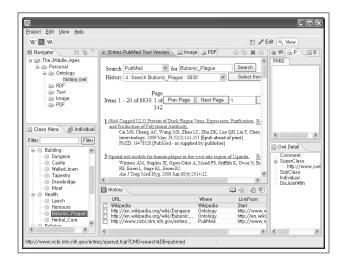


Figure 15: A result from PubMed

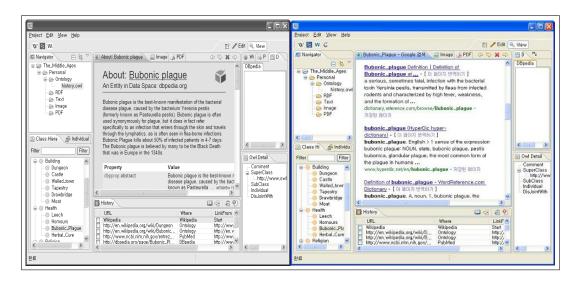


Figure 16: A result from DBpedia (left figure) and a result from Google (right figure)

needs to provide the students with an ontology which contains basic concepts. Using our system, the teacher browses Wikipedia articles and occasionally the teacher selects a certain class of the ontology and browses its Wikipedia article. When reading the article, the teacher can find essential concepts even though they are not defined in the basic ontology vet. By clicking a certain word in the article, the teacher can browse the Wikipedia articles related to the word. The teacher can save the navigational history as an OWL document where the articles are saved as the individuals of Article class. For example, when reading Mitosis article, the teacher also reads articles about DNA, RNA and virus. Although the basic ontology does not contain the concepts, the teacher wants to make the students read the articles. So the teacher saves navigational history about DNA, RNA and virus article. Using the system, the students can load the ontology of navigational history where DNA, RNA and virus are saved as the individuals of Article class when the teacher saves the navigational history. They select a certain class of the ontology and request information about the concept from Wikipedia, Google, PubMed and DBpedia. A student selects Virus class and asks virus article from Wikipedia. (Left figure in Fig. 17) The student wants to find out news, blogs, and images related to virus. So, the student asks information from Google after selecting Virus class. (Right figure in Fig. 17)

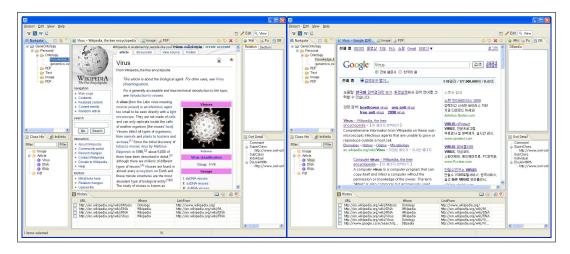


Figure 17: Wikipedia article (left figure) and Google searching results (right figure)

If the student wants to find journals related to virus, the student can request PubMed to search journals after selecting Virus class using our system. The student can also obtain structured data of Wikipedia article (i.e., an RDF document) about virus from DBpedia. On top of these, students can find semantically-related document to the ontology on the web using a semantic search engine [19] where the ontology can serve as searching context. Fig. 18 is the screenshot of the semantic search engine, where individuals of the loaded ontology are shown. The semantic search engine extracts terms from the ontology that are its individual names such as DNA, RNA, Mitochondria, and Ribosome. They are used as keywords for searching. For example, [19] can generate a query string from these keywords and find Web documents or Semantic Web documents using the query string.

4.3 Interacting with an adaptive hypermedia system

Our system can interact with AHA! which is an adaptive hypermedia system [27] that has been developed in order to offer links and content suitable to each user so that the user cannot be disoriented during navigation in a hyperspace. AHA! server consists of domain model, user model and adaptation model. Domain model describes the structure of domain and user model

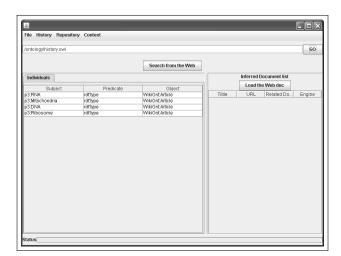


Figure 18: Individual information

contains the information about users. Adaptation model describes how adaptation can be made depending on a situation. Our system can use an ontology that consists of concepts of the domain model. The classes of the ontology are the concepts of the domain model and the hierarchy of the classes is the same as that of the concepts. When a user accesses a web page in AHA! server, the user can also access additional information related to the web page which is automatically provided by our system.

For example, an author creates a tutorial about markup language in AHA! server. The domain model consists of concepts about markup language such as HTML, XHTML, XML, etc. Each concept has a corresponding web page. In our system, there is an ontology that consists of classes corresponding to the concepts. A user is accessing a web page about xhtml in AHA! server that is one of the concepts in domain model. (Fig. 19)

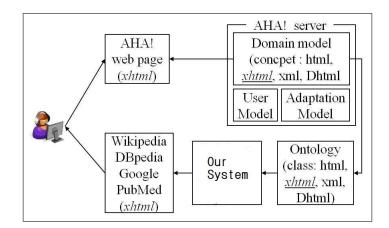


Figure 19: Interaction between our system and AHA! system

The user wants to read additional information related to xhtml. So, the user clicks "additional information" on the top of the web page that is available from the AHA! server. Clicking "additional information" makes our system show a Wikipedia article about XHTML by loading the ontology that consists of the concepts of the current domain model and receiving the concept xhtml which the user is accessing to. (Fig. 20)

Similarly, the user can also access the information from DBpeida, PubMed, and Google.

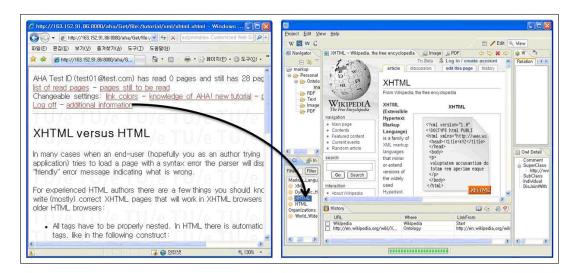


Figure 20: Accessing additional information from Wikipedia

5 Conclusion and Future works

In this paper, we presented an ontology-based system that can help users access various types of information sources. Using our system, users can browse Wikipedia articles, Google searching results, PubMed journals, and DBpedia data about a certain concept defined in an ontology in a single user interface. They can also find out other concepts semantically related to the concept by traversing ontological hierarchies between the concepts. Our system provides an environment for collaborative ontology engineering that enables users to create, edit and maintain ontologies and to reuse the ontologies created by other users. The users build new ontology projects and open ontologies created by others and edit them. In addition, they store information about articles navigated on the Wikipedia and ontologies are automatically created.

The contributions of the paper are as follows. First, the proposed system supports the integrated and uniform way of searching and browsing information contained in different data sources using ontologies. Second, users can create and maintain personalized navigation histories that can be utilized for searching. Third, the system allows users to build ontologies collaboratively.

We are currently working on ways by which information about metadata and ontologies can be utilized in our system as [26].

Bibliography

- [1] Wikipedia, 2009. Available from http://en.wikipedia.org/wiki/WikiPedia.
- [2] Gene Wikipedia, 2006. Available from http://www.bioinformatics.org/genewiki/wiki.
- [3] Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., Ives, Z., 2007. DBpedia: A Nucleus for a Web of Open Data. Proc. 6th Int. Semantic Web Conf. and 2nd Asian Semantic Web Conf. (ISWC+ASWC 2007), Busan, Korea. p.722-735.
- [4] DBpedia, 2009. Available from http://dbpedia.org.
- [5] PubMed, 2009. Available from http://www.ncbi.nlm.nih.gov/pubmed.
- [6] RDF, 2004. Available from http://www.w3.org/RDF.

- [7] OWL, 2004. Available from http://www.w3.org/TR/owl-ref.
- [8] Ontology, 2002. Available from http://km.aifb.uni-karlsruhe.de/projects/owl/index.html
- [9] Ontoverse, 2005. Available from http://www.ontoverse.org.
- [10] Magpie, 2005. Available from http://projects.kmi.open.ac.uk/magpie/main.html.
- [11] Pathway, 2007. Available from http://pathway.screenager.be.
- [12] Gollum, 2009. Available from http://gollum.easycp.de/en.
- [13] Indywiki, 2008. Available from http://indywiki.sourceforge.net.
- [14] GoPubMed, 2009. Available from http://www.gopubmed.com.
- [15] Tao, G., Embley, D. W., 2007. Seed-Based Generation of Personalized Bio-ontologies for Information Extraction. In Proceedings of the Advances in Conceptual Modeling - Foundations and Applications, p.74-84.
- [16] Swoogle,2007. Available from http://swoogle.umbc.edu.
- [17] Sun, D., Park, S., Jung, H., 2009. Knowledge management using Wikipedia. Semantic Web Applications and Tools for Life Sciences, Proc. CEUR, Vol. 435, Edinburgh, United Kingdom.
- [18] Sun, D., 2009. Design and Implementation of a Wikipedia Browser using Semantic Web Technologies, master's thesis, Korea University, Korea.
- [19] Ahn, J., Jung, H., Kim, H., Sun, D., Park, S., 2008. A system for contextual search. Proc. IEEE Int. Workshop on Semantic Computing and Applications, p.96-98.
- [20] Mainz, I., Weller, K., Paulsen, I., Mainz, D., Kohl, J., von Haeseler, A., 2008. Ontoverse: Collaborative Ontology Engineering for the Life Sciences. Informations-Wissenschaft und Praxis, Vol. 2, p.91-99.
- [21] Google, 2009. Available from http://www.google.com.
- [22] Groth, K., Lannerö, P., 2006. Context browser: ontology based navigation in information spaces. Proc. 1st Int. Conf. on information interaction in Context, IIiX: Vol. 176, p.75-78.
- [23] Jung, H., Kim, H., Min, K., Park, S., 2009. The Ontology-based Web Navigation Guidance System. The Journal of Korean Association of Computer Education, Vol. 12(5).
- [24] Schroeder, M., Burger, A., Kostkova, P., Stevens, R., Habermann, B., Dieng-Kunts, R., 2006. Sealife: A Semantic Grid Browser for the Life Sciences Applied to the Study of Infectious Diseases. Vol. 120, p.167-178.
- [25] Bonomi, A., Mosca, A., Palmonari, M., Vizzari, G., 2008. Integrating a Wiki in an Ontology Driven Web Site: Approach, Architecture and Application in the Archaeological Domain. 3rd Semantic Wiki Workshop.
- [26] Bhagdev, R., Chapman, S., Ciravegna, F., Lanfranchi, V., Petrelli, D., 2008. Hybrid Search: Effectively Combining Keywords and Semantic Searches. Springer Berlin / Heidelberg, p.554-568.
- [27] De Bra, P., Brusilovsky, P., Houben, G., 1999. Adaptive Hypermedia: From Systems to Framework. ACM Computing Surveys, Vol. 31(4).

- [28] Velart, Z., Saloun, P., 2007. Ontology Based Course Navigation. Proc. 18th Conf. on Hypertext and Hypermedia, p.151 152.
- [29] Sendhilkumar, S., Geetha, T. V., 2008. Personalized Ontology for Web Search Personalization. Annual Bangalore Compute Conference. Proc. 1st Bangalore Annual Compute Conf., p.1-7.

Digital Control of a Waste Water Treatment Plant

R. Vilanova, J.D. Rojas, V.M. Alfaro

Ramón Vilanova, José David Rojas

Department de Telecomunicació i Enginyeria de Sistemes Universitat Autònoma de Barcelona 08193, Bellaterra, Spain, E-mail: ramon.vilanova@uab.cat, josedavid.rojas@uab.cat

Víctor M. Alfaro

Escuela de Ingeniería Eléctrica Universidad de Costa Rica San José, 11501-2060 Costa Rica. E-mail: victor.alfaro@ucr.ac.cr

Abstract: The Activated Sludge Process (ASP) is arguably the most popular bioprocess utilized in the treatment of polluted water. The ASP is described by means of a nonlinear model and results on a Two-Input Two-Output multivariable system. In this paper a discrete time digital control is proposed where the design of a decentralized controller is faced. Local controllers are given the form of a Two-Degree-of-Freedom PI controller tuned using the data-driven Virtual-Reference Feedback tuning approach.

Keywords: ASP, Process Control, PID, Data-Driven Control, VRFT.

1 Introduction

Water pollution represents one of the most serious environmental problems due to the discharge of nutrients into receiving waters. Hence, stricter standards for the operation of wastewater treatment plants (WWTPs) have been imposed by authorities. In order to meet these standards, improved control of WWTPs is needed. Wastewater treatment control has begun a gradual progress towards the use of more advanced technology, in the face of more stringent modern water quality standards. Several approaches have been reported in the literature that attempt to control the WWTPs process. Among others, the Activated Sludge Process (ASP) is arguably the most popular bioprocess utilized in the treatment of polluted water, using microorganisms present within the treatment plant in the biological oxidation of the wastewater. The simplified but still realistic and highly non-linear four-state multivariable model considered here is the ASP as presented in [1].

The ASP is described by means of a nonlinear model and results on a Two-Input Two-Output (TITO) multivariable system. In this paper a discrete time digital control is proposed where the design of a decentralized controller is faced. The paper designs the local controllers as discrete time Proportional-Integral (PI) controllers. The controllers are synthesized using the Virtual Reference Feedback Tuning, which is a model-free based approach where just purely data taken from the system is considered, therefore there is no need for a mathematical model of the system. On the basis of these data the discrete time PI controllers are tuned.

2 The two-degree-of-freedom Virtual Reference Feedback Tuning

The Virtual Reference Feedback Tuning (VRFT) is a one-shot data-based method for the design of feedback controllers. The original idea was presented in [2], and then formalized by

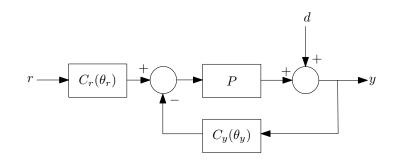


Figure 1: Two degrees of freedom structure

Lecchini, Campi and Savaresi (see [3–5]). In this section, an outline of the two-degree-of-freedom case is presented. The design methodology is presented in [5], the control structure is presented in Fig. 1. The objective of this method is to minimize the criterion in (1).

$$J_{MR}(\theta_r, \theta_y) = \|(\Psi_M(z; [\theta_r, \theta_y]) - M(z))W_M(z)\|_2^2 + \|(\Psi_S(z; \theta_y) - S(z))W_s(z)\|_2^2$$
(1)

with

$$\Psi_M(z; [\theta_r, \theta_y]) = \frac{P(z)C_r(z; \theta_r)}{1 + P(z)C_y(z; \theta_y)} \qquad \Psi_S(z; \theta_y) = \frac{1}{1 + P(z)C_y(z; \theta_y)}$$
(2)

and M(z) being the target input-to-output transfer function and S(z) the target sensitivity function. In the VRFT framework a plant model is not available, and is not intended to find one. Instead, a batch of input/output data is taken from an experiment on the plant (namely input u(t) and output y(t)). So, in order to find the parameters of the controllers (θ_r and θ_y) the signals $\bar{r}(t)$, $\bar{d}(t)$ and $\bar{y}(t)$ are defined. These signals are called "virtual" because they are not really measured, but constructed from the input/output data available and the desired closed-loop relations as follows:

- $\bar{r}(t)$ is the virtual reference, so that $y(t) = M(z)\bar{r}(t)$
- $\bar{d}(t)$ is the virtual perturbation, so that $y(t) + \bar{d}(t) = S(z)\bar{d}(t)$
- $\bar{y}(t)$ is the virtual-perturbed output of the plant, so that $\bar{y}(t) = y(t) + \bar{d}(t)$

This signals are the ones that would be found if u(t) and y(t) had been measured in closed-loop and if the closed-loop dynamics were given by M and S i.e., if the perfect controllers were set in the loop. On the basis of these "virtual" signals the controller's parameters are found by minimizing the following alternative identification cost function:

$$J_{VR}^{N}(\theta_{r},\theta_{y}) = \frac{1}{N} \sum_{t=1}^{N} \left[\Gamma_{M}(t;[\theta_{r},\theta_{y}]) \right]^{2} + \frac{1}{N} \sum_{t=1}^{N} \left[\Gamma_{S}(t;[\theta_{r},\theta_{y}]) \right]^{2}$$

where

$$\Gamma_M(t; [\theta_r, \theta_y]) = L_M(z)(u(t) - C_r(z; \theta_r)\bar{r}(t) + C_y(z; \theta_y)y(t))$$
(3)

$$\Gamma_S(t; [\theta_r, \theta_y]) = L_S(z)(u(t) + C_y(z; \theta_y)\bar{y}(t))$$
(4)

and $L_M(z)$ and $L_S(z)$ are appropriate filters to be chosen so (3) becomes an approximation to (1). If the controllers are linear in the parameter $(C_r(z;\theta_r) = \beta_r(z)^T \theta_r \text{ and } C_y(z;\theta_y) = \beta_y(z)^T \theta_y)$ the cost criterion (3) becomes a standard quadratic optimization problem. In [5] the authors use the concept of "ideal controller" to derive the structure of filters L_M and L_S . The *ideal controllers* C_{r0} and C_{y0} are the ones that, if used in the control loop, would solve (1) exactly, that is

$$C_{y0} = \frac{1-S}{SP} \qquad C_{r0} = \frac{M}{SP} \tag{5}$$

When comparing (1) and (3) using the Parseval Theorem the expression of the filters L_M and L_S that must make the identification problem (3) match the control problem (1) are found to be:

$$|L_M|^2 = |M|^2 |S|^2 |W_M|^2 \frac{1}{\Phi_u} \qquad |L_S|^2 = |S-1|^2 |S|^2 |W_S|^2 \frac{1}{\Phi_u}$$
(6)

3 2-DoF PI structure for the VRFT

In order to apply the VRFT framework to the Activated Sludge Process, the structure of the controllers has to be decided before the optimization is carried out. In [6], a decentralized PI structure is used in the same plant with good results for both, reference tracking and disturbance rejection. In this paper, a discretized version of the PI controller is used as the chosen structure for the VRFT controllers. Using a two-degree-of-freedom PI as in Fig. 1, the continuous time version of the controller is:

$$C_r(s) = K_c \left(\beta + \frac{1}{T_i s}\right) \qquad C_y(s) = K_c \left(1 + \frac{1}{T_i s}\right) \tag{7}$$

When applying the bilinear transformation $s = \frac{2}{T_s} \frac{z-1}{z+1}$, the controllers are

$$C_r(z) = \frac{K_c \left(\beta + \frac{T_s}{2T_i}\right) + K_c \left(\frac{T_s}{2T_i} - \beta\right) z^{-1}}{1 - z^{-1}}$$

$$\tag{8}$$

$$C_y(z) = \frac{K_c \left(1 + \frac{T_s}{2T_i}\right) + K_c \left(\frac{T_s}{2T_i} - 1\right) z^{-1}}{1 - z^{-1}}$$
(9)

From the VRFT point of view, (8) and (9) can be seen simply as linear-in-the-parameters controllers with two parameters as follows:

$$C_r(z) = \frac{\alpha_1 + \alpha_2 z^{-1}}{1 - z^{-1}} \qquad C_y(z) = \frac{\gamma_1 + \gamma_2 z^{-1}}{1 - z^{-1}} \tag{10}$$

Since the continuous time controllers (7) have three adjustable parameters (K_c, T_i, β) , one of the four parameters in its discrete time equivalent (10) should depend of the other three. From (8) to (10) it can be found that

$$\alpha_1 + \alpha_2 = \gamma_1 + \gamma_2 = \frac{K_c T_s}{T_i} \tag{11}$$

and we have that

$$\gamma_2 = \alpha_1 + \alpha_2 - \gamma_1 \tag{12}$$

Then the discrete time controllers (10) are now

$$C_r(z) = \frac{\alpha_1 + \alpha_2 z^{-1}}{1 - z^{-1}} \qquad C_y(z) = \frac{\gamma_1 + (\alpha_1 + \alpha_2 - \gamma_1) z^{-1}}{1 - z^{-1}}$$
(13)

Once the parameters of the controllers (13) $(\alpha_1, \alpha_2, \gamma_1)$ are found, one can recuperate the PI (7) parameters using:

$$K_c = \gamma_1 - \frac{1}{2} \left(\alpha_1 + \alpha_2 \right) \tag{14}$$

$$T_i = T_s \frac{\gamma_1 - \frac{1}{2} \left(\alpha_1 + \alpha_2\right)}{\alpha_1 + \alpha_2} \tag{15}$$

$$\beta = \frac{\alpha_1 - \frac{1}{2} (\alpha_1 + \alpha_2)}{\gamma_1 - \frac{1}{2} (\alpha_1 + \alpha_2)}$$
(16)

4 Activated Sludge Process (ASP) Description

The mathematical model considered in this paper is given in [1]. The ASP process comprises an aerator tank where microorganisms act on organic matter by biodegradation, and a settler where the solids are separated from the wastewater and recycled to the aerator. The layout is shown in Fig. 2. The component balance for the substrate, biomass, recycled biomass and dissolved oxygen provide the following set of non-linear differential equations:

$$\frac{dX(t)}{dt} = \mu(t)X(t) - D(t)(1+r)X(t) - rD(t)X_r(t)$$
(17)

$$\frac{dS(t)}{dt} = -\frac{\mu(t)}{Y}X(t) - D(t)(1+r)S(t) + D(t)S_{in}$$
(18)

$$\frac{dDO(t)}{dt} = -\frac{K_o\mu(t)}{Y}X(t) - D(t)(1+r)DO(t) + K_{La}(DO_s - DO(t)) + DO(t)DO_{in}$$
(19)

$$\frac{dX_r(t)}{dt} = D(t)(1+r)X(t) - D(t)(\beta+r)X_r(t)$$
(20)

$$\mu(t) = \mu_{max} \frac{S(t)}{k_S + S(t)} \frac{DO(t)}{k_{DO} + DO(t)}$$
(21)

where X(t) - biomass, S(t) - substrate, DO(t) - dissolved oxygen, DO_s - maximum dissolved oxygen, $X_r(t)$ - recycled biomass, D(t) - dilution rate, S_in and DO_in - substrate and dissolved oxygen concentrations in the influent, Y - biomass yield factor, μ - biomass growth rate, μ_{max} - maximum specific growth rate, k_S and k_{DO} - saturation constants, $K_{La} = \alpha W$ - oxygen mass transfer coefficient, α - oxygen transfer rate, W - aeration rate, K_o - model constant, r and β - ratio of recycled and waste flow to the influent. The influent concentrations are set to $S_{in} = 200 \text{ mg/l}$ and $DO_{in} = 0.5 \text{ mg/l}$. With respect to the control problem definition, the waste water treatment process is considered under the assumption that the dissolved oxygen, DO(t), and substrate, X(t), are the controlled outputs of the plant, whereas the dilution rate, D(t), and aeration rate W(t) are the two manipulated variables. The initial conditions and kinetic parameters are taken as in [1] and [6].

5 Discrete-time VRFT tuned PI controller applied to the ASP

Using the non-linear model presented in (17) to (21), the data in Fig. 3(a) and Fig. 3(b) was collected. Using only this data without the information of the non-linear model, or any linear approximation, the parameters of the controller are calculated according to Section 2 using the PI structure specified in Section 3. The VRFT controllers are tested against two decentralized

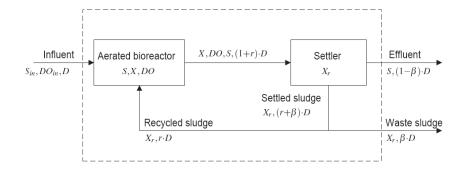


Figure 2: Activated Sludge Process layout. Taken from [7]

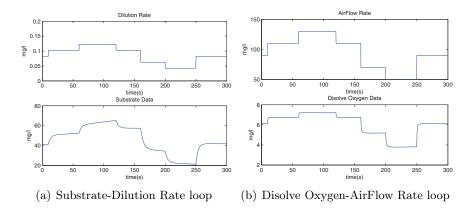


Figure 3: Data used to find the VRFT controller

PI controllers which parameters are computed using IMC [8], based on considering a First-Order (FO) controlled-process given by

$$P(s) = \frac{K_p}{Ts+1} \tag{22}$$

The identified models obey to $K_1 = 437.1$ and $T_1 = 2.7h$ for the first loop and $K_2 = 0.03$ and $T_2 = 0.51h$ for the second loop. These values are as in [6]. The controllers were discretized using the bilinear transformation (sampling time of 15min). In the case of the VRFT, the controllers are directly found in discrete time. The closed-loop specifications are given in terms of the desired time constants of the controlled system for each loop. The desired time constant of the Substrate-Dilution Rate loop is approximately $T_1 = 2.7h$ while the Dissolved Oxygen-Airflow Rate loop constant time is approximately $T_2 = 0.51h$. The closed-loop constant time for each variable is giving in terms of the dimensionless variables τ_{c1} and τ_{c2} via $T_{c1} = \tau_{c1}T_1$ and $T_{c2} = \tau_{c2}T_2$. If, for example, $\tau_{c1} < 1$, the Substrate-Dilution Rate closed-loop is expected to be faster than in open loop. It is worth to note that for the VRFT tuning it is possible to specify a different time constant for the disturbance attenuation transfer function (in fact S(z)). In this case, the VRFT tuned PI controllers are found with a time constant for the corresponding S(z) transfer function that is half the one specified for the reference to output relation by using τ_c (therefore $\tau_c/2$). If $\tau_{c1} = \tau_{c2} = 1$, the resulting controller parameters are $K_{c1} = 0.0023$, $K_{c2} = 33.33$, $T_{i1} = 2.7h$ and $T_{i2} = 0.51h$ for the IMC and $K_{c1} = 0.0042, K_{c2} = 25.41, T_{i1} = 3.43h, T_{i2} = 0.42h, \beta_1 = 0.6602$ and $\beta_2 = 0.74$ for the VRFT. The response to a change in the set points of both loops is giving in Fig. 4. The IAE value represents the Integrated Absolute Value of the Error. As it can be seen, the responses for the VRFT provide smaller IAE as well as less demanding control actuation, computed here as the Total Variation (TV) or the sum of the control movements from

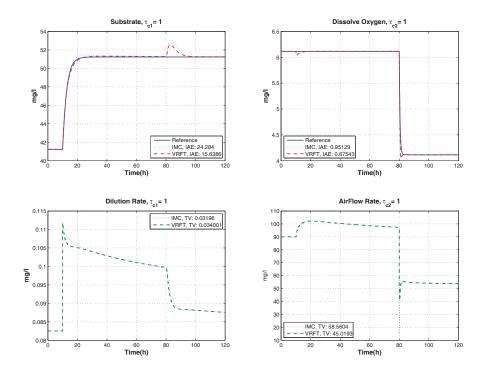


Figure 4: VRFT and IMC controllers responses to set-points step changes, $\tau_{c1} = \tau_{c2} = 1$

one sampling time to the other.

In case the closed-loop bandwidth is increased and we set $\tau_{c1} = \tau_{c2} = 0.5$ the resulting controller parameters are $K_{c1} = 0.0046$, $K_{c2} = 66.67$, $T_{i1} = 2.7h$ and $T_{i2} = 0.51h$ for the IMC and $K_{c1} = 0.0065$, $K_{c2} = 33.38$, $T_{i1} = 2.68h$, $T_{i2} = 0.35h$, $\beta_1 = 0.7752$ and $\beta_2 = 1.03$ for the VRFT. In this case the IMC controller presents and undesirable oscillating behavior in the Dissolved Oxygen loop, as shown in Fig. 5. The output of the controllers was saturated to 0, in case it went below this value. Also the control effort of the second controller is quite lower in the case of the VRFT with a performance nearly 50% better. Also a simulation was carried out for a disturbance in the inflow concentration. The results are depicted in Fig. 6. Again the results are quite similar, and the multivariable characteristic is tackled in a satisfactory way. As it can be seen, this data driven methodology is suitable for the control of the ASP process and it allows to skip the modeling step in order to find a good controller.

6 Conclusions

This paper has presented the application of a purely data based approach for tuning of discrete time PI controllers. The main advantage of the proposed method is that it does not rely on the usual linear model approximation of the system to be controlled. Just an experiment that provides input-output data from the system is needed. The performance of the tuning approach has been tested on a non-linear multivariable system and performance compared with that of the well known IMC method. As designed performance is more demanding, the resulting control system exhibits better results than its IMC counterpart.

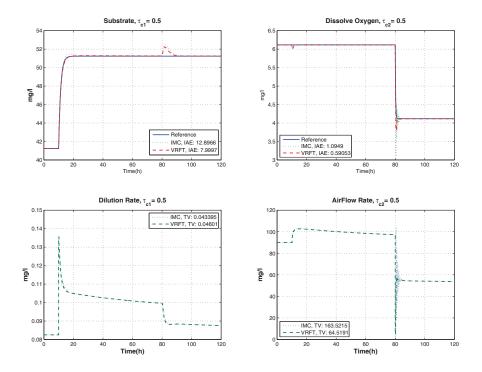


Figure 5: VRFT and IMC controllers responses to set-points step changes, $\tau_{c1} = \tau_{c2} = 0.5$

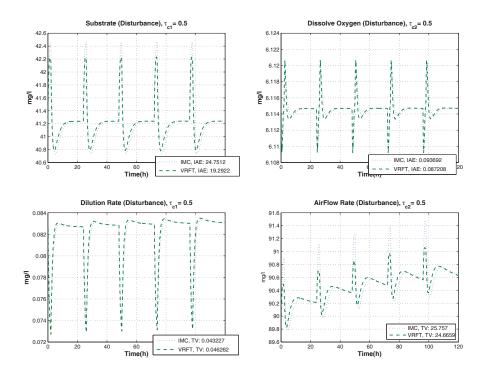


Figure 6: Responses changing the specification for the VRFT and comparison with the IMC controller, $\tau_{c1} = \tau_{c2} = 0.5$

Acknowledgment

This work has received financial support from the AECI-PCI program A/025100/09 and from the Spanish CICYT program under grant DPI2007-63356. Research work by J.D. Rojas has received financial support from the Universitat Autònoma de Barcelona. Support from the Universidad de Costa Rica is greatly appreciated.

Bibliography

- F. Nejjari, A. Benhammou, B. Dahhou, and G. Roux, "Non-linear multivariable adaptive control of an activated sludge wastewater treatment process," *Int. J. Adapt. Control Signal Process.*, pp. 347–365, 1999.
- [2] G. Guardabassi and S. Savaresi, "Virtual reference direct design method: an off-line approach to data-based control system design," *Automatic Control, IEEE Transactions on*, vol. 45, pp. 954–959, May 2000.
- [3] M. C. Campi, A. Lecchini, and S. M. Savaresi, "Virtual reference feedback tuning: a direct method for the design of feedback controllers," *Automatica*, vol. 38, no. 8, pp. 1337 – 1346, 2002.
- [4] A. Lecchini, M. Campi, and S. Savaresi, "Sensitivity shaping via virtual reference feedback tuning," in *Decision and Control, 2001. Proceedings of the 40th IEEE Conference on*, vol. 1, pp. 750–755 vol.1, 2001.
- [5] A. Lecchini, M. Campi, and S. Savaressi, "Virtual reference feedback tuning for two degree of freedom controllers," *International Journal of Adaptative control and Signal Processing*, vol. 16, no. 5, pp. 355–371, 2002.
- [6] R. Vilanova, R. Katebi, and V. Alfaro, "Multi-loop PI-based control strategies for the Activated Sludge Process," in *Emerging Technologies and Factory Automation*, 2009. ETFA 2009. IEEE International Conference on, September 2009.
- [7] S. Caraman, M. Sbarciog, and B. Marian, "Predictive Control of aWastewater Treatment Process," *International Journal of Computers, Communications & Control*, vol. 2, no. 2, pp. 132–142, 2007.
- [8] M. Morari and E. Zafirou, *Robust Process Control*. Prentice-Hall International, 1989.

Erratum to the paper

A Metrics-based Diagnosis Tool for Enhancing Innovation Capabilities in SMEs (Int. J. of Computers, Communications & Control ISSN 1841-9836, E-ISSN 1841-9844 Vol. V (2010), No. 5, pp. 919-928)

by

J. Sepulveda, J. Gonzalez, M. Camargo, M. Alfaro

Juan Sepulveda, Miguel Alfaro Department of Industrial Engineering, University of Santiago of Chile 3769 Ecuador Ave. Santiago, Chile. PO Box 10233

Mauricio Camargo Nancy-Universite / ERPI (Equipe de Recherche des Processus Innovatifs) 8, rue Bastien Lepage 54010 Nancy Cedex, France

Javier Gonzalez
Nancy-Universite / ERPI (Equipe de Recherche des Processus Innovatifs)
8, rue Bastien Lepage 54010 Nancy Cedex, France and Department of Industrial Engineering,
University of Santiago of Chile
3769 Ecuador Ave. Santiago, Chile. PO Box 10233

In our paper, we presented, a metrics-based diagnosis tool for measuring and enhancing the innovation capabilities in SMEs, along with a set of preliminary results from case-based studies at the local industry. Then we indicated that the paper proposed a new method by studying the competences of SMEs in concepts tied to innovation.

We would like to remark that a more accurate title of the article is "Application of A Metricsbased Diagnosis Tool for Enhancing Innovation Capabilities in Chilean SMEs", as explained in section 4.

The mathematical model of flowsort used to classify the enterprises in Chile was developed by Nemery [15] as indicated in the Bibliography, and the subsequent tool for application was the result of a join work between ERPI-INPL (National Polytechnic Institute of Lorraine) and the University of Portsmouth.

Gonzalez [16] was co-supervised graduate student of both USACH and INPL.

Addendum to the References

[16] Gonzalez, J. Final Report. Research Master's Program in Innovation Management and Industrial Design, Nancy, France, September, 2009.

Author index

Aderounmu G.A., 204 Ajayi A.O., 204 Akinboro S.A., 204 Alfaro V.M., 367 Aseri T.C., 214 Bartha A., 222 Borne P., 246 Chen J., 258 Dumitrescu D., 222 Hao Z., 227 Haoliang S., 236 Harbaoui Dridi I., 246 He D., 258 Hu J., 258 Hwang C., 349 Jung H., 349 Kammarti R., 246 Kim H., 349 Krassovitskiy A., 266 Ksouri M., 246 Lefranc G., 278 Leighton F., 278 Li J., 286 Lixiang L., 236 Lupu V., 297 Mallapur J.D., 305 Manvi S.S., 305 Novakovic J., 317 Ogundoyin I.K., 204 Olajubu E.A., 204 Olaniyan O.M., 204 Osorio R., 278 Pârv B., 328

Pan Q., 286 Park S., 349 Parpucea I., 328 Perera I., 337 Rankov S., 317 Rao D.H., 305 Rojas J.D., 367 Singla N., 214 Socaciu T., 328 Sun D., 349 Tiliute D.E., 297 Vilanova R., 367 Wand S., 286 Xiaohui H., 236 Xie S., 286 Yu N., 227 Zhong S., 227