



NATIONAL AND KAPODISTRIAN UNIVERSITY OF ATHENS
DEPARTMENT OF INFORMATICS & TELECOMMUNICATIONS

ABSTRACTS OF DOCTORAL DISSERTATIONS



Athens 2011

Volume 6

NATIONAL AND KAPODISTRIAN UNIVERSITY OF ATHENS
DEPARTMENT OF INFORMATICS & TELECOMMUNICATIONS

ABSTRACTS OF DOCTORAL DISSERTATIONS

The Committee of Research and Development

A. Eleftheriadis
M. Koubarakis
E. Manolakos (Chair)
T. Theoharis

ISSN: 2241-2115

Copyright © 2011

Volume 6

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
Panepistimiopolis, 15784 Athens, Greece

PREFACE

This volume contains the extended abstracts of the Doctoral Dissertations conducted in the Department of Informatics and Telecommunications, National and Kapodistrian University of Athens and completed in the time period December 2009 to December 2010.

The goal of this volume is to demonstrate the breadth and quality of the original research performed by our Ph.D. students and to facilitate the dissemination of their research results. We are glad to present the sixth collection of this kind and expect this initiative to continue in the years to come. The submission of an extended abstract in English is required by all our graduating Ph.D. students.

We would like to thank the students who contributed to this volume and hope that this has also been a positive experience for them. Finally, we would like to thank graduate student Mr. Nikos Bogdos for his help in putting together this publication. The painting shown in the cover is "*Pathways*" by artist Dionisios Cabolis.

The Committee of Research and Development

A. Eleftheriadis

M. Koubarakis

E. Manolakos (Chair)

T. Theoharis

Athens, October 2010

Table of Contents

Preface	3
Table of Contents	5
Doctoral Dissertations	
Vassilis Dalakas , <i>Development, Modelling and Simulation of Advanced Techniques to Estimate and Improve the Performance of Satellite Telecommunications Systems in Non-Linear Channels.</i>	9
Lambros E. Dermentzoglou , <i>A Defect Oriented Approach for Testing RF Front-Ends of Wireless Transceivers.</i>	21
Katerina Glezou , <i>Development of Learning Environments with Use of Logo programming language in teaching praxis.</i>	33
Emmanouil N. Kafetzakis , <i>Effective Capacity Theory for Modeling Systems with Time-Varying Servers, with an Application to IEEE 802.11 WLANs.</i>	45
Dimitris Kanakidis , <i>Secure Optical Systems based on Chaotic Carriers.</i>	57
Alexandros Kapsalis , <i>Optical Microring Devices for Optical Networks Applications: Investigation, Design and Characterization.</i>	69
Christos Kareliotis , <i>BPEL scenario execution: QoS-based dynamic adaptation and exception resolution.</i>	81
Dimitrios Kogias , <i>Study and Design of Algorithms for Information Dissemination in Unstructured Networking Environments</i>	93
Christos Konaxis , <i>Algebraic algorithms for polynomial system solving and applications</i>	103
George Kontolemakis , <i>A generic product ontology based on software agents incorporating negotiation and decision support techniques.</i>	115
Alexandros Makris , <i>A multimedia content modeling and classification methodology using visual information for the protection of sensitive user groups.</i>	127

Gerasimos Mileounis , <i>Nonlinear Signal Processing and its Applications to Telecommunications.</i>	139
Alexandros Papadimitriou , <i>Adaptive Educational Hypermedia Systems on the Web for the Didactics of Science and Technology.</i>	151
Zafeiro G. Papadimitriou , <i>Performance Analysis of Wireless Single Input Multiple Output Systems (SIMO) in Correlated Weibull Fading Channels.</i>	161
Eleni Patouni , <i>Decision Management and Object--oriented Protocol and Services Reconfiguration in Future Internet Autonomic and Heterogeneous Telecommunication Environments.</i>	173
Dimitrios Pierrakos , <i>A Web Usage Mining Framework for Web Directories Personalization.</i>	185
Gerasimos G. Pollatos , <i>Static and dynamic graph algorithms with applications to infrastructure and content management of modern networks.</i>	197
Constantinos Rizogiannis , <i>Algorithms for Space-Time Equalization of Wireless Channels.</i>	209
Theodoros Rokkas , <i>Technoeconomic analysis of Next Generation Networks.</i>	221
Odysseas Sekkas , <i>Context Information Management for Pervasive Computing.</i>	233
Dimitris Skyrianoglou , <i>Quality of Service Provision for IP Traffic over Wireless Local Area Networks.</i>	245
Teta Stamati , <i>Transformational Government and Electronic Government Adoption Model.</i>	257
Panagiotis Tsakanikas , <i>Image processing methods and algorithms for accurate protein spot detection in 2-dimensional gel electrophoresis (2DGE).</i>	273
Vassileios Tsetsos , <i>Semantic Information Management for Pervasive Computing.</i>	285
Leonidas M. Tzevelekas , <i>Efficient algorithms for topology control and information dissemination/ retrieval in large scale Wireless Sensor Networks.</i>	297

Stylios Tzikopoulos , <i>Analysis and Retrieval of Mammographic Images.</i>	309
Georgios Vamvakas , <i>Processing and Recognition of Handwritten Documents.</i>	319
Konstantinos Xenoulis , <i>Information Theory and Signal Processing for (nonlinear) Communication Channels.</i>	331

Development, Modelling and Simulation of Advanced Techniques to Estimate and Improve the Performance of Satellite Telecommunications Systems in Non-Linear Channels

Vassilis Dalakas *

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
vdalakas@di.uoa.gr

Abstract. Non-linear amplifiers that are used near saturation in satellites, cause severe distortions of the transmitted signal. This dissertation proposes novel signal processing techniques that counteract these distortions to improve performance. Firstly, the fact that the non-linearities introduced by memoryless bandpass amplifiers preserve the symmetries of the M -ary Quadrature Amplitude Modulation (M -QAM) constellation is exploited, to present a Cluster-Based Sequence Equalizer (CBSE). The proposed equalizer exhibits enhanced performance compared to other techniques, including Volterra equalizers and neural network equalizers. Moreover, this gain in performance is obtained at a substantially lower computational cost. Secondly, a new constellation shaping technique, which efficiently and effectively reduces the Peak to Average Power Ratio (PAPR) of Orthogonal Frequency Division Multiplexing (OFDM) systems, is proposed. The proposed technique requires minimal implementation complexity, while it offers considerable performance gains. Closed form analytical expressions for the distribution of the PAPR and the Bit Error Rate (BER) are derived and their accuracy is verified via simulations. Finally, a detailed comparative study of two single-carrier frequency-division multiple access schemes is presented, namely localized FDMA scheme (LFDMA) and interleaved FDMA scheme (IFDMA), versus orthogonal scheme (OFDMA), for a satellite up-link based on the digital video broadcasting via satellite (DVB-S) standard. Considering two state-of-the-art high power amplifiers, operating in the K- and S-bands, the performance of synchronous and asynchronous LFDMA, IFDMA and OFDMA is evaluated in a multi-user environment including inter-block interference.

Keywords: CBSE equalization, PAPR reduction, SC-FDMA, OFDMA, satellite systems.

* Dissertation Advisor: Sergios Theodoridis, Professor

1 Introduction

The role of a satellite is to receive a signal from an earth station or another satellite (uplink) and, acting as a simple repeater, to transmit it to another earth station or satellite (downlink) [1]. Fig. 1 illustrates a typical satellite communication system [1]. New generation satellites have regenerative payloads [2, 3] with on-board processing. This means that the baseband transmitted signal is available on-board, via demodulation, and hence uplink and downlink can be treated separately.

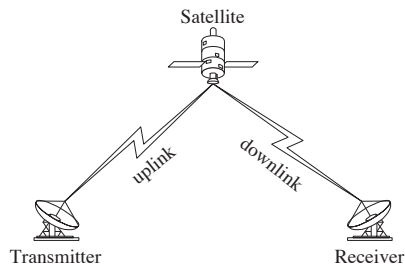


Fig. 1. Satellite communication system.

The need to maximally exploit on-board resources in a satellite communication system often imposes driving a high power amplifier (HPA), such as the travelling wave tube amplifier (TWTA), at or near its saturation point, resulting in a nonlinear distortion of the signal, and rendering the overall link nonlinear. To overcome nonlinear distortions, constant modulus constellation symbols (e.g., 4-QAM) are commonly used [2]. However, large QAM signal constellations have to be adopted whenever high bandwidth efficiency is required [4], resulting in severe nonlinear distortions. Two approaches have been proposed for solving the problem of correct reception of the transmitted signal in those cases: (a) Equalization [5–7] and (b) Predistortion or power amplifier linearization [5, 8–10]. Equalization refers to processing the signal at the receiver side in order to recover the transmitted data, thus post-cancelling the link’s nonlinear (amplifier) and linear (multipath) distortions. Conventional linear equalizers combat only the Inter-Symbol Interference (ISI), introduced by the propagation channel, while nonlinear equalizers aim also at equalizing the non-linear effects of the HPA. The main drawback of the equalization approach is the additional cost and the computational load it entails for each terminal. On the other hand, predistortion techniques aim at pre-cancelling the nonlinear effects via modelling the inverse of the amplifier characteristic and pre-distorting the data *prior to* the amplification stage. The overall characteristic then becomes linear. The advantage of this approach lies in the fact that only a single system is needed for cancelling the HPA nonlinearity at the satellite, compared to using an equalizer in each

terminal. On the other hand, its main drawback is that the predistorter must be on-board, so it cannot be applied to the satellites already on orbit. Moreover, in case multipath is present, an equalizer at the terminal side is still needed.

The non-linearities introduced by memoryless bandpass amplifiers preserve the symmetries of the M -QAM constellation. A CBSE that takes advantage of these symmetries is presented. The proposed equalizer exhibits enhanced performance compared to other techniques, including the conventional linear transversal equalizer, Volterra equalizers and Radial Basis Function (RBF) network equalizers [5]. Moreover, this gain in performance is obtained at a substantially lower computational cost.

In recent years, the increasing commercial demand for higher data rates has led to the utilization of OFDM in several well known standards, including the 2nd generation of Digital Video Broadcasting by Satellite (DVB-S2) [11], Digital Video Broadcasting by Satellite Handheld (DVB-SH) [12] and 3rd Generation Partnership Project (3GPP) [13]. The main technical advantage for such a choice is OFDM's robustness in the presence of frequency selective fading channels commonly encountered in wireless broadband communication systems [14]. Moreover, OFDM simplifies the equalization process at the receiver's side. A major drawback of the OFDM technology is the high PAPR of the transmit signal [14]. The presence of high PAPR becomes even more critical when non-constant envelope modulation techniques, such as M -QAM, are used for transmission in order to increase the overall system capacity. Hence, the transmit power amplifiers have to operate with a large input power back-off, (IBO), from their peak power which leads to poor power efficiency [14].

In the second chapter of the dissertation, a novel, low computational complexity constellation shaping technique is presented and its performance is analysed and evaluated in order to counteract high PAPR. In the past the problem of high PAPR in OFDM has been extensively studied and various techniques have been proposed and analysed for its reduction [15]. The novelty of the proposed technique is that, in contrast to previous art, constellation shaping takes place *after* the Inverse Fast Fourier Transform (IFFT) operation instead of before, thus, directly modifying the OFDM symbol. This makes the technique independent of the modulation format, the number of subcarriers, N , or the input to the IFFT signal combination. Furthermore, handshake between transmitter and the receiver is not necessary, power increase of the transmit signal is avoided through power normalization while this procedure can be easily reversed at the receiver. Analytical closed form expressions for the PAPR distribution and the BER performance of the proposed technique are provided.

Finally, the dissertation presents a thorough and detailed comparison study of multiple access techniques for the DVB-S2. This study was done for the European Space Agency in collaboration with the European Satellite Communications Network of Excellence (SatNEx) [16]. To the best of our knowledge, similar studies for state-of-the-art satellite based systems, such as DVB-S type, are not available in the open technical literature. Orthogonal Frequency Division Multiple Access (OFDMA) is the multiple access scheme that naturally extends

OFDM to simultaneously serve multiple users. Single Carrier Frequency Division Multiple Access (SC-FDMA) schemes are employed as alternative access schemes, which offer reduced PAPR as compared to OFDMA's high PAPR [17]. Although SC-FDMA utilize single carrier modulation at the transmitter and frequency domain equalization at the receiver and typically achieve lower PAPR, they have similar transmitter structure and BER performance, as compared to an OFDMA system. The key difference in the transmitters of the two schemes is the presence of an additional discrete Fourier transform (DFT) in SC-FDMA. Among the various SC-FDMA schemes, the most popular are: (i) Localized FDMA (LFDMA); and (ii) Interleaved FDMA (IFDMA) [17]. The application of z -FDMA¹ schemes for state-of-the-art satellite multi-user systems is considered. In particular, z -FDMA access scheme performances, for a DVB-S satellite link operating in dual frequency bands (K- and S-Bands) in the presence of non-linear HPA and synchronization error are presented.

This dissertation summary has the following structure. The preservation of M -QAM symmetries by memoryless nonlinear amplifiers is demonstrated in section 2, where the new equalization algorithm is referred. In section 3 a short description of the proposed PAPR reduction technique is given while in section 4 the system model is illustrated and a summary of the comparative results is offered. Section 5 gives some concluding remarks. All the results presented here are already published [5, 18–20].

2 Cluster-Based Equalizer for Satellite Communication Channels with M -QAM Signaling

There are two technologies for the HPA on board satellites: Travelling Wave Tube Amplifiers (TWTAs) and Solid State Power Amplifiers (SSPA).

- TWTAs can generally be considered as memoryless. They are characterized by an AM/AM conversion and an AM/PM conversion. These are commonly modelled by a Saleh model [21].
- SSPAs have intrinsically memory. It is common to model an SSPA with memory by a memoryless non-linearity (see [22] for the type of the non-linearity) followed by a linear IIR filter [8].

Here we will deal only with TWT amplifiers, due to their dominant use in satellites. According to Saleh's model [21], an input

$$x(t) = A \cos(2\pi f_c t + \theta) \quad (1)$$

into a bandpass amplifier produces an output of the form [23, 24]:

$$z(t) = g(A) \cos[2\pi f_c t + \theta + \Phi(A)] \quad (2)$$

¹ For the conciseness of the presentation, from now on and unless otherwise stated, the notation z -FDMA ($z \in \{L, I, O\}$) will be used as a common representation of the three multiple access schemes considered in this work, i.e., LFDMA, IFDMA and OFDMA.

where the nonlinear gain function $g(A)$ is commonly referred to as the *AM/AM characteristic* and the nonlinear phase function $\Phi(A)$ is called the *AM/PM characteristic*. These are expressed as

$$g(A) = \frac{\alpha_a A}{1 + \beta_a A^2} \quad (3)$$

$$\Phi(A) = \frac{\alpha_p A^2}{1 + \beta_p A^2} \quad (4)$$

The adopted signaling scheme, namely rectangular M -ary QAM, may be viewed as a form of combined digital amplitude and digital phase modulation. In view of eqs. (1–4), the baseband complex envelope of the TWTA output is given by

$$\begin{aligned} \tilde{z}(t) &= g[A(t)] e^{j\{\theta(t) + \Phi[A(t)]\}} \\ &= [A(t) e^{j\theta(t)}] \left\{ \frac{g[A(t)]}{A(t)} e^{j\Phi[A(t)]} \right\} \\ &\triangleq \tilde{x}(t) G(|\tilde{x}(t)|) \end{aligned} \quad (5)$$

where $\tilde{\cdot}$ denotes complex envelope. In words, *the output of the TWTA is the product of the input signal with a factor that depends only on the input amplitude*. The result is an amplitude change and a phase rotation of the input signal constellation points. Eq. (5) implies that *the change is the same for all constellations points that share the same energy level*. The M symbols in the input constellation can be grouped in two possible ways (see Fig. 2a for the example of 16-QAM):

1. in I circles on the complex plane, where I is the number of the energy levels (for the 16-QAM case, $I = 3$), and
2. in $M/4$ squares (four points in each square), that are centered on the origin.

It is not difficult to see that *the above symmetries (1, 2) of the constellation are preserved* by the amplifier. This is a consequence of the fact that the angles between the constellation points that lie on the same energy circle remain unaltered (see Fig. 3). Thus, the resulting points continue to form squares centered on the origin, as it was the case prior to the application of the nonlinearity.

These symmetries can be efficiently exploited to reduce the total number of cluster centers to be estimated directly from the training sequence in the CBSE equalizer, thus leading to a significant gain in performance, compared to Volterra and NN-based techniques, and at a significantly lower computational cost.

The performance of the proposed equalizer is compared with two of the most widely used non-linear equalizers: a Volterra series equalizer and an RBF equalizer. The algorithms are compared in terms of the resulting BER and their computational requirements. CBSE outperforms its competitors in every case. The Bayesian equalizer performs almost equally well, however, it is far more expensive in terms of computational requirements. It is of interest to note that

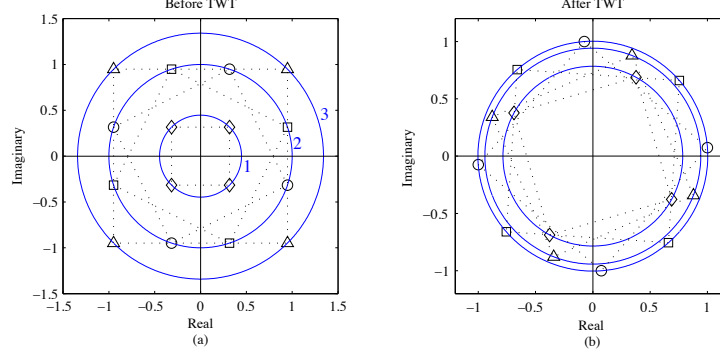


Fig. 2. 16-QAM constellation at the (a) input and (b) output of the TWTA. The 3 energy levels and the 4 squares formed by the 16 constellation points are illustrated.

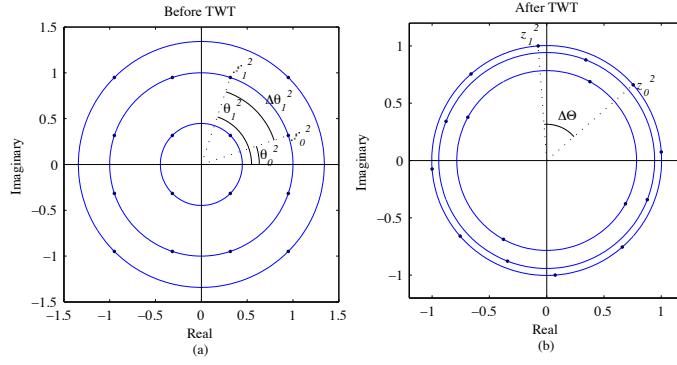


Fig. 3. 16-QAM constellation at the (a) input and (b) output of the TWTA. Angles between equal modulus symbols are shown: $\Delta\theta = \Delta\theta_1^2$.

even in the case of 0 dB IBO (full power efficiency [1]) with 16-QAM where other methods fail [7], the proposed equalizer still offers some gain (Fig. 4).

Table 1 shows the total number of real operations required for the processing of a received block consisting of 20 training samples (per energy zone) and 500 data symbols, for a two-taps channel with the 16-QAM signaling scheme. Observe that the superior performance of the CBSE equalizer is attained at a substantially lower complexity, especially in terms of real multiplications/divisions.

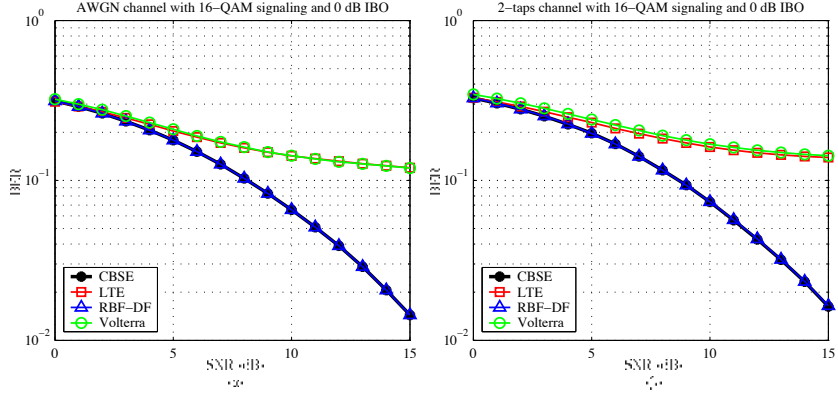


Fig. 4. BER performance for 16-QAM at 0 dB IBO. (a) AWGN and (b) 2-taps channel.

Method	Mul/Div		Add/Sub		$(\cdot)^2$		$\exp(\cdot)$
	Training	Decision	Training	Decision	Training	Decision	
CBSE	32	0	356	512000	—	256000	—
RBF-DF		128000		1016000		512000	128000
LTE	1620	6000	1800	5000	360	—	—
Volterra	19080	114000	16920	77000	2520	—	—

Table 1. The total number of real operations for each equalizer, for a 2-taps channel ($L = 2$) with 16-QAM input, needed to process a packet of 60 training and 500 information samples; $L_{\text{LTE}} = 3$, $p = 3$, $L_V = 21$, $M_V = 36$.

3 The Novel PAPR Reduction Technique

To obtain a PAPR lower than that of the original time domain vector \mathbf{x} , a new time domain vector \mathbf{y} is needed so that the following holds

$$\text{PAPR}(\mathbf{y}) < \text{PAPR}(\mathbf{x}) \Leftrightarrow \frac{P_{\max}^y}{P_{\text{av}}^y} < \frac{P_{\max}^x}{P_{\text{av}}^x} \quad (6)$$

where

$$P_{\max}^x = \max_{0 \leq n \leq N-1} |x_n|^2, \quad (7)$$

$$P_{\text{av}}^x = \frac{1}{N} \sum_{n=0}^{N-1} |x_n|^2 \quad (8)$$

and similarly for P_{\max}^y , P_{av}^y . It can be easily verified that if

$$P_{\max}^y = P_{\max}^x + P_1 \quad (9)$$

and

$$P_{\text{av}}^y = P_{\text{av}}^x + P_2 \quad (10)$$

with

$$\frac{P_1}{P_2} < \frac{P_{max}^x}{P_{av}^x} \quad (11)$$

then Eq. (6) always holds. Thus, if a new vector $\mathbf{y} = \mathbf{x} + \alpha \mathbf{u}$ is defined with α positive and $\mathbf{u} = [x_1/|x_1| \dots x_n/|x_n| \dots x_N/|x_N|]^T$ then $\text{PAPR}(\mathbf{y})$ (as shown in the dissertation) is always less than or equal to the initial PAPR, independently of the vector \mathbf{x} . If α is known at the receiver there is no need of any side information or transmit power increase. The parameter α can be selected during the communication system design according to the desired level of BER or CCDF performance.

Analytical expressions for the achievable PAPR reduction and equivalent BER having α as the main parameter of interest are presented and verified by means of computer simulations (Fig. 5). The performance of the proposed tech-

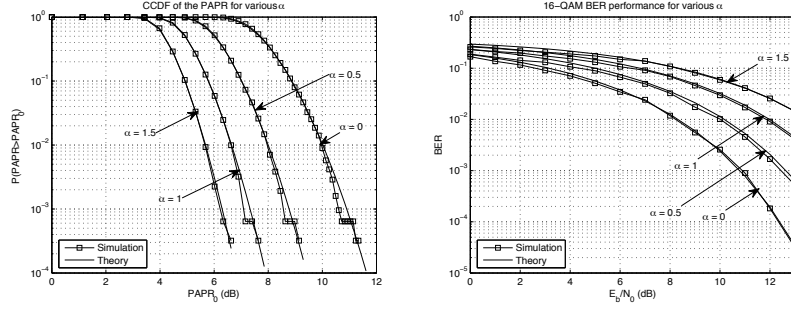


Fig. 5. CCDF and BER performances of the proposed technique for various α for 16-QAM and $N = 256$.

nique in terms of PAPR reduction is superior as compared to other known similar techniques at the expense of a BER performance degradation. Performance evaluation comparisons were made with the TI technique, as a representative from the family of constellation shaping techniques, and with the AC technique, as a rival with comparable simplicity. The obtained performance evaluation results have shown that the proposed technique is an attractive alternative for PAPR reduction regardless of the number of subcarriers N and the modulation format.

4 SC-FDMA vs. OFDMA in DVB-S

The block diagram of the system model under consideration, shown in Fig. 6, represents the up-link of a typical multi-user satellite communication system. Both synchronous and asynchronous signal reception has been considered for two state-of-the-art HPAs operating in the K- and S-band. Performance evaluation results have shown that, although, IFDMA outperformed the other two schemes

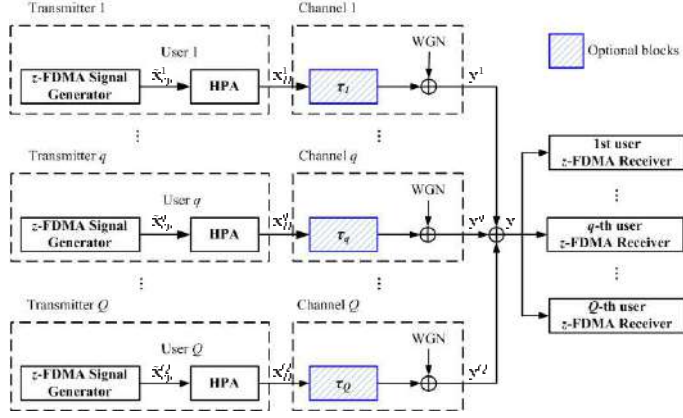


Fig. 6. Block diagram of the system under consideration.

in terms of Total Degradation for a synchronous system, for asynchronous reception it is the most sensitive to degradation. OFDMA, due to its large PAPR, has been found as the most sensitive to non-linearity. On the contrary, LFDMA had only slightly inferior performance as compared to IFDMA for synchronous reception while it outperformed the other two access schemes in the asynchronous scenario examined, i.e., in the presence of IBI and non-linearity.

Table 2. Advantages and disadvantages for each access scheme in a non-linear channel and in the presence of IBI. A (+) means advantage, a (−) means disadvantage and (≈) means almost equal to the best.

	OFDMA	LFDMA	IFDMA
Non-linearity	−	≈	+
$i = 0 \rightarrow \tau_C/T_s = \frac{1}{8}$	−	≈	+
$i = 6 \rightarrow \tau_C/T_s = \frac{1}{64}$	−	+	≈
$i = 12 \rightarrow \tau_C/T_s = \frac{7}{32}$	−	+	−
$i = 18 \rightarrow \tau_C/T_s = \frac{17}{64}$	−	+	−

5 Conclusions

This dissertation proposes novel signal processing techniques that deal efficiently with the non-linear distortions introduced by the satellite amplifiers. A CBSE for satellite channels is proposed exploiting the fact that TWT memoryless non-linearities respect the symmetries underlying the signaling scheme, thus leading to a significant gain in performance, compared to Volterra and NN-based techniques, and at a significantly lower computational cost. A simple time domain

constellation shaping technique that achieves PAPR reduction in OFDM with minimal complexity is introduced. Theoretical expressions are derived for key performance measures, i.e., BER and CCDF, while simulation results verify the accuracy of the derived analytical expressions. The performance of the proposed technique in terms of PAPR reduction is superior as compared to other known similar techniques at the expense of a BER performance degradation. Finally, a thorough comparison study is presented for two SC-FDMA schemes, LFDMA and IFDMA versus OFDMA, for a satellite up-link, based on DVB-S. Both synchronous and asynchronous signal reception has been considered for two state-of-the-art HPAs operating in the K- and S-band. As a further research topic it will be interesting to investigate the total degradation performance of the proposed PAPR reduction technique in a DVB-S scenario combined with the CBSE and LFDMA.

References

1. G. Maral and M. Bousquet, *Satellite Communication Systems*, 2nd ed. Chichester: John Wiley & Sons, Ltd., 1993.
2. M. Ibnkahla, Q. M. Rahman, A. I. Sulyman, H. A. Al-Asady, J. Yuan, and A. Safwat, "High-speed satellite mobile communications: Technologies and challenges," *Proceedings of the IEEE*, vol. 92, no. 2, pp. 312–339, February 2004.
3. ESA SkyPlex. [Online]. Available: <http://telecom.esa.int/telecom/www/object/index.cfm?fobjectid=12143>
4. F. Xiong, "Modem techniques in satellite communications," *IEEE Communications Magazine*, vol. 32, no. 8, pp. 84–98, August 1994.
5. G. E. Corazza, Ed., *Digital Satellite Communications*, 1st ed. Springer, 2007.
6. E. Biglieri, A. Gersho, R. D. Gitlin, and T. L. Lim, "Adaptive cancellation of non-linear intersymbol interference for voiceband data transmission," *IEEE Journal on Selected Areas in Communications*, vol. 2, no. 5, pp. 765–777, September 1984.
7. S. Bouchired, D. Roviras, and F. Castanié, "Equalization of satellite mobile channels with neural network techniques," *Space Communications*, vol. 15, no. 4, pp. 209–220, 1999.
8. F. Langlet, H. Abdulkader, D. Roviras, A. Mallet, and F. Castanié, "Comparison of neural network adaptive predistortion techniques for satellite down links," in *Proceedings of IJCNN'01*, Washington, DC, USA, July 2001.
9. F. Langlet, D. Roviras, A. Mallet, and F. Castanié, "Mixed analog/digital implementation of MLP NN for predistortion," in *Proceedings of IJCNN'02*, Hawaii, USA, May 2002.
10. F. Langlet, H. Abdulkader, and D. Roviras, "Predistortion of non-linear satellite channels using neural networks: Architecture, algorithm and implementation," in *Proceedings of EUSIPCO'02*, Toulouse, France, September 2002.
11. ETSI, "Digital Video Broadcasting (DVB): Second generation framing structure, channel coding and modulation system for Broadcasting, Interactive Services, News Gathering and other broadband satellite applications," Tech. Rep. EN 302 307 v1.1.2, June 2006. [Online]. Available: <http://www.etsi.org>
12. —, "Digital Video Broadcasting (DVB): Framing Structure, channel coding and modulation for Satellite Services to Handheld devices (SH) below 3 GHz," Tech. Rep. EN 302 583 v1.0.0, December 2007. [Online]. Available: <http://www.etsi.org>

13. 3rd Generation Partnership Project (3GPP), "Evolved universal terrestrial radio access (e-utra)," Tech. Rep. 3GPP TR 36.942, January 2009. [Online]. Available: <http://www.3gpp.org>
14. R. van Nee and R. Prasad, *OFDM for Wireless Multimedia Communications*, 1st ed. Norwood, MA, USA: Artech House, 2000.
15. S. H. Han and J. H. Lee, "An overview of peak-to-average power ratio reduction techniques for multicarrier transmission," *IEEE Wireless Communications [see also IEEE Personal Communications]*, vol. 12, no. 2, pp. 56–65, April 2005.
16. "Study of Enhanced Multicarrier (OFDM) Digital Transmission Techniques for Broadband Satellites," Tech. Rep. ESTEC Contract. N. 21072/07/NL/AD, January 2009. [Online]. Available: <http://www.esa.org>
17. H. G. Myung and D. J. Goodman, *Single Carrier FDMA: A New Air Interface for Long Term Evolution*, 1st ed. John Wiley & Sons, Ltd., 2008.
18. E. Kofidis, V. Dalakas, Y. Kopsinis, and S. Theodoridis, "A novel efficient cluster-based mlse equalizer for satellite communication channels with m-qam signaling," *EURASIP Journal on Applied Signal Processing*, 2006, article ID 34343.
19. V. Dalakas, A. Rontogiannis, and P. Mathiopoulos, "Time domain constellation shaping technique for peak-to-average power ratio reduction," *IET Communications*, vol. 3, no. 7, pp. 1144–1152, 2009. [Online]. Available: <http://link.aip.org/link/?COM/3/1144/1>
20. V. Dalakas, P. T. Mathiopoulos, F. D. Cecca, and G. Gallinaro, "Comparison study of single-carrier FDMA schemes vs. OFDMA for DVB-S systems," *IEEE Transactions on Broadcasting*, 2010, **(Under review)**.
21. A. A. M. Saleh, "Frequency-independent and frequency dependent nonlinear models of TWT amplifiers," *IEEE Transactions on Communications*, vol. 29, no. 11, pp. 1715–1720, November 1981.
22. C. Rapp, "Effects of HPA-nonlinearity on a 4-DPSK/OFDM-signal for a digital sound broadcasting system," in *Proceedings of the 2nd European Conference on Satellite Communications ECSC-2*, Liege, Belgium, October 1991.
23. A. L. Berman and C. E. Mahle, "Nonlinear phase shift in travelling wave as applied to multiple access communication satellites," *IEEE Transactions on Communications Technology*, vol. 198, no. 1, pp. 37–48, February 1970.
24. J. B. Minkoff, "Wideband operation of nonlinear solid state power amplifiers – Comparison of calculations and measurements," *AT&T Bell Lab. Tech. Journal*, vol. 63, no. 2, pp. 231–248, 1984.

A Defect Oriented Approach for Testing RF Front-Ends of Wireless Transceivers

Lambros E. Dermentzoglou*

National and Kapodistrian University of Athens
Department of Informatics & Telecommunications
dermetz@di.uoa.gr

Abstract. In this dissertation the problem of testing wireless transceiver's RF Front-Ends is addressed. The proposed approach constitutes a cost effective, reliable and efficient solution for characterizing a complex system as a faulty or fault free, based on robust defect oriented build-in testing circuits. Testing techniques and related Built-In Self-Test circuits were proposed for the effective fault diagnosis of integrated differential Low Noise Amplifiers, Mixers and Voltage Controlled Oscillators for both receiver and transmitter parts. These individual test circuits were finally combined to form a fully integrated test solution for RF Front-Ends of wireless transceivers.

Keywords: Built-In Self-Test, Defect-Oriented Testing, Low Noise Amplifiers, Mixers, Voltage Controlled Oscillators

1 Introduction

The subject of this dissertation lies on the field of testing analog circuits and high frequency systems on chip. Its main contribution is the proposal of a cost efficient and effective technique for testing complex RF systems, such as the RF front-end of the wireless transceivers, focusing on a defect oriented test approach.

Testing cost becomes a major concern being a large portion of the total production cost. Especially in the case of high frequency/RF integrated circuits (IC) the cost of testing is prohibitive. In this area, high cost dedicated automatic test equipments (ATEs) are used to measure performance characteristics of the circuit under test (CUT) and compare them against predefined limits that are called specifications. Although these measurements are simple, they require a variety of test resources, which along with the long test application time, increase further the manufacturing cost. Today, almost all analog circuits pass through external functional (specification based) testing procedures to ensure performance and quality. However, a key problem is that it is not always possible for the ATE to have a direct access to all or even part of the internal signals of an IC. This is mainly the case of System-on-Chip (SoC) or System-on-Package (SOP) designs. Even if some internal signals can be routed to

* Dissertation Advisor: Aggeliki Arapoyanni, Assoc. Professor

become available to the external tester, frequency limitations due to lower speed of external I/O pads may not permit their direct observation. Consequently, incorporating Built-In-Test (BIT) structures to the circuit seems today a decent compromise between area overhead and total manufacturing cost in order to improve testability and test access speed [1].

However, BIT schemes are not always suitable for the implementation of direct measurement techniques, due to the high hardware overhead that is required [2]. To overcome the cost and inabilities of functional testing, the concept of alternate test was proposed [3]. The objective of the alternate test methodology is to find a suitable test stimulus and to predict circuit characteristics accurately from the corresponding alternate test response. Although many alternate test techniques exploit BIT schemes to support testing [4], still the elaboration of test responses is accomplished off chip.

Among the different approaches proposed for built-in testing of analog circuits and Systems-on-Chip, two are the main streams [5]. The first one is considering each building block as a standalone circuit under test. A test signature is generated according to the specific requirements of the CUT and the comparison of the CUT's response with the signature provides a metric for the operational status of the circuit [6-7]. This principle, although highly effective since the test procedure is optimized for the specific CUT, it normally imposes high test circuitry overhead.

The second commonly used technique exploits a loop-back test path between the transmitter and the receiver in transceiver modules. The test stimuli are injected through the transmitter's baseband and the signature of the CUT is evaluated in receiver's baseband interface [8-9]. Although these tests are capable of detecting gain, noise or linearity distortion with practically no extra circuitry requirements, they are frequently more susceptible to fault masking due to the distance between the control and observation points inside the loopback path.

In this dissertation, defect oriented test techniques and suitable test circuits are combined together in order to formulate an efficient, easy to implement and cost effective test strategy for wireless transceiver's RF Front-Ends. In this context, its contribution lies on two fields, a) the development of suitable test techniques for the main building blocks of the front end and b) the proposal of a complete test architecture for the system under test.

2 Defect Oriented Test Techniques for RF Front-End Circuits

In this dissertation we have proposed novel techniques for testing individually the main building blocks of the RF front-end. These techniques are easy to implement and their test outputs are digital signals which can be further post-processed by cheap digital testers. Finally the test cycle, both for each building block and for the system under test requires no more than a few microseconds to be completed and has the minimum requirements regarding the number of I/O test pads, compared to the conventional techniques.

2.1 Test Technique and BIST Circuit for Voltage Controlled Oscillators

Voltage Controlled Oscillators (VCOs) are commonly used in phase locked loops (PLL) and frequency synthesizers to produce a precise and controllable reference frequency. In RF synthesizers, LC-tank oscillators are preferred due to their superior phase noise performance [10-11].

In this section, a new DFT technique and a test-circuit are proposed for the testing of high frequency LC-tank differential voltage controlled oscillators and CMOS differential ring oscillators. This circuit generates a single digital *Fail/Pass* output signal which can be easily processed by a standard digital tester or alternatively by a JTAG TAP controller [12-14].

The fault model under consideration in our work includes single resistive short and bridging faults (up to the value of 500Ohms [15]), resistive (more than 100kOhms) as well as capacitive open faults plus single parametric faults that cover parametric variations that exceed $\pm 10\%$ of the passive devices' nominal value [16], over all possible PVT conditions (process, power supply and temperature variations).

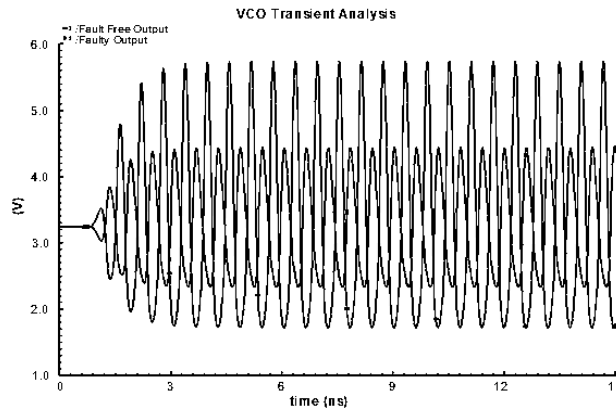


Fig. 1. Voltage Controlled Oscillator output waveforms in a faulty case.

The test strategy is based on the observed difference between the output waveforms of the VCO circuit under consideration in the fault-free and faulty cases, when the fault injection technique is followed. In Fig. 1 the faulty output waveforms are shown in the presence of a short-circuit. According to these waveforms, a compression of the oscillation amplitude at the infected output of the VCO is easily observable in the faulty case. This behavior has been verified for all faults of the previously described fault model. Thus, a circuit that can sense and discriminate the amplitude change of the oscillator outputs between the fault-free and the faulty case can serve as an embedded test vehicle. Towards this direction, we propose a suitable DFT technique and a test circuit.

The test circuit of the proposed DFT solution, embedded in the VCO unit, is presented in Fig. 2.

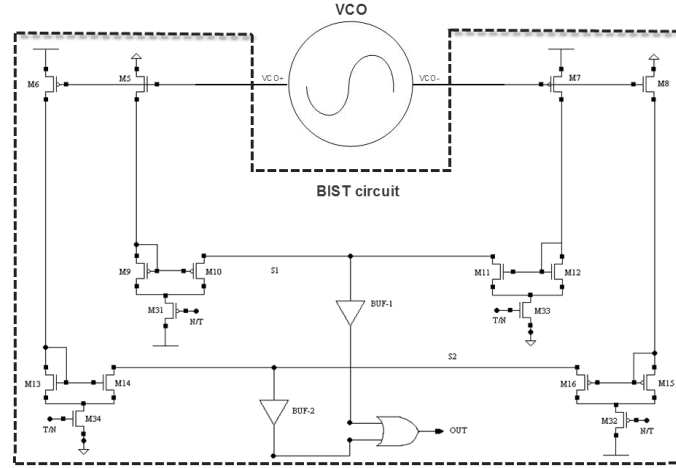


Fig. 2. The VCO along with the proposed test circuit.

In this configuration two pairs of PMOS/NMOS “sensing” transistors (M_7/M_5 and M_6/M_8) sink a certain amount of current, which depends on the output oscillation voltages. The gates of each pair of sensing transistors are directly connected to the output nodes of the VCO. Two tailing current mirrors (M_{11}/M_{12} and M_9/M_{10}), are driven by the first pair of the “sensing” transistors (M_7/M_5). Accordingly, the second pair of sensing transistors (M_6/M_8) drives the other pair of mirrors (M_{13}/M_{14} and M_{15}/M_{16}). In the fault-free case the current mirrors have been suitably adjusted so that, the currents of M_{11}/M_{12} and M_{13}/M_{14} dominate over these generated by the mirrors M_9/M_{10} and M_{15}/M_{16} , driving the nodes S_1 and S_2 towards ground. In the contrary case of a fault existence, the current of the mirrors M_9/M_{10} or M_{15}/M_{16} dominates over that generated by the mirrors M_{11}/M_{12} or M_{13}/M_{14} , rising the node S_1 or S_2 towards V_{dd} , and providing a fault indication signal.

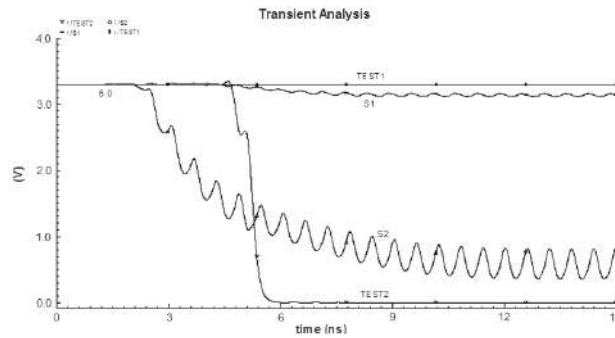


Fig. 3. Response of the test circuit in a faulty case.

In the presence of a fault, the compression of the amplitude, on at least one of the oscillator outputs, results in a reduction of the current that flows through the corresponding PMOS “sensing” transistor (M_7 or M_6). Consequently the corresponding NMOS current mirror (M_{11}/M_{12} or M_{13}/M_{14}) fails to discharge its dedicated node (S_1 or S_2) and thus the output of the corresponding buffer (BUF-1 or BUF-2) remains high, indicating fault detection as shown in Fig. 3.

In order to demonstrate the effectiveness of the proposed DFT scheme, the topology has been designed using a $0.35\mu\text{m}$ Si-Ge BiCMOS technology. The DFT scheme has been evaluated for the hard faults, while its effectiveness in the presence of soft faults has also been validated. The fault coverage results are summarized in Table 1. As it can be observed in this table, the DFT scheme is highly effective presenting an overall fault coverage performance of 91.8%.

Table 1. Fault coverage results for the BIST-VCO.

Fault Type	Fault Coverage (%)
Drain opens (DO)	(4/4) 100
Gate opens (GO)	(4/4) 100
Source open (SO)	(4/4) 100
Gate to Drain shorts (GDS)	(3/3) 100
Gate to Source shorts (GSS)	(4/4) 100
Drain to Source shorts (DSS)	(4/4) 100
Varactor-Inductor opens (VIO)	(6/10) 60
Inductor shorts (IS)	(2/2) 100
Varactor shorts (VS)	(2/2) 100
VCO Output shorts (OS)	(4/4) 100
Inductor variations (IV) $\pm 10\%$	(4/4) 100
Varactor variations (VV) $\pm 10\%$	(4/4) 100
Overall Fault coverage	(45/49) 91.8

2.2 Test Technique and BIT Circuit for Low Noise Amplifiers

A test technique for Low Noise Amplifiers along with the accompanied build-in circuit has also been proposed in the present work. The proposed BIT is capable to detect faults related to output amplitude alterations (attenuations or over-amplifications) and discriminate faulty from fault free LNA circuits (both differential and single ended), without deteriorating the overall performance of the circuit under test [17-20]. The fault model under consideration in this work includes parametric faults (passive and active devices’ parameter deviations outside specified limits) as well as catastrophic faults (resistive and capacitive opens, resistive shorts between devices’ terminals and resistive bridgings between circuit nodes).

The DLNA-BIT topology outline is presented in Fig. 4. The BIT circuit is driven directly from the outputs of the DLNA and provides a single digital PASS/FAIL signal. It consists of two main subcircuits: a) the first one is the Amplitude Alterations Detector (AAD) and b) the second is the Timing Difference Discriminator (TDD). The first subcircuit monitors the outputs of the DLNA and provides two digital

signals, TEST1 and TEST2, which perform a transition from V_{DD} to ground. In case of output amplitude alterations, due to a fault in the CUT, the two response signals of AAD present a relative transition timing difference. The second subcircuit detects the existence or not of those timing differences and discriminates faulty from fault free circuits.

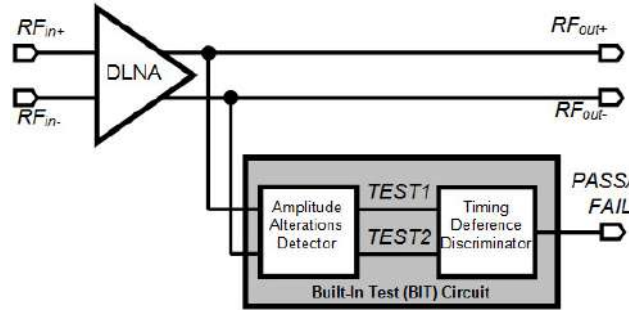


Fig. 4. Proposed LNA test technique.

The topology of the Amplitude Alterations Detector (AAD) is based on the push-pull BIST circuit presented in Fig. 2. Based on the detailed description provided earlier, as long as the DLNA is fault free, the AAD responds with a transition of TEST₁ and TEST₂ signals from V_{DD} to ground. In presence of a fault in the DLNA at least one of these signals remains permanently high. However, there are some cases where, in presence of a fault, both TEST₁ and TEST₂ signals turn to ground with a time delay, depending on the strength of the fault. In these cases, the AAD needs to cooperate with the Time-Difference Discriminator (TDD) test circuit in order to effectively identify those circuits as faulty.

The TDD subcircuit is illustrated in Fig. 5. It consists of a pair of NOR gates followed by a pair of D Flip-Flops. The outputs of the Flip-Flops drive a delay stage which is composed of a couple of buffers and capacitors. Finally, the PASS/FAIL signal is provided by a single NAND gate. The two capacitors, CAP₁ and CAP₂ ($C_{CAP1}=C_{CAP2}$), are used to insert identical delays to the TEST₁ and TEST₂ signals. The delay introduced to these signals is equal to the maximum delay that may be inserted by “acceptable” device mismatches. Obviously, this must be larger than half the DLNA signal period.

The functionality of the TDD circuit is analyzed as follows. The TEST₁ and TEST₂ signals drive the NOR gates. Since both signals are initially high, the outputs of the NOR gates are initially low. The output of each NOR gate triggers the CLK input of a D Flip-Flop. The D inputs of the Flip-Flops are permanently high, tied to V_{DD} . Before test mode activation the Flip-Flops outputs are preset to low with the use of a reset signal. Thus, the PASS/FAIL signal is initially high. As reset signal the complement of the TEST_EN signal is used. In the presence of a fault in the DLNA or “acceptable” device mismatches, one of the TEST₁, TEST₂ signals turns to low earlier

than the other. Without loss of generality, let us consider that this is the $TEST_1$ signal. Consequently, the output of the corresponding NOR_1 gate rises to high triggering the pertinent Flip-Flop. The output of the Flip-Flop goes high and the same stands for the $Delayed_TEST_1$ signal after a time delay that is determined by the capacitance value attached to it. The $Delayed_TEST_1$ signal drives the second NOR gate. Depending on the delay time for the falling edge of the $TEST_2$ signal, with respect to $TEST_1$, the second Flip-Flop may be also triggered or not. In the first case, the falling edge of the $TEST_2$ arrives earlier than the rising edge of the $Delayed_TEST_1$ signal (this is a small delay on $TEST_2$ related to device mismatches). Then, the output of the NOR_2 gate goes high (since both $TEST_2$ and $Delayed_TEST_1$ signals are low) and the second Flip-Flop is triggered raising its output to high. Consequently, both inputs of the $NAND$ gate are permanently high resulting in a low $PASS/FAIL$ response which indicates that the DLNA is fault free. In the second case, the falling edge of the $TEST_2$ signal arrives later than the rising edge of the $Delayed_TEST_1$ signal (this is a greater delay on $TEST_2$ related to a fault in the DLNA). Then, the output of the corresponding NOR_2 gate remains permanently low and the same stands for the output of the related Flip-Flop. Consequently, the second input of the $NAND$ gate is low and the $PASS/FAIL$ signal remains high indicating the presence of the fault in the DLNA.

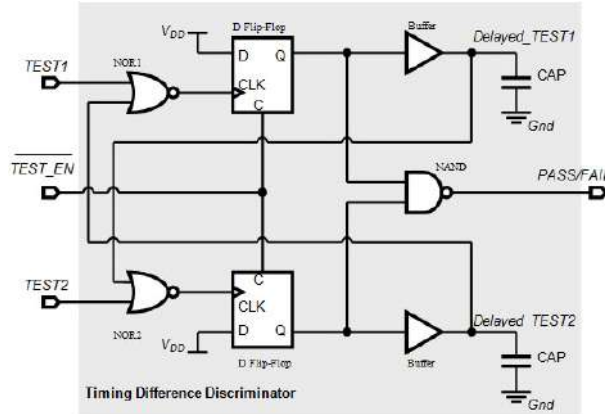


Fig. 5. The Timing Difference Discriminator.

Simulations on a typical DLNA, designed in a $0.35\mu m$ Si-Ge BiCMOS technology, have demonstrated an overall fault coverage of the proposed test technique and the BIT circuit over 90%. The fault coverage results are summarized in Table 2. The proposed test technique can be applied in either single ended or differential Low Noise Amplifiers. Moreover the test circuit offers a high fault coverage at the expense of a very low silicon area cost. Consequently, suspicious LNA circuits can be easily identified early in the production cycle (e.g. at the wafer or die level) reducing the total manufacturing cost.

Table 2. Fault coverage results for the LNA-BIT.

Fault Type	Fault Coverage (%)
Resistive Shorts	16/20
Resistive Bridgings	19/20
Resistive/Capacitive Opens	26/26
Transistor Width Parametric	8/8
Transistor Length Parametric	4/4
Inductors Parametric	4/8
Capacitors Parametric	8/8
Overall Fault Coverage	85/94 (90.4%)

2.3 Test Technique and BIST Circuit for RF Mixers

Mixers are vital parts in every wireless transceiver, regardless of the selected architecture. In this section we present a defect oriented BIST technique for RF front-end Mixers [21]. According to this, the Mixer is operated as a homodyne circuit and the generated DC voltage at its output is used as test observable. This voltage can further be used to control the oscillation frequency of a simple voltage controlled oscillator. Deviations of the oscillation frequency from the expected range of values indicate a defective Mixer. The simplicity of the proposed BiST scheme makes it an efficient solution for identification of defective Mixers (especially embedded ones in System-on-Chip applications) early in the production cycle (e.g. at the wafer level) reducing the total manufacturing cost.

The BIST circuit adopts the use of the Local Oscillator (LO) signal as test stimulus signal at the inputs of the Mixer. During the test operation, the signal input of the Mixer is disconnected from the signal driver (e.g. the LNA or baseband amplifier) and connected to the LO output with the use of proper analog switches. In Fig. 6 the above topology is illustrated for the case of an RF differential Mixer in a receiver.

The self-mixing of the LO signal forces the mixer to operate as a homodyne mixer (zero IF), generating at its “IF” outputs (IF+, IF-) a DC level (zero IF frequency) accompanied by the higher order mixing products. Simple RC Low-Pass Filters (LPFs) are added to reject these high frequency components and the DC outputs (VC+, VC-) are used as control signals to a ring voltage controlled oscillator (VCO) in order to control its oscillation frequency. Finally, the output signal of the VCO is used as the clock signal for a simple digital Counter. The above BiST scheme is based on the observation that the presence of a defect in the Mixer changes the DC levels of the IF outputs. This in turn will alter the oscillation frequency of the VCO from its nominal value in the defect free case. Consequently, the number of counts in the Counter, within a predetermined test phase time interval, will also deviate from the pertinent defect free value. In other words, the Counter’s value is considered as a test signature and in case that this deviates from its defect free value, the corresponding circuit under test is characterized as defective. The test signature can be exported outside the chip through a scan-out port (SO) for comparison. In case that the RF

front-end is embedded in a SoC with other digital circuits, the standard scan facilities of those circuits can be exploited.

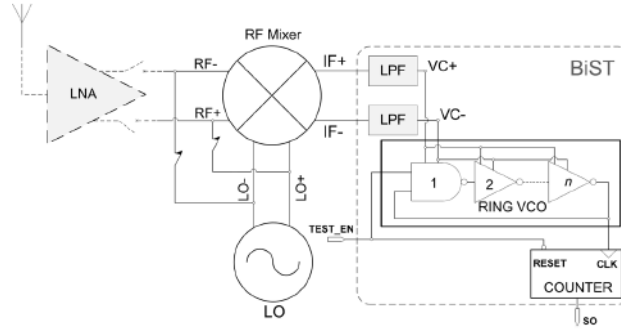


Fig. 6. Receiver's Mixer test architecture.

The proposed test strategy can be easily extended for testing both mixers in transceiver circuits using a single embedded BiST. In a similar approach as in the receiver's case above, a set of analog switches disconnect the inputs of the transmitter's Mixer from the outputs of the baseband amplifier and connect the LO outputs activating the test path. Then, the outputs of the transmitter's Mixer are multiplexed with the outputs of the receiver's Mixer, using a 2:1 differential multiplexer. The multiplexer output signals are low-pass filtered to be used as the control signals of the ring VCO.

The overall Mixer-BiST performance is summarized in Table 3. The silicon area of the BiST circuit is estimated to be 16% of the RF Mixer area.

Table 3. Fault coverage results for the Mixer-BiST.

Type of Defect	Fault Coverage
Resistive Shorts	16/22
Resistive Bridgings	39/41
High Ohmic Resistive Opens	28/28
Overall Fault Coverage	83/91 (91.2%)

3 RF Front-End Test Architecture

The overall architecture for testing wireless transceiver's RF front-ends is illustrated in Fig. 7. According to this, each and every one of the participating circuits is tested individually, applying defect oriented test techniques like those presented earlier. The signal of the local oscillator is used as the test stimulus, minimizing the overall need for external signals and input pins. The circuits under test (LNA, $VCO_{S_{RX-TX}}$,

Mixers_{RX-TX}) are successively tested through a complete network of switches so that the test path can be fully isolated during the normal operation of the transceiver. For the LNA and the Mixers, the signal of the Local Oscillator (VCO) is used as test stimulus.

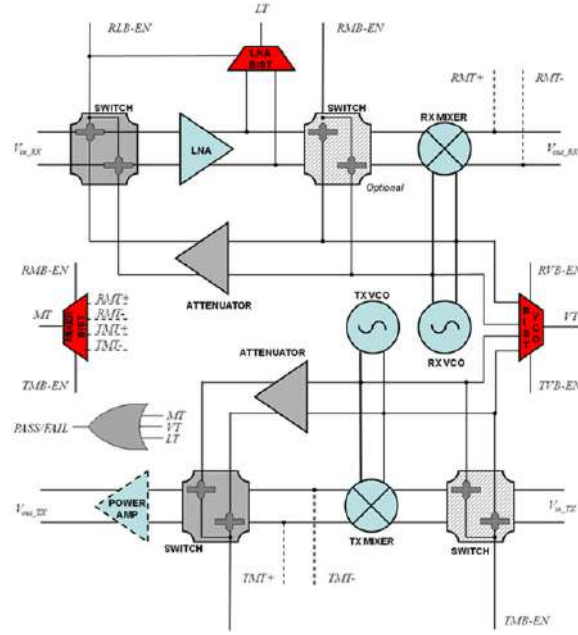


Fig. 7. Test architecture of a wireless transceiver RF front-end.

In more details, during the test phase a set of control signals (LB-EN, MB-EN, VB-EN) activates consecutively the test paths and the relevant build-in test circuits which are attached on the circuits under test that constitute the RF front-end system. These control signals are easily provided using a scan chain facility.

The test outputs (MT, VT, LT) drive a three-input OR gate which further generates a single digital test signal (*PASS/FAIL*). In case that all the individual test outputs are “0”, the system under test is considered as a fault free one. In the opposite case where at least one test signal is “1”, the system is characterized as faulty.

4 Contribution of the dissertation – Conclusions

The testing of RF circuits and systems is a challenging procedure, and in today’s nanoscale era it defines in great extend the overall manufacturing cost. In this dissertation we:

- Proposed a novel technique for testing wireless transceiver's RF front-end circuits, based on a defect oriented approach. The resulted overall architecture can efficiently utilize the available digital and analog system resources during the test phase (for cost reduction) and effectively distinguish the defective from the defect-free structures.
- Developed defect oriented design for testability circuits for the basic modules that constitute the RF front-end of the wireless transceivers (LNA, Mixers, VCOs).

Based on the results obtained from the analysis, this dissertation does not only answer the question "Can the RF system testing get more effective and cost efficient?" but it goes one step further, to the more fundamental, "Is it worthwhile to test RF systems?", and the answer is, yes.

As future work, a defect-oriented approach for testing transmitter's RF amplifiers may be considered as well as a technique for self-calibration of the defective circuits, in order to increase manufacturing yield.

References

1. Milor, L.S.: A Tutorial Introduction to Research on Analog and Mixed-Signal Circuit Testing. In: IEEE Transactions on Circuits and Systems – II: Analog and Digital Signal Processing, vol. 45, no. 10, pp. 1389–1407 (1998).
2. Akbay, S.S., Halder, A., Chatterjee, A., Keezer, D.: Low-Cost Test of Embedded RF/Analog/Mixed-Signal Circuits in SOPs. In: IEEE Transaction on Advanced Packaging, vol. 27, no. 2, pp. 352–363 (2004).
3. Variyam, P., Cherubal, S., Chatterjee, A.: Prediction of Analog Performance Parameters using Fast Transient Testing. In: IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems, vol. 21, no. 3, pp. 349–361 (2002).
4. Voorakaranam, R., Cherubal, S., Chatterjee, A.: A Signature Test Framework for Rapid Production Testing of RF Circuits. In: IEEE Design Automation and Test in Europe Conference, pp. 186-191 (2002).
5. Bushnell, M. L., Agrawal, V. D.: Essentials of Electronic Testing for Digital, Memory and Mixed-Signal VLSI Circuits. Kluwer Academic Publishers (2001).
6. Slamani, M., Kaminska, B.: Multifrequency Analysis of Faults in Analog Circuits. In: IEEE Design and Test of Computers, vol. 12, no. 2, pp. 70-80 (1995).
7. Huang, J.-L., Ong, C.-K., Cheng, K.-T.: A BIST Scheme for On-Chip ADC and DAC Testing. In: Proc. of the Design Automation and Test in Europe Conference (DATE), pp. 216–220 (2000).
8. Dabrowski, J., Gonzalez-Bayon, J.: Mixed Loop-back BIST for RF Digital Transceivers. In: Proc. of the International Symposium on Defect and Fault Tolerance in VLSI Systems (DFT), pp. 220–228 (2004).
9. Yoon, J.-S., Eisenstadt, W.R.: Embedded Loopback Test for RF ICs. In: IEEE Trans. on Instrumentation and Measurement, vol. 54, no. 5, pp. 1715–1720 (2005).
10. Craninckx, J., Steyaert, M.: A fully integrated CMOS DCS-1800 frequency synthesizer. In: Proc. of IEEE Dig. Tech. Papers, pp. 372–373 (1998).
11. Rofougaran, A., Rael, J., Rofougaran, M., Abidi, A.: A 900 MHz CMOS LC-oscillator with quadrature outputs. In: Proc. of IEEE Int. Solid-State Circuits Conference (ISSCC), pp. 392-393 (1996).

12. Dermentzoglou, L., Tsiatouhas, Y., Arapoyanni, A.: A Design for Testability Scheme for CMOS LC-Tank Voltage Controlled Oscillators. In: *Journal of Electronic Testing: Theory and Applications*, vol. 20, no. 2, pp. 133—142 (2004)
13. Dermentzoglou, L., Tsiatouhas, Y., Arapoyanni, A.: A novel scheme for testing radio frequency voltage controlled oscillators. In: *10th IEEE International Conference on Electronics, Circuits and Systems 2003 (ICECS 2003)*, pp. 595—598 (2003)
14. Dermentzoglou, L., Tsiatouhas, Y., Arapoyanni, A.: A built-in self-test scheme for differential ring oscillators”, 6th International Symposium on Quality of Electronic Design, 2005 (ISQED 2005), pp. 448—452 (2005)
15. Tang, J.J., Lee, K.J., Liu, B.D.: A Practical Current Sensing Technique for IDDQ Testing. In: *IEEE Transactions on VLSI Systems*, vol. 3, no. 2, pp. 302—310 (1995).
16. Variyam, P.N., Chatterjee, A.: Digital-Compatible BIST for Analog Circuits Using Transient Response Sampling. In: *IEEE Design & Test of Computers*, pp. 106—115 (2000).
17. Dermentzoglou, L., Arapoyanni, A., Tsiatouhas, Y.: A Built-In-Test Circuit for RF Differential Low Noise Amplifiers. In: *IEEE Transactions on Circuits and Systems I*, Vol 57, no. 7, pp. 1549-1558 (2010)
18. Dermentzoglou, L., Tsiatouhas, Y., Arapoyanni, A.: A Design for Testability Technique for Differential RF Low Noise Amplifiers. In: *XX Conference on Design of Circuits and Integrated Systems (DCIS 2005)*
19. Dermentzoglou, L., Tsiatouhas, Y., Arapoyanni, A.: An Embedded Test Circuit for RF Single Ended Low Noise Amplifiers. In: *14th IEEE International Conference on Electronics, Circuits and Systems, (ICECS 2007)*, pp. 1119—1122 (2007)
20. Dermentzoglou, L., Tsiatouhas, Y., Arapoyanni, A., Karagounis, A.: A Built-In Test Circuit for Single Ended RF Low Noise Amplifier. In: *17th North Atlantic Test Workshop (NATW 2008)*
21. Dermentzoglou, L., Tsiatouhas, Y., Arapoyanni, A.: A Build-In Self-Test Technique for RF Mixers. In: *13th IEEE International Symposium on Design and Diagnostics of Electronic Circuits and Systems (DDECS2010)*, pp. 88—92 (2010)

Development of Learning Environments with Use of Logo programming language in teaching praxis

Katerina Glezou*

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
kglezou@di.uoa.gr

Abstract. The thesis focuses on the development of learning environments with the use of Logo programming language in didactic praxis to promote exploratory and collaborative learning in primary and secondary education and teacher training. A framework is proposed for designing, developing and implementing scenarios and teaching activities by using microworlds developed in Logo-like environments as teaching-learning tools focusing on teaching Informatics and investigating their contribution to the teaching-learning process through authentic teaching examples. Activities and microworlds were designed and developed to support the teaching and learning process (i) promoting learning through exploration and collaboration and (ii) providing a scaffolding to those involved (teachers and learners) during the engagement with activities in classroom. Also online learning environments are proposed as alternatives online teaching suggestions in introducing Logo and Logo-like environments in modern Learning Management Systems / Course (LMS / LCMS) Moodle and LAMS. In particular, a) an online introductory Logo course on Moodle platform where the teaching-learning material for each level of education is grouped and b) sequences of learning activities on introduction to Logo into LAMS platform. Finally, the online learning environment of social networking “The Logo in education: A learning community of practice” is proposed, which is formed in the framework of development and operation of the homonymous educational online social network.

Keywords: Logo, Logo-like environments, Logo programming, microworlds, learning activities, microworlds, investigation, exploration, collaboration

1 Introduction

The introduction and exploitation of Information and Communication Technologies (ICTs) in classroom remains an open, composite and multifactoral issue. The creation of interesting and demanding environments encouraging the active and constructive participation of students is a great challenge for teachers. The planning of a

* Dissertation Advisor: Maria Grigoriadou, Professor

learning environment includes extensive decision making for planning, which should be the result of conscious thought rather than an unconscious choice ([27], [4], [6]). Learning occurs through a process of continuous changes in the individual's cognitive structures and is directly linked to the effects of the sociocultural environment [28]. At the same time, the context in which learning takes place [24], as well as the tools' mediation [28] play a crucial part providing opportunities for active, exploratory and personally significant learning for the individual.

Logo is considered an important tool in the hands of teachers and students for the development of their exploration skills, creativity skills and problem solving skills and for the cultivation of logical-algorithmic reasoning ([24], [16], [2], [18], [17], [3], [26], [5]). The trainees become, at the same time, users and designers as they design and construct tools and objects for the solution of problems. This double role of the trainees leads directly to the notion of constructionism. Constructionism involves two interweaving types of construction: knowledge construction through construction of artifacts with personal meaning ([16], [18]).

Logo-like environments can be used to plan and develop microworlds that offer students the possibility to express and exploit their thoughts, ideas and instincts and support the process of building knowledge by creating learning environments rich in speculation and opportunities for experimentation ([17], [26]).

The microworld concept has been present for over four decades now and the exploitation of microworlds in education has triggered the interest and attention of many researchers and instructors, who plan, experiment with and explore alternative constructionist approaches in various thematic fields ([18], [5], [6], [19], [20], [21], [22], [23], [1], [2], [25]).

In this context, the research focuses on the development of technologically supported learning environments through the use of Logo as a programming language and philosophy of education, that supports exploratory and collaborative learning in primary and secondary education and teacher training.

Central research topic was the design, development, implementation and evaluation of teaching scenarios, activities and microworlds developed in Logo-like learning environments to promote exploratory and collaborative learning and exploring their contribution to the teaching-learning process. Also, an important research topic was the development of online learning environments as alternative online teaching suggestions in introducing Logo and Logo-like environments towards promoting communication,

interaction and collaboration among members of the educational community.

The remainder of this paper is organized as follows. In Section 2, a framework for designing, developing and implementing scenarios and teaching activities by using microworlds developed in Logo-like environments is presented. In Section 3, online learning environments are presented as alternatives online teaching suggestions in introducing Logo and Logo-like environments in modern Learning Management Systems / Course (LMS / LCMS) Moodle and LAMS. Following, in Section 4, the online learning environment of social networking “The Logo in education: A learning community of practice” is presented. Finally, conclusions are given in Section 5, with the main points of the research and its contribution in the specific research area.

2 A framework for designing, developing and implementing scenarios and teaching activities by using microworlds developed in Logo-like environments

A framework is proposed for designing, developing and implementing scenarios and teaching activities by using microworlds developed in Logo-like environments (for example MicroWorlds Pro, Xelonokosmos/E-Slate) as teaching-learning tools focusing on teaching Informatics and investigating the contribution of these in the teaching-learning process through authentic teaching examples.

Educational scenarios, activities and microworlds were designed and developed to support the teaching - learning process in primary and secondary education and teacher training (i) promoting learning through exploration and collaboration and (ii) providing a scaffolding to those involved (teachers and learners) during the engagement with activities in classroom ([8], [11], [12], [14], [15]).

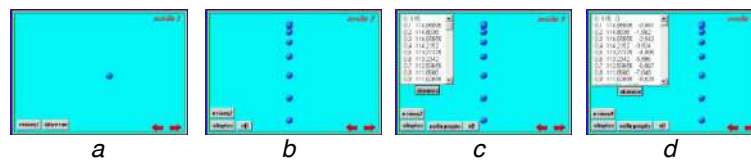


Fig. 1. Snapshots of the “Free fall simulation development” preconstructed microworld.

The basic axes for the design of educational scenarios and activities suggested are: a) structuring exploratory roles, b) supporting the process of active knowledge building, c) exploiting students' previous knowledge, experiences and intuitions, d) developing new student-teacher roles, e) creating collaborative learning environments, and f) using a cross-thematic approach.



Fig. 2. Snapshots of the “Free fall simulation development” students’ project work microworld.

The Investigation Course functions as a framework for the introduction and exploitation of microworlds in the classroom, focusing on alternative forms of exploration, knowledge structuring, expression, collaboration and communication for students and teachers [15]. The activities of the Investigation Course offer rich opportunities for experimentation, formulation and testing hypotheses, interpreting and shaping ideas by placing emphasis on the development of high level reasoning and problem solving skills.

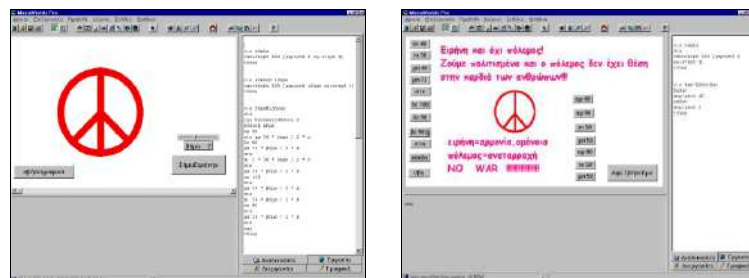


Fig. 3. Snapshots of “The Peace Symbol” students’ microworlds.

The following are stressed as particular focal points of the Investigation Course: a) emphasis is given to the process and not to the end product, b) students’ intuition is cultivated and exploited, c) students’ thoughts and ideas are visualized, d) students exploit their mistakes and are led to the depenalization of the mistake, e) new

problem solving strategies develop, such as the analysis of the problem in different parts, f) students make original artifacts of personal interest and meaning. Taking into consideration the particular students' previous knowledge and experience level, the starting point and the task's course are different each time, having as basic guiding axes the fact that we gradually move to the writing of an increasingly difficult code: a) to the familiarization with simple Logo commands, b) to using the simple and composite repetition command, c) to defining procedures, d) to defining superprocedures, e) to introducing the concept of variable and the definition of parametric procedures and f) to defining parametric Logo superprocedures.



Fig. 4. Snapshots of a microworld, while working out activity with a gradually increasing degree of complexity (Stages 2-3).

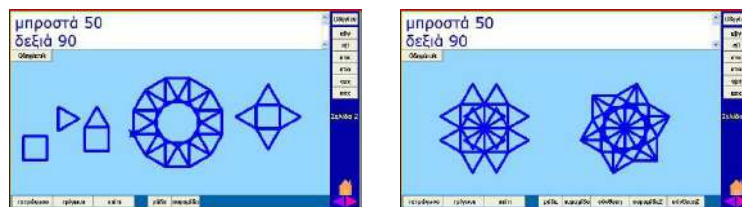


Fig. 5. Snapshots of a microworld, while working out activity with a gradually increasing degree of complexity (Stages 4-5).

It also proposes an alternative constructionist training approach for the introduction to Logo programming by using a structured series of activities and preconstructed reusable microworlds developed in the multimedia programming environment MicroWorlds Pro. The proposed approach and material aims to scaffold the gradual familiarization of the trainees with Logo –programming language and philosophy of education- and the programming environment by applying in action the constructionist reasoning ([9], [10], [13]). The training material is characterized by a gradual increase in complexity

and difficulty degree, and exploits the gradually acquired experience of the trainees by providing them with tools that they are in position to use.

The preconstructed microworlds functioned as “objects to think with”, a good starting point and a solid ground for explorations-modifications-extensions, as a vehicle for collaboration and led to various alternative constructions of personal and social meaningful artifacts. The use and re-use of preconstructed microworlds for the creation of new artifacts with a gradually increasing degree of complexity encourages the systematization of knowledge and bridges the gap between the simple and the more complex ([8], [9], [10], [13]).

3 Development of online learning environments

Online learning environments are proposed as alternatives online teaching suggestions in introducing Logo and Logo-like environments in modern Learning Management Systems / Course (LMS / LCMS) Moodle and LAMS. In particular, a) an online introductory Logo course on Moodle web platform where is grouped the teaching-learning material for each level of education and b) sequences of learning activities on introduction to Logo into LAMS platform.

The special features of online learning environments create new conditions for learning and demonstrate a variety of new possibilities for alternative forms of communication, interaction and cooperation by supporting collaborative learning. Both online learning environments are emerging as convenient and friendly web development tools of learning environments that support active participation and promote communication, interaction and collaboration between stakeholders (teachers and learners). They are different platforms with specific features each and a total benchmarking is not feasible. The teacher emerges as the catalyst acting on the teaching strategy that will follow and support through the tool and in this way will add value to the pedagogical use of technology. Nevertheless it is considered that by using the Learning Management System / Courses Moodle more emphasis is given on content delivery, while by using the Learning Activity Management System LAMS more emphasis is placed on interaction and cooperation.

The network primarily concerns teachers of Informatics and Computer Science and in parallel, teachers of various specialties, cognitive subjects and all educational levels who are interested in or/and experimenting with the usage of Logo programming language in the teaching praxis. As it is denoted in the “LogoinEdu” subtitle “Learn - Construct - Collaborate - Communicate” the ulterior objective of “LogoinEdu” is to function as a learning community of practice, as a forum for the dialogue and mutual support between members of the educational community focusing on the pedagogical exploitation of Logo and Logo-like environments attempting to improve the teaching-learning process. The network members are invited to interact in the spirit of Social Constructionism: “Let’s function as a community of practice and learning and exchange views, experiences, practices and tools, such as microworlds, websites, lesson plans, worksheets, codes and all kinds of resources necessary for our teaching practice, with the purpose to upgrade the teaching-learning process.” as it is characteristically mentioned in the network pages.



Fig. 8. Snapshot of home page of the “Logo in education: A learning community of practice” educational online social network.

5 Conclusions

The research presented contributes to the fields of didactics of informatics, and especially of didactics of Logo programming as well as of computer-supported collaborative learning. The main

contribution of the work lies in the provision of a framework and in the development of Logo-based learning environments that support the construction of knowledge and promote synchronous and asynchronous communication and collaboration.

It proposes a framework for designing, developing and implementing scenarios and teaching activities by using microworlds developed in Logo-like environments as teaching-learning tools focusing on teaching Informatics and investigating the contribution of these to the teaching-learning process through authentic teaching examples. Activities and microworlds were designed and developed to support the teaching and learning process (i) promoting learning through exploration and collaboration and (ii) providing a scaffolding to those involved (teachers and learners) during the engagement with activities in classroom.

Online learning environments are proposed as alternatives online teaching suggestions in introducing Logo and Logo-like environments in modern Learning Management Systems: a) an online introductory Logo course on Moodle platform where is grouped the teaching-learning material for each level of education and b) sequences of learning activities on introduction to Logo into LAMS platform. Finally, the online learning environment of social networking "The Logo in education: A learning community of practice is proposed. The network is functioning as a learning community of practice and as a step for dialogue and mutual support of the educational community in an effort to upgrade the teaching - learning process.

The studies conducted, revealed encouraging and positive results for the above mentioned environments in serving their underlying objectives and in supporting the learning process.

The structured teaching/training material exploited in gradual steps and according to the acquired experience of the students/trainees could be considered especially effective in introducing Logo programming and in gradual familiarization with the programming environment; it may be adapted/extended to individual needs and may be used in different learning contexts.

The preconstructed microworlds functioned as a good starting point, as a solid ground for explorations-modifications-extensions, as a vehicle for collaboration and led to various alternative constructions of personal and social meaningful artifacts. The use of preconstructed microworlds for the construction of new artifacts with a gradually increasing degree of complexity encourages the systematization of knowledge and bridges the gap between the simple and the more complex.

The proposed framework and corresponding learning environments support exploratory and collaborative learning and contribute to the active involvement of stakeholders (teachers and learners), the construction of knowledge in programming concepts and cultivation of programming, expression and collaboration skills.

References

- [1] Brouwer, N., Muller, G. & Rietdijk, H. (2007). Educational Designing with MicroWorlds. *Journal of Technology and Teacher Education*, 15 (4), pp. 439-462. Chesapeake, VA: AACE.
- [2] Clements, D. H., & Meredith, J.S. (1993). Research on Logo: Effects and efficacy. *Journal of Computing in Childhood Education*, 4, 263-290.
- [3] Dagiene, V. (2003). *A set of Logo problems for learning algorithms*. In Proceedings of Eurologo 2003. Edited by Cnotinfor, Lda. Porto, August.168-177.
- [4] Dimitracopoulou, A. & Komis, V. (2005). Design principles for the support of modelling and collaboration in a technology-based learning environment. *Int. J. Cont. Engineering Education and Lifelong Learning*, Vol. 15, Nos. 1/2, 30–55.
- [5] diSessa, A. (1995). Epistemology and Systems Design, In diSessa, A. - Hoyles C., *Computers and Exploratory Learning*, Springer Verlag, 15-29.
- [6] diSessa, A. (2000). Changing minds: Computers, learning, and literacy. Cambridge, MA: MIT Press.
- [7] Glezou, K., Grigoriadou M., & Samarakou, M., (2010). Educational Online Social Networking in Greece: A Case Study of a Greek Educational Online Social Network. *The International Journal of Learning*, Volume 17, Issue 3, pp. 399-420.
- [8] Glezou, K. & Grigoriadou M., (2010). Engaging Students of Senior High School in Simulation Development. *INFORMATICS IN EDUCATION*, 2010, Vol. 9, No. 1, pp. 37-62.
- [9] Glezou, K. & Grigoriadou M., (2010). Teacher Training in Logo Programming by using Preconstructed Reusable Microworlds. *The International Journal of Learning*, Volume 17, Issue 1, pp. 347-364.
- [10] Glezou, K. & Grigoriadou, M. (2009). An Alternative Instructional Approach for Introductory Courses to Logo Programming. In *Proceedings of IADIS International Conference CELDA 2009*, pp. 419-424. Rome, Italy.
- [11] Glezou, K. & Grigoriadou, M. (2009). Supporting Student Engagement in Simulation Development. In C. O'Malley, D. Suthers, P. Reimman, A. Dimitracopoulou (Eds.) *Proceedings of 8th International Conference on Computer Supported Collaborative*

Learning CSCL2009: Computer Supported Collaborative Learning Practices, pp. 414-418. Rhodes.

- [12] Glezou, K. & Grigoriadou M. (2009). Design Principles of Training Material for Introductory Courses to Programming and Logo by using preconstructed microworlds. In *Proceedings of World Conference on Educational Multimedia, Hypermedia and Telecommunications 2009 (ED-MEDIA 2009)*, pp. 1606-1614. Chesapeake, VA: AACE.
- [13] Glezou, K. & Grigoriadou, M. (2008). Simulation Development by Students: An Alternative Cross-Thematic Didactical Approach. In *Proceedings of World Conference on Educational Multimedia, Hypermedia and Telecommunications 2008 (ED-MEDIA 2008)*, pp. 4108-4117. Chesapeake, VA: AACE.
- [14] Glezou, K. & Grigoriadou, M. (2007). A novel didactical approach of the decision structure for novice programmers. In Ivan Kalas (ed.) *Proceedings of 11th European Logo Conference (Eurologo 2007)*, Bratislava.
- [15] Glezou, K., Grigoriadou M. & Verginis, I. (2009). Rethinking the "Investigation Course" in Primary School. In Ignacio Aedo, Nian-Shing Chen, Kinshuk, Demetrios Sampson, Larissa Zaitseva (Eds.) *Proceedings of 9th IEEE International Conference on Advanced Learning Technologies (ICALT 2009)*, 554-555, Riga, Latvia.
- [16] Harel, I. & Papert, S. (1991). Constructionism: Research Reports & Essays, 1985-1990 by the Epistemology & Learning Research Group. Norwood: Ablex Publishing Corporation, US.
- [17] Hoyles, C., Noss, R. & Adamson, R. (2002). Rethinking the microworld idea. *Journal of Educational Computing Research*, 27(1&2), pp. 29-53.
- [18] Kafai, Y. & Resnick, M. (Eds.). (1996). Constructionism in practice: Designing, thinking, and learning in a digital world. Mahwah, NJ: Lawrence Erlbaum Associates.
- [19] Kalas, I. (2006). Discovering Informatics Fundamentals Through Interactive Interfaces for Learning. In R. T. Mittermeir (Ed.), *ISSEP 2006, LNCS4226*, pp. 13-24.
- [20] Komis, V. (2005). *Introduction in Didactics of Informatics*. Athens: Kleidarithmos Publications. (In Greek).
- [21] Kynigos, C. (2007). Half-baked Microworlds in use in Challenging Teacher Educators' Knowing, *International Journal of Computers for Mathematical Learning*. Kluwer Academic Publishers, Netherlands, 12 (2), 87-111.
- [22] Kynigos, C. (2007). Half-Baked Logo Microworlds as Boundary Objects in Integrated Design, *Informatics in Education*, 2007, Vol. 6, No. 2, 1-24, Institute of Mathematics and Informatics, Vilnius.
- [23] Louca, L., Druin, A., Hammer, D. & Dreher, D. (2003). Students' collaborative use of computer-based programming tools in science: A Descriptive Study. In B. Wasson, St. Ludvigsen, & Ul. Hoppe (Eds.). *Designing for change in Networked Learning Environments: Proceedings of the CSCL 2003* (pp. 109-118). The Netherlands: Kluwer Academic Publishers.

- [24] Papert, S. (1980). *Mindstorms: Children, Computers, and Powerful Ideas*. Basic Books, New York.
- [25] Rieber, L.P. (2004). Microworlds. In *Handbook of research for educational communications and technology (2nd ed.)*, D. Jonassen (Ed.), Mahwah, NJ: Lawrence Erlbaum Associates, 583-603.
- [26] Resnick, M., Maloney, J., Monroy-Hernández, A., Rusk, N., Eastmond, E., Brennan, K., Millner, A., Rosenbaum, E., Silver, J., Silverman, B., Kafai, Y., (2009). Scratch: Programming for All, November 2009, *Communications of the ACM*, 52(11), pp. 60-67.
- [27] Vosniadou, S. (2005). *Planning learning environments supported by modern technologies*, Athens, Gutenberg Publications. (In Greek).
- [28] Vygotsky, L.S. (1978). *Mind in Society: The development of Higher Psychological Processes*, Harvard University Press, Cambridge, Massachusetts.

Effective Capacity Theory for Modeling Systems with Time-Varying Servers, with an Application to IEEE 802.11 WLANs

Emmanouil N. Kafetzakis*

National Center for Scientific Research “Demokritos”
Institute of Informatics & Telecommunications
Department of Informatics & Telecommunications
National & Kapodistrian University of Athens
mkafetz@di.uoa.gr

Abstract. Many network applications rely on stochastic QoS guarantees. With respect to loss-related performance, the Effective Bandwidth/Capacity theory has proved useful for calculating loss probabilities in queues with complex input- and server-processes and for formulating simple admission control tests to ensure associated QoS guarantees. This success has motivated the application of the theory for delay-related QoS too. However, up to now this application has been justified only heuristically for queues with variable service rate. The thesis fills this gap by establishing rigorously that the Effective Bandwidth/Capacity theory may be used for the asymptotically correct calculation and enforcement of delay tail-probabilities in systems with variable rate servers too. Subsequently, the thesis applies the general results to IEEE 802.11 WLANs, by representing each IEEE 802.11 station as an On/Off server and employing the Effective Capacity function for this model. Comparison of analytical results with simulation validates the effectiveness of the On/Off IEEE 802.11 model for loss- and delay-related QoS. Finally, the thesis uses the Effective Capacity of an IEEE 802.11 station in yet another way, namely as a design tool. Indeed, the specific form of the IEEE 802.11 Effective Capacity function highlights the role of certain parameters of the IEEE 802.11 backoff window distributions. These parameters, when appropriately tailored, allow better delay-related (and loss-related) performance, while maintaining the standard saturation throughput of IEEE 802.11 WLANs.

Keywords: admission control, Effective Capacity, IEEE 802.11, quality of service, server modeling, tail-probabilities

1 Introduction

Many demanding network applications rely on stochastic Quality of Service (QoS) guarantees. With respect to loss-related performance, the asymptotic theory based on the notions of Effective Bandwidth and Effective Capacity has proved successful for

* Dissertation Advisors: Kimon Kontovasilis, Research Director N.C.S.R. “Demokritos” and Ioannis Stavrakakis, Professor N.K.U.A.

calculating low loss probabilities in queueing systems with complex time-varying input and server processes and for formulating simple admission control tests to enforce associated QoS guarantees (see, e.g., [12, 4, 5]).

This success has motivated the application of the theory to the calculation and enforcement of delay-related QoS too. However, up to now this application has only been justified on the basis of heuristic arguments when the queue is served at a variable rate (see, e.g., [1, 14]). The thesis fills this gap, by formally establishing that the Effective Bandwidth/Capacity theory may be applied for asymptotically correct calculation and/or enforcement of delay tail-probabilities in systems with variable rate servers too. In particular, the heuristically suggested linkage between the exponential decay rates of the buffer content and delay probability tails through the server's Effective Capacity function is formally shown to apply [7].

Due to the prevalence of wireless networking, systems with time-varying servers are becoming all the more important. Indeed, a wireless station can be regarded as a time-varying data server, due to rate fluctuations at the physical or at the medium access control layer. In this context, the thesis proceeds with an application of the general results to IEEE 802.11 WLANs. In doing so, the thesis first establishes that an IEEE 802.11 mobile station can be regarded as a Semi-Markovian data server of the On/Off type, with known distributions for the On and Off periods, and subsequently derives the Effective Capacity function of this On/Off server [8]. The general results can then be used for computing buffer overflow and delay violation probabilities in WLANs, and for employing simple traffic control policies to enforce related QoS guarantees.

Finally, the thesis illustrates the usage of the Effective Capacity function of the IEEE 802.11 stations as a design tool: Towards this end, the form of the said function highlights certain parameters of the backoff window distributions, which, if appropriately tailored, may lead to higher Effective Capacity values, hence to better delay-related (or loss-related) performance.

The rest of this thesis summary is organized as follows: Section 2 firstly reviews pre-existing Eff. Bandwidth/Capacity results about buffer content tail-probabilities and then justifies the applicability of Eff. Bandwidth/Capacity theory in connection with delay tail-probabilities. Section 3 briefly presents the analytical model used to characterize the IEEE 802.11 Distributed Coordination Function (DCF) as an On/Off server. This model is used to derive the Eff. Capacity of the IEEE 802.11 protocol. Subsequently, it discusses computational and algorithmic issues related to the application of the general theory of Section 2 with the particular Eff. Capacity function of this On/Off model. Section 4 describes how the Eff. Capacity function of IEEE 802.11 stations may be used for an informed choice of parameters for the backoff window distributions, towards better performance. Section 5 provides validation of the IEEE 802.11 model, through comparison of the analytical results with simulations. Finally, the thesis summary is concluded in Section 6.

2 Effective Bandwidth and Effective Capacity Theory

Effective Bandwidth and Effective Capacity theory offers a linkage between input load, system capacity and QoS requirements and it was developed by a great number of con-

tributions from various researchers. The theory was originally developed for queueing systems with constant server capacity.

When the server's capacity is time-varying independently from the input, the theory can be generalized, by defining an *Effective Capacity function* to capture the server's burstiness. Although this generalization has been studied for some years (see, e.g., [5, 6]) it did not attract much attention until recently [13, 14, 1] when the importance of wireless systems grew considerably. This is because most such systems feature a variable service rate and Effective Capacity is ideal for modeling such settings.

For a quick review of the Eff. Bandwidth/Capacity theory, consider at first a single-server queue, fed by a traffic stream that produces an amount of data $V(t)$ within a time-window $(-t, 0]$ and let c be the constant service rate. According to the Eff. Bandwidth theory, provided that $V(t)$ has stationary increments and satisfies some additional mild technical conditions (see, e.g., [7]), the probability that the queue size Q exceeds a certain threshold b has at all times an exponential upper bound of rate $\theta > 0$. Specifically, the tail of the queue-length distribution satisfies

$$a_V(\theta) < c \Rightarrow \lim_{b \rightarrow \infty} b^{-1} \log \Pr\{Q > b\} \leq -\theta,$$

where

$$a_V(s) \triangleq \frac{1}{s} \lim_{t \rightarrow \infty} \frac{1}{t} \log \mathbb{E} \left[e^{sV(t)} \right], \quad s \in \mathbb{R}, \quad (1)$$

determines the Effective Bandwidth function.

Now consider a single-server queue with time-varying capacity, where the capacity fluctuations are independent from the input. Let the input traffic be as previously and denote by $C(t)$ the amount of data that can be processed within the time-window $(-t, 0]$. Assuming the same technical conditions for the input and output processes (in particular stationary increments), the probability that the queue size Q exceeds a certain threshold b has at all times an exponential upper bound of rate $\theta > 0$, viz.,

$$a_V(\theta) < a_C(-\theta) \Rightarrow \lim_{b \rightarrow \infty} b^{-1} \log \Pr\{Q > b\} \leq -\theta, \quad (2)$$

where $a_V(\theta)$ stands for the Effective Bandwidth function in (1) and

$$a_C(s) \triangleq \frac{1}{s} \lim_{t \rightarrow \infty} \frac{1}{t} \log \mathbb{E} \left[e^{sC(t)} \right], \quad s \in \mathbb{R}, \quad (3)$$

determines the Effective Capacity function.

Given that *strict* monotonicity holds for at least one of the Eff. Bandwidth and Eff. Capacity functions,

$$a_V(\theta) \leq a_C(-\theta) \Leftrightarrow \lim_{b \rightarrow \infty} b^{-1} \log \Pr\{Q > b\} \leq -\theta. \quad (4)$$

Equivalence (4) always holds in the context of IEEE 802.11 WLAN, since the IEEE 802.11 Eff. Capacity is always a strictly increasing function.

Assume now that the queue has a finite size q , and that we want to provide stochastic QoS by limiting the overflow probability to a value $\leq e^{-\epsilon}$. By using the tail percentile $\Pr\{Q > q\}$ of the respective infinite queue as a proxy for the overflow probability, (2) implies that (asymptotically) $\theta = \epsilon/q$ and the input traffic must satisfy

$$a_V(\epsilon/q) < a_C(-\epsilon/q).$$

This inequality directly suggests an Admission Control (AC) scheme by bounding the input traffic.

Finally, let $\theta^* \triangleq \sup\{\theta \in \mathbb{R} \mid a_V(\theta) \leq a_C(-\theta)\}$. If there exist $\theta_o > 0$ such that $a_V(\theta_o) < a_C(-\theta_o)$ then the asymptotic exponential decay rate of the tail-probabilities of Q equals to θ^* [7], i.e.,

$$\lim_{b \rightarrow \infty} b^{-1} \log \Pr\{Q > b\} = -\theta^*. \quad (5)$$

The conceptual simplicity of the Eff. Bandwidth/Capacity theory makes it an attractive choice for coping with delay-related QoS as well. For First-Come-First-Served (FCFS) queueing systems with a constant service rate c this is directly possible, because delay probabilities of the form $\Pr\{D > d\}$ are equal to the queue length probabilities $\Pr\{Q > cd\}$. However, this simple equivalence does not hold when the service rate is time-varying.

The thesis (see the relevant results in [7]) formally establishes that the Eff. Bandwidth/Capacity theory may be applied for the asymptotically correct calculation and enforcement of delay tail-probabilities in the general setting with variable service rate. Specifically,

$$a_V(\theta) < a_C(-\theta) \Rightarrow \lim_{d \rightarrow \infty} d^{-1} \log \Pr\{D > d\} \leq -\theta a_C(-\theta). \quad (6)$$

As with the queue content, we now consider admission control for ensuring delay-related QoS guarantees. Let $u_C(s) \triangleq \lim_{t \rightarrow \infty} t^{-1} \log \mathbb{E}[e^{sC(t)}]$, $s \in \mathbb{R}$, be the asymptotic cumulant generator of $C(t)$. By the definition of the Eff. Capacity function in (3),

$$a_C(s) = u_C(s)/s, \quad s \in \mathbb{R}. \quad (7)$$

In order to ensure that the decay rate of the delay tail-probabilities is bounded below by some $\xi = \theta a_C(-\theta) = -u_C(-\theta) > 0$, the admission control $a_V(\theta) < a_C(-\theta)$ in (6) must be tested for

$$\theta(\xi) = -u_C^{-1}(-\xi). \quad (8)$$

The value of the parameter ξ to employ in the tests is determined in a way analogous to the one used for loss-related QoS requirements. This time the QoS specification dictates that the delay should not exceed some given threshold τ with probability higher than $e^{-\epsilon}$. Provided that both τ and ϵ are large maintaining a finite ratio, the QoS specification leads to $-\epsilon/\tau \geq \tau^{-1} \log \Pr\{D > \tau\} \approx \lim_{d \rightarrow \infty} d^{-1} \log \Pr\{D > d\}$, so $\xi = \epsilon/\tau$ should be used in the admission control tests.

Finally, the thesis establishes the linkage between the decay rate θ^* of the buffer content tail-probabilities and the decay rate ξ^* of the delay tail-probabilities through the server's Effective Capacity function, viz.,

$$\xi^* = \theta^* a_C(-\theta^*).$$

3 Effective Capacity of IEEE 802.11 WLAN

A simple, but accurate analytical model for the saturation throughput computation of the IEEE 802.11 protocol was provided in [2]. The analysis focuses on the saturation

condition, where every station has always a packet to send. The analysis also assumes that the number of stations n under contention is known and constant and the probability of a collision seen by a packet being transmitted on the channel, named conditional collision probability p , is constant and independent from the number of retransmission suffered. In order to compute the station's throughput, its behaviour is studied through a Markov chain model. This model yields the probability that a station transmits in a random time-slot, referred to as transmission probability τ .

In [3], a more general model, permitting the usage of arbitrary backoff window distributions, is proposed, generalizing and supplementing [2]. This model takes into account more details of the IEEE 802.11 protocol. By employing the more general analysis of [3] and assuming an infinitive number of retries, one obtains

$$\tau = \left[1 + (1 - p) \left(\frac{\bar{W}_0}{1 - B_0} - 1 + \sum_{i=1}^{m-1} p^i \bar{W}_i + \frac{p^m \bar{W}_m}{1 - p} \right) \right]^{-1} \quad (9)$$

where \bar{W}_i , $i = 0, \dots, m$ is the mean backoff window at the i^{th} stage, m is the back-off stage beyond which the upper window margin does not grow anymore and B_0 is the probability that a backoff window drawn at the 0^{th} stage is zero. The expression $\bar{W}_0/(1 - B_0) - 1$ is a modified mean backoff window at 0^{th} stage, when the nonzero backoff window drawn is examined after being initially decremented by one, for synchronization purposes.

Assuming that the Markov chains of the mobile stations are independent, one also obtains that

$$1 - p = (1 - \tau)^{n-1}, \quad (10)$$

because a packet will not suffer a collision exactly when all other stations do not attempt to transmit when the station emitting the packet does so. Equations (9) and (10) can be solved uniquely for p and τ .

As already noted, the preceding analysis assumes that all stations are saturated. This thesis employs the values of p and τ obtained from (9) and (10), to calculate the Eff. Capacity of an IEEE 802.11 station, effectively assuming that all other stations are saturated. This approximation is on the safe side (i.e., 'conservative') since assuming the other stations saturated corresponds to the worst case, and has the merit that in this way the Eff. Capacity of a station can be computed without regard to input traffic details. As will be discussed later, the saturation assumption can be waived and the model be applicable to all network load settings.

Due to the CSMA/CA access algorithm, the system can be modeled as a Semi-Markov server model featuring four states [8], as depicted in Fig. 1. State `bc` corresponds to the backoff procedure when the backoff counter is nonzero. State `ov` models overhead time before and after the transmission. It can be proved that the overhead time before and after transmission is possible to be merged in one state. State `tr` corresponds to the transmission (active) period and State `dc` models the idle slot needed for the initial decrement of the backoff window so that other stations realize that a successful transmission is over.

Overheads and transmissions always occur in pairs thus transitions from State `ov` to State `tr` occur with probability one. After the transmission is over, State `ov` is visited

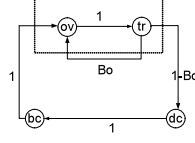


Fig. 1. Semi-Markov chain for IEEE 802.11 MAC.

again with probability B_0 (probability that backoff counter at the 0th stage is zero) so the station transmits a packet successfully, one more time. With probability $1 - B_0$ the station enters the backoff procedure (State bc) after the backoff counter has been decremented by one (State dc). The service rate in State tr is equal to the nominal bit rate \hat{r} of the IEEE 802.11 channel, while in all other states the service rate is zero.

Sojourn time distributions in every state are characterized by the respective moment generators. The moment generator of the sojourn time in State tr (i.e., payload transmission time) is

$$\gamma_{tr}(\omega) \triangleq \mathbb{E} \left[e^{P/\hat{r}} \right],$$

where P denotes the payload size. Note that for constant payload size the transmission time is deterministic. The moment generators of time distributions for States ov and dc are given by

$$\gamma_{ov}(\omega) \triangleq e^{\omega t_{ov}}, \quad \gamma_{dc}(\omega) \triangleq e^{\omega t_{slot}}, \quad (11)$$

where

$$t_{ov} \triangleq (RTS + CTS + PHY_{hdr} + ACK)/r_{signal} + MAC_{hdr}/\hat{r} + 3SIFS + DIFS,$$

and

$$t_{coll} \triangleq RTS/r_{signal} + EIFS + t_{slot}$$

are the deterministic overhead time of the transmission and the collision duration, respectively. The quantities r_{signal} , t_{slot} , RTS , CTS , $SIFS$, $DIFS$, $EIFS$, ACK , MAC_{hdr} , PHY_{hdr} denote respectively the transmission rate used for signaling operations, the slot time of the system, the RTS packet size, the CTS packet size, the SIFS time, the DIFS time, the EIFS time, the ACK packet size, the MAC header and the PHY header. The values of all these parameters are determined by the IEEE 802.11 standard [11].

Finally, the backoff time distribution is characterized by the moment generator

$$\gamma_{bc}(\omega) = \frac{g_0(\gamma_s(\omega)) - B_0}{\gamma_s(\omega)(1 - B_0)} \left[\sum_{l=0}^{m-1} \left((1-p)p^l e^{l\omega t_{coll}} \prod_{j=1}^l g_j(\gamma_s(\omega)) \right) + \frac{(1-p)(pe^{\omega t_{coll}})^m \prod_{j=1}^m g_j(\gamma_s(\omega))}{1 - pg_m(\gamma_s(\omega))e^{\omega t_{coll}}} \right], \quad (12)$$

where $g_i(z)$ stands for the probability generator function of the backoff window distribution associated with the i^{th} stage (beyond stage m , the backoff windows maintain the

same probability generator $g_m(z)$ and $\gamma_s(\theta)$ denotes the moment generator of the time needed for the reduction by one of the backoff counter.

$$\gamma_s(\omega) = P_{\text{coll}}e^{\omega t_{\text{coll}}} + P_{\text{empty}}e^{\omega t_{\text{slot}}} + P_{\text{succ}} \frac{(1 - B_0)\gamma_{\text{tr}}(\omega)}{1 - B_0\gamma_{\text{tr}}(\omega)} e^{\omega t_{\text{slot}}}, \quad (13)$$

where

$$P_{\text{succ}} = (n-1)\tau(1-\tau)^{n-2}, \quad P_{\text{empty}} = (1-\tau)^{n-1}, \quad P_{\text{coll}} = 1 - P_{\text{succ}} - P_{\text{empty}}, \quad (14)$$

are the probabilities of a successful transmission, an empty slot and a collision respectively, observed by a station backing-off (which observes $n-1$ other stations).

Note that the generator function $g_i(z)$ depends on the distribution of the backoff window at the i^{th} stage. By definition, $g_0(0) = B_0$ and $g'_i(1) = \bar{W}_i$. For the uniform backoff window distribution, described in the standard,

$$g_i(z) \triangleq \sum_{l=0}^{w_i-1} \frac{1}{w_i} z^l = \frac{1}{w_i} \frac{z^{w_i} - 1}{z - 1}, \quad w_i = 2^{\min\{i, m\}} w_0, \quad i \geq 0, \quad (15)$$

where $w_0 - 1$ is the maximum value of the backoff counter at the backoff stage zero.

Given the model just described, a straightforward extension of the Eff. Bandwidth theory for Semi-Markovian models [9] yields the Eff. Capacity function. In fact, it is possible to show that, in terms of the Eff. Capacity, the model is equivalent with an On/Off server model [8] with an On period characterized by the moment generator

$$\gamma_{\text{on}}(\omega) = \gamma_{\text{tr}}(\omega) = \mathbb{E} \left[e^{\omega P/\hat{r}} \right] \quad (16)$$

and an Off period whose moment generator is equal to

$$\gamma_{\text{off}}(\omega) = \gamma_{\text{ov}}(\omega) \left(B_0 + (1 - B_0)\gamma_{\text{bc}}(\omega)\gamma_{\text{dc}}(\omega) \right) \quad (17)$$

where $\gamma_{\text{on}}(\cdot)$, $\gamma_{\text{ov}}(\cdot)$, $\gamma_{\text{bc}}(\cdot)$, $\gamma_{\text{dc}}(\cdot)$ as in (16), (11), and (12).

Using this alternative On/Off representation the Eff. Capacity function is given by (7) where $u_C(s)$ is the unique negative solution of

$$f(s, u_C(s)) = 0, \quad f(s, u) \triangleq \log \gamma_{\text{on}}(\hat{r}s - u) + \log \gamma_{\text{off}}(-u) = 0. \quad (18)$$

The formulation of (9) assumes saturation condition. The dependence of $\gamma_{\text{off}}(\cdot)$ on the saturation assumption is only through the conditional collision probability p , used in (12) and the probabilities P_{succ} , P_{empty} and P_{coll} employed by (13). Under non-saturation condition these parameters retain their meaning, but take different values. Thus, if each mobile station assesses these probabilities by direct measurement, rather than computing them through (9), (10) and (14), the model works well in all settings, lightly loaded ones included.

For the construction of a loss-related traffic control mechanism is not necessary to numerically solve (18). Using the admission control condition $a_V(\theta) \leq a_C(-\theta)$, $\theta \geq 0$

and the monotonicity of the related functions, we can prove that in order to accept a traffic stream the following inequality must be satisfied:

$$\log \gamma_{\text{on}}(-\hat{r}\theta^* + \theta^* a_V(\theta^*)) + \log \gamma_{\text{off}}(\theta^* a_V(\theta^*)) \leq 0,$$

where $\theta^* = \epsilon/q$, according to (2). This condition simplifies greatly the computational aspects of the loss-related AC scheme.

Calculations for the delay-related AC test are also simple. According to the results of Section 2, given the QoS specification ξ , one must first determine $\theta(\xi)$ in (8) and then check if the left-hand side inequality in (6) holds. Since $u_C(-\theta(\xi)) = -\xi$, (18) suggests that $\theta(\xi)$ is the unique solution in θ of $f(-\theta, -\xi)$, thus

$$\theta(\xi) = \xi/\hat{r} - (\log \gamma_{\text{on}})^{-1}(-\log \gamma_{\text{off}}(\xi))/\hat{r}. \quad (19)$$

This requires only a single evaluation of the function $\gamma_{\text{off}}(\cdot)$ at the argument ξ , keeping the computational complexity low. Moreover, when the payload of the transmitted packets has a constant value P , (16) yields $(\log \gamma_{\text{on}})^{-1}(x) = \hat{r}x/P$, so (19) simplifies further to the closed form solution $\theta(\xi) = \xi/\hat{r} + \log \gamma_{\text{off}}(\xi)/P$. Note that, as long the conditions¹ in the WLAN remain unchanged, a single evaluation of $\theta(\xi)$ suffices to enable an arbitrary number of admission control tests (6), each of them being invoked whenever the mobile station is about to engage a new traffic flow.

4 Tuning the backoff window distributions for improved Effective Capacity

We now investigate appropriate choices of the backoff window distributions employed by the IEEE 802.11 MAC protocol, so as to obtain an Eff. Capacity function greater than the one corresponding to the standard distributions. A greater Eff. Capacity function signifies improved performance.

The mean rate of the On/Off model for an IEEE 802.11 mobile station is

$$\bar{r}_C \triangleq u'_C(0) = a_C(0) = \frac{\hat{r}E[T_{\text{on}}]}{E[T_{\text{on}}] + E[T_{\text{off}}]}. \quad (20)$$

In view of (20), one might attempt to obtain a greater Eff. Capacity by reducing $E[T_{\text{off}}]$. As seen in the thesis, this corresponds to reducing the mean window sizes $E[W_i]$ at all backoff stages $i \geq 0$. However, when the saturation-based variant of the model is used, (9) indicates that a reduction of the mean window sizes affects the transmission probability τ and the conditional collision probability p , increases contention on the shared channel and may ultimately negate the intended effect, due to the impact of p on $\gamma_{\text{off}}(\cdot)$ through (12) and (13). The same phenomenon occurs also under non-saturated environments and affects the measured values of p , P_{succ} , P_{empty} and P_{coll} employed by the other variant of the model. Similarly, increasing the probability B_0 of sampling a null

¹ Number of active stations in the WLAN and (if the measurement-assisted variant of the model is used), loading conditions at other stations.

window at stage zero decreases $\gamma_{\text{off}}(\cdot)$ through (17), but also results in repeated successful transmissions in other competing stations, indirectly increasing $\gamma_{\text{off}}(\cdot)$ through the third term in (13) and (12).

For the reasons just described, in the following analysis it is assumed that any changes in the backoff window distributions leave the mean window sizes and B_0 invariant, thus also maintain the same value of the mean server rate \bar{r}_C . In order to examine the effect of higher order properties of the backoff window distributions, we use a Taylor series expansion of $u_C(\theta)$ around $\theta = 0$ and remember that $u_C(0) = 0$, to obtain

$$a_C(\theta) = u_C(\theta)/\theta = u'_C(0) + \frac{u''_C(0)\theta}{2} + O(\theta^2),$$

for small values of the parameter θ . Employing (18), differentiating twice, setting $\theta = 0$ and remembering that $(\log \gamma_i)'(0) = E[T_i]$ and $(\log \gamma_i)''(0) = \text{Var}[T_i]$ ($i = \text{on}, \text{off}$) leads to

$$u''_C(0) = \frac{\text{Var}[T_{\text{on}}](\hat{r} - \bar{r}_C)^2 + \text{Var}[T_{\text{off}}]\bar{r}_C^2}{E[T_{\text{on}}] + E[T_{\text{off}}]}. \quad (21)$$

Greater values of the Eff. Capacity function for negative arguments correspond to smaller values of $u''_C(0)$. Quantities in (21) relating to the On period do not depend on the backoff window distributions.

The mean value of the Off period can be obtained by differentiating (17) at $\omega = 0$. One gets

$$E[T_{\text{off}}] = t_{\text{ov}} + (1 - B_0)(t_{\text{slot}} + E[T_{\text{bc}}]), \quad (22)$$

where T_{bc} is the time spent in backoff mode, with moment generator as in (12) and mean equal to

$$E[T_{\text{bc}}] = \frac{p}{1-p}t_{\text{coll}} + E[T_s] \left(\frac{E[W_o]}{1-B_0} - 1 + \sum_{l=1}^{\infty} p^l E[W_l] \right). \quad (23)$$

$E[T_s] = \gamma'_s(0)$ is the mean of the time needed for the reduction by one of the IEEE 802.11 backoff counter. Equations (22) and (23) verify the earlier claim stating that, when the mean backoff window sizes and the probability B_0 remain invariant, the mean duration of the Off period and, by (20), \bar{r}_C also remain invariant. Thus, in view of (21), the only way of reducing the value of $u''_C(0)$ is through a smaller value of $\text{Var}[T_{\text{off}}]$. By differentiating twice (17) at $\omega = 0$ and collecting terms,

$$E[T_{\text{off}}^2] = B_0 t_{\text{over}}^2 + (1 - B_0)((t_{\text{over}} + t_{\text{slot}} + E[T_{\text{bc}}])^2 + \text{Var}[T_{\text{bc}}]),$$

so, using also (22),

$$\text{Var}[T_{\text{off}}] = E[T_{\text{off}}^2] - E[T_{\text{off}}]^2 = (1 - B_0)(B_0(t_{\text{slot}} + E[T_{\text{bc}}])^2 + \text{Var}[T_{\text{bc}}])$$

and reducing $\text{Var}[T_{\text{off}}]$ can only be achieved by reducing $\text{Var}[T_{\text{bc}}]$. It is shown in the thesis that a reduction of $\text{Var}[T_{\text{bc}}]$ may occur only through smaller variances for the backoff window sizes.

We now describe a specific way of adjusting the backoff window distributions, in order to reduce the variance in all backoff stages by a uniform percentage: The range

of the standard uniform distributions in (15) is narrowed from $[0, w_i - 1]$ to $[(w_i - 1)/2 - 1/2 - \delta_i, (w_i - 1)/2 + 1/2 + \delta_i]$. The same mean value is maintained, equal to $(w_i - 1)/2$. The parameter δ_i (a positive integer) is chosen so that the variance of the modified distribution is a fraction $\alpha < 1$ of the standard distribution's variance, so

$$\frac{4(\delta_i + 1)^2 - 1}{12} = \text{Var}[W_{i,\text{modified}}] = \alpha \text{Var}[W_{i,\text{standard}}] = \alpha \frac{w_i^2 - 1}{12},$$

which yields $\delta_i = \frac{1}{2} \sqrt{\alpha(w_i^2 - 1)} + 1 - 1$, $i \geq 1$. This result is rounded to the nearest integral value.

Extra arrangements must be made for the modification at the 0th backoff stage, because at this stage, besides the mean, one must also preserve B_0 , the probability of sampling a null backoff window.

Although that best results are achieved by selecting the smallest possible variances for all backoff stages, this is not an advisable strategy. The reason is that the IEEE 802.11 model has been constructed on the assumption that the competing mobile stations operate independently. Specifically, it is assumed that the collision of an observed station with one or more competing stations does not affect the probability with which this station will collide at the end of the next backoff stage. However, if deterministic backoff windows are used, when two or more stations collide at the end of some backoff stage $i \geq 0$, they will draw the same backoff window at the stage $i + 1$ and collision at the end of all stages from that point onwards will be certain. It follows that the backoff window distributions should retain a sufficient degree of randomness.

5 Validation of the IEEE 802.11 model for tail-related QoS

The Eff. Capacity model has been validated against ns-2 [10] simulation results, under various forms of traffic load and number of competing terminals. For details, see [7, 8]. Here we limit ourselves in two results, in the interest of further highlighting the concepts already discussed. In both of the results the system parameter values correspond to Frequency Hopping Spread Spectrum (FHSS) PHY layer [11]. Also, in all cases, the payload size (see (16)) was chosen constant and equal to $P = 8184$ bits.

Fig. 2 illustrates the accuracy of the Eff. Capacity, by comparing curves of the function (dashed, dotted and solid lines), computed with the use of the saturation-based model, against simulation results (marks).

In the simulation runs used for producing Fig. 2, the values of the Eff. Capacity function were indirectly measured, by feeding a “tagged” IEEE 802.11 station with traffic of known profile, sampling the probability with which the station's buffer exceeded a given threshold and exploiting the linkage (see (5)) between the Eff. Bandwidth, the Eff. Capacity and the probability tail just mentioned. All terminals, besides the tagged one, were operating under saturation conditions. The match between theory and simulations validates the model and indicates its suitability for estimating tail-probabilities or, equivalently, for taking AC decisions in IEEE 802.11 WLANs.

Fig. 3 depicts curves of the queue tail-probabilities versus the tail threshold (in semilog scale) for a network with 10 stations, of which 9 are saturated. The queue of the unsaturated station has been observed under two kinds of traffic load, CBR and Poisson,

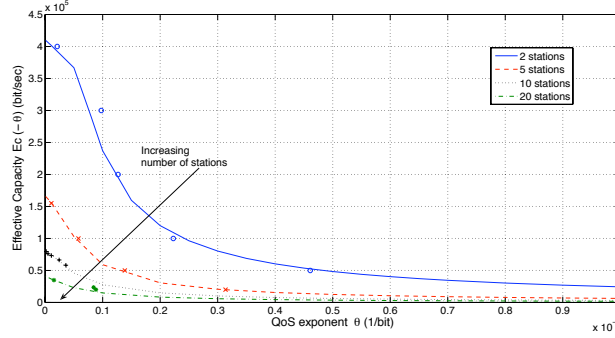


Fig. 2. Curves of $a_C(-\theta)$ vs. θ , for different values of the number of stations n .

both featuring the same mean rate of 79.84 kbps. The slope θ^* of the queue tail for the model-derived curve in each loading case was determined according to the theory (see (5)). As shown in the figure, the simulation-derived queue tails decay exponentially

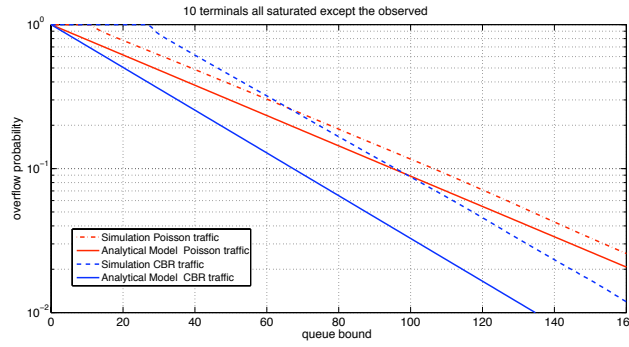


Fig. 3. Modeling and simulation results of queue tail-probabilities for 10 stations, one unsaturated, for two types of load.

and the decay rates agree well with the analytical results.

6 Conclusions

The thesis provided the formal justification for the use of the Eff. Bandwidth/Capacity theory in delay-related performance contexts. Specifically, it was established rigorously that the theory is capable of providing an asymptotically tight approximation to delay tail-probabilities. The thesis also formalized the association of the asymptotic exponential decay rate of the queue content probabilities with its counterpart for the delay

probabilities, through the server's Eff. Capacity function. The asymptotic approximation to the delay tail-probabilities was complemented by associated admission control schemes that are useful for providing delay-related QoS guarantees. The general results were applied to the important setting of IEEE 802.11 WLANs, by modeling each IEEE 802.11 station as an On/Off server and then using the Eff. Capacity function corresponding to this model. Computational and algorithmic details relating to the application of the general theory with the particular Eff. Capacity function of this On/Off model were also discussed. Comparison of the analytical results with simulation validated the effectiveness of the On/Off IEEE 802.11 model in providing tail-related QoS guarantees. Finally, the particular form of the Eff. Capacity function of an IEEE 802.11 station suggested an appropriate modification of the backoff window distributions for reduced variance, without affecting the mean backoff window sizes. This modifications results in a greater Eff. Capacity function, hence better performance.

References

1. A. Abdrabou and W. Zhuang. "Stochastic Delay Guarantees and Statistical Call Admission Control for IEEE 802.11 Single-Hop Ad Hoc Networks". *IEEE Trans. Wireless Commun.*, 7(10):3972–3981, October 2008.
2. G. Bianchi. "Performance Analysis of the IEEE 802.11 Distributed Coordination Function". *IEEE JSAC*, 18(3):535–547, 2000.
3. G. Bianchi and I. Tinnirello. "Remarks on IEEE 802.11 DCF Performance Analysis". *IEEE Commun. Lett.*, 9(8):765–767, 2005.
4. C. Chang. "Stability, Queue Length, and Delay of Deterministic and Stochastic Queueing Networks". *IEEE Trans. Autom. Control*, 39(5):913–931, 1994.
5. C. Chang and J. Thomas. "Effective Bandwidth in High-Speed Digital Networks". *IEEE JSAC*, 13(6):1091–1100, 1995.
6. C. Chang and T. Zajic. "Effective bandwidths of departure processes from queues with time varying capacities". In *Proc. IEEE INFOCOM*, pages 1001–1009, 1995.
7. E. Kafetzakis, K. Kontovasilis, and I. Stavrakakis. "Effective Capacity-based Stochastic Delay Guarantees for Systems with Time-Varying Servers, with an Application to IEEE 802.11 WLANs". *Elsevier Performance Evaluation*.
8. E. Kafetzakis, K. Kontovasilis, and I. Stavrakakis. "A Novel Effective Capacity-Based Framework for Providing Statistical QoS Guarantees in IEEE 802.11 WLANs". *submitted to Elsevier Computer Communications*, 2010.
9. K. Kontovasilis and N. Mitrou. "Effective bandwidths for a class of non Markovian fluid sources". In *Proc. ACM SIGCOMM, Computer Communication Review*, volume 27, pages 263–274, 1997.
10. The ns-2 network simulator. www.isi.edu/nsnam/ns, 1998.
11. "IEEE Standard for Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications, Nov. 1997. P802.11".
12. G. Vecianna, G. Kesidis, and J. Walrand. "Resource Management in Wide-Area ATM Networks Using Effective Bandwidths". *IEEE JSAC*, 13(6):1081–1090, 1995.
13. D. Wu and R. Negi. "Effective Capacity: A Wireless Link Model for Support of Quality of Service". *IEEE Trans. Wireless Commun.*, 2(4):630–643, 2003.
14. X. Zhang, J. Tang, H. Chen, S. Ci, and M. Guizani. "Cross-Layer-Based Modeling for Quality of Service Guarantees in Mobile Wireless Networks". *IEEE Commun. Mag.*, 44(1):100–106, 2006.

Secure Optical Communication Systems based on Chaotic Carriers

Dimitris Kanakidis *

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
dkanax@di.uoa.gr

Abstract. In this Thesis, the generation of a chaotic carrier by semiconductor lasers is theoretically studied. Three different encoding techniques are employed and the performance of an optical chaotic communications system for different receiver configurations is evaluated. It is proved that chaotic carriers allow the successful encoding and decoding of messages with up to 10 Gb/s bit rate. Focusing on the Chaotic Modulation encoding method it is shown that the normalized decoding method provides significantly different results in the estimation of the system's performance compared to the method of photocurrent subtraction, proposed in this study. Finally, the effects of transmission in optical fiber are studied for a chaotic communications system when CM and CSK encoding techniques are employed for various message bit-rates. The downgrading of the system's performance due to fiber transmission impairments is confronted by the employment of various dispersion management techniques.

Keywords: Chaotic Cryptography, Synchronization, Semiconductor Lasers, Optical fibers, Dispersion management, Secure Communications.

1 Introduction

The development of optical telecommunications systems has generated a particular interest in secure optical transmission [1], [2]. Owing to the high dimensionality and wide bandwidth of a chaotic signal, several optoelectronics and all-optical systems based on fast chaotic dynamics have been proposed as a possible alternative to classical encryption techniques relying on numerical algorithms. Chaotic cryptography performed in the physical layer rather than the application layer seems extremely attractive, since it exploits the deterministic nature of chaos, showing at the same time a strong dependence on minimal variations of the system's parameter values. Communication systems utilizing chaotic carriers are an extension of the conventional communication systems used today. While in a conventional communication system the message is modulated in the transmitter upon a periodic carrier, in chaotic systems an information message is transmitted using a chaotic

* Dissertation Advisor: Dimitris Syvridis, Professor

signal as a broadband carrier. The chaotic carrier in which data is encrypted can be either electronic or optical using corresponding electronic [3], or optic oscillators working in the chaotic regime. In the case of the semiconductor laser-based oscillators the operation in the chaotic regime can be achieved by applying optical feedback [4], optoelectronic feedback [2], or optical injection [5], and their chaotic behavior appears either in the amplitude or in the wavelength regime. Concerning the encoding of the information in the chaotic carrier, the schemes that have been proposed are the chaotic masking (CMS), chaotic modulation (CM), chaotic shift-keying (CSK), On-Off shift keying (OOSK) and chaos shift-keying using chaos in wavelength [1], [6]. The philosophy of the decoding process however is always the same, based on a very good synchronization between the transmitter and the receiver system.

In this Thesis, the generation of a chaotic carrier by semiconductor lasers is theoretically studied. Three different encoding techniques (Chaotic Modulation, Chaotic Shift Keying, and Chaotic Masking) are employed and the performance of an optical chaotic communications system for different receiver configurations is evaluated. It is proved that chaotic carriers allow the successful encoding and decoding of messages with up to 10 Gb/s bit rate.

Focusing on the Chaotic Modulation encoding method, the performance of a back-to-back chaotic communications system is also studied when two different decoding methods are followed. In both closed- and open-loop configurations, the normalized decoding method that is widely used in literature provides significantly different results in the estimation of the system's performance compared to the more realistic method of photocurrent subtraction, proposed in this study.

Finally, the effects of transmission in optical fiber are studied for a chaotic communications system when CM and CSK encoding techniques are employed for various message bit-rates. The downgrading appearing to the system's performance due to fiber transmission impairments is confronted by the employment of various dispersion management techniques.

2 Numerical modeling

For the numerical analysis, the optical carrier is a chaotic signal generated by a Fabry-Perot semiconductor laser diode subjected to an external optical feedback (master laser - ML). This chaotic waveform is optically injected into another semiconductor laser diode (slave laser - SL) that operates at the exactly same conditions as the ML (fig. 1a), and forces it to synchronize with the ML, producing at its output a chaotic carrier identical to that of the ML. The above setup, representing a closed-loop scheme, is theoretically described by the well-known Lang-Kobayashi equations [2]

$$\begin{aligned} \frac{dE_{i,r}(t)}{dt} = & \left(1 - j\alpha_{i,r}\right) \left(G_{i,r}(t) - \frac{1}{\tau_{p,i,r}}\right) \frac{E_{i,r}(t)}{2} \\ & + \gamma E_{i,r}(t - \tau) e^{j\omega_0 \tau} + \kappa_r E_{ext}(t) + \sqrt{2\beta N_{i,r}(t)} \xi_{i,r}(t) \end{aligned} \quad (1)$$

$$\frac{dN_{t,r}(t)}{dt} = \frac{I}{e} - \frac{1}{\tau_{n_{t,r}}} N_{t,r}(t) - G_{t,r}(t) |E_{t,r}(t)|^2 \quad (2)$$

$$G_{t,r}(t) = \frac{g(N_{t,r} - N_{0,t,r})}{(1 + s|E_{t,r}(t)|^2)} \quad (3)$$

where $E_{t,r}(t)$ is the complex slowly varying amplitude of the electric field at the oscillation frequency ω_0 and $N_{t,r}(t)$ the carrier number within the cavity. The subscript symbols t,r refer to the transmitter and the receiver system, respectively. The two lasers were considered to be similar to each other and therefore most of the internal parameters are the same. The spontaneous emission process is also taken into account through a complex Gaussian white noise term $\xi(t)$ of zero mean value and correlation

$$\langle \xi_t(t) \cdot \xi_r^*(t') \rangle = 2\delta_{t,r} \delta(t-t') \quad (4)$$

The light propagation along the fiber connecting the transmitter and the receiver has been modeled by integrating the well known generalized non-linear Schrödinger equation [7],

$$j \frac{\partial E}{\partial z} = -\frac{j}{2} \alpha E - \gamma |E|^2 E + \frac{1}{2} \beta_2 \frac{\partial^2 E}{\partial t^2} + \frac{j}{6} \beta_3 \frac{\partial^3 E}{\partial t^3} \quad (5)$$

where E is the total field envelope at the pump wave optical frequency, α is the fiber attenuation coefficient, $\gamma = \frac{2\pi n_2}{\lambda A_{\text{eff}}}$ is the nonlinear coefficient of the fiber responsible

for Self Phase Modulation (SPM), where n_2 is the Kerr coefficient, λ is the carrier wavelength, and A_{eff} is the fiber cross section; and constants β_2 , β_3 are the second-order chromatic dispersion and the third-order chromatic dispersion respectively and are related to the zero-dispersion wavelength and dispersion slope of the fiber. The transmission module additionally consists of in-line EDFAs and Gaussian optical band-pass filters. The former have been modeled under the assumption that the gain profile of the amplifier is practically constant in the wavelength region of the propagating wave, while its value is set to $G = \exp(\alpha L_c)$ in order to compensate for the fiber loss in an amplification period. Additionally, Amplified Spontaneous Emission (ASE) noise is taken into account introducing noise power per unit frequency $P_n = n_{\text{sp}} (G - 1) h \nu$, where n_{sp} is the spontaneous emission factor. The optical band-pass filters, used to suppress the ASE noise added by the EDFAs, are selected to be Gaussian of 1.1 nm bandwidth.

3 Performance Characterization of High Bit-Rate Optical Chaotic Communication Systems in a Back to Back Configuration

A comparative study of three data encoding techniques in optical chaotic communication systems is presented in this chapter. The chaotic carrier is generated by a semiconductor laser with optical feedback and the data are encoded on it by Chaotic Modulation (CM), Chaotic Masking (CMS) or Chaotic Shift Keying (CSK) (fig.1). In all cases the receiver, which is directly connected to the transmitter,

consists of a semiconductor laser similar to that of the transmitter subjected to the same optical feedback. The performance of this back to back configuration is numerically tested by calculating the Q-factor of the eye diagram of the received data for different bit rates from 1 to 20 Gb/s.

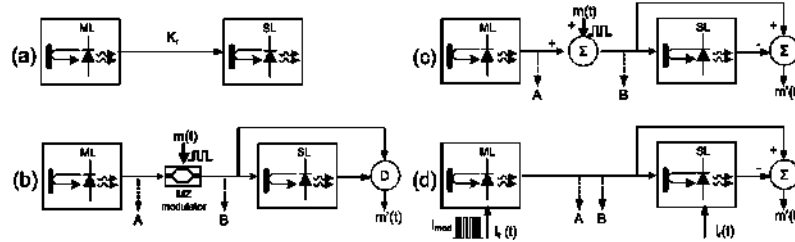


Fig 1. (a) The master laser (ML) output is optically injected to the slave laser (SL) and causes it to synchronize. (b) CM encryption method, (c) CMS encryption method, and (d) CSK encryption method.

Fig. 2 shows the Q-factor values calculated for different repetition rates, for the three encryption methods. These values are taken for the decoded message, as well as after filtering the decoded message with a fifth order Butterworth low-pass filter. Without any usage of filters, the CM encryption method has an obvious advantage over the other two methods. It exhibits very high Q-factor values (~ 18) at low values of bit-rate (1 Gb/s), while for higher bit-rates that end up to 20 Gb/s, the Q-factor value gradually decreases to 5. On the contrary, CMS method exhibits very low Q-factor values (~ 1.5), which remain practically constant as the bit-rate increases (fig. 2b). Finally, in the CSK method the Q-factor degrades in respect to the bit-rate increase, with values confined below 4.

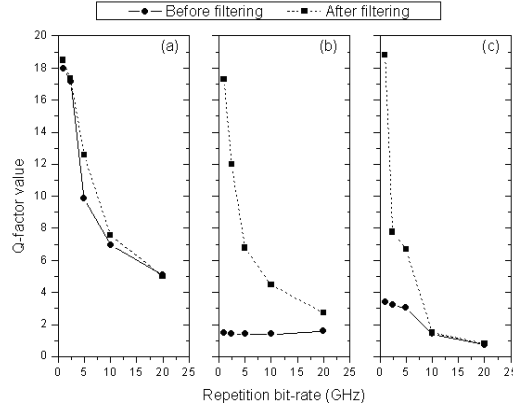


Fig 2. Estimation of Q-factor values for different repetition bit-rates of the decoded message, employing (a) the CM method, (b) the CMS method, and (c) the CSK method. The synchronization error of the master – slave system is set to 0.2%.

The higher Q-factor values extracted from the CM method in comparison to the other two methods can be explained by the nature of the encoding process. For the optical feedback system under investigation, complete synchronization entails not only synchronization of the laser field amplitude waveform, but also the synchronization of the laser field phase. In the CM method the message is being applied by modulating the transmitter's chaotic carrier itself, according to the expression: $(1 + m(t)) \cdot E_{ML} \cdot e^{i\phi_{ML}}$, resembling the typical coherent AM modulation scheme. Therefore, the phase of the chaotic carrier with the message encoded on it, which is injected to the receiver, is the same with that of the chaotic carrier without any information. Thus, the presence of signal on the chaotic carrier does not cause a significant perturbation in the synchronization process of the system. On the contrary, in the CMS method the message is a totally independent electric field, which is added to the chaotic carrier, according to the expression: $E_{ML} \cdot e^{i\phi_{ML}} + m_E(t) \cdot e^{i\phi_{message}}$. Therefore, the phase of the total electric field injected now to the receiver consists of two independent components. The phase of the message acts in this case as a perturbation in the phase matching condition of a well-synchronized system and thus, the phase difference between the transmitter and the receiver diverges from zero. In the CSK scheme with one receiver, although the phase exhibits the same behavior as in the CM method, the Q-factor values are degraded due to the synchronization error of the system, induced by the different injection currents between the ML and the SL. An additional drawback of the CSK method is that the laser cannot be modulated properly for frequencies above the relaxation oscillation frequency.

In order to improve the efficiency of the studied encoding schemes, a fifth-order Butterworth low-pass filter is employed. The cut off frequency of the filter is optimized for the different repetition bit-rate values of the encoding message. As it can be seen from fig. 2, for the 1 Gb/s bit-rate message all methods exhibit a similar behavior, providing a very high Q-factor value (17 to 19), since the high frequency oscillations have been entirely removed.

By increasing the synchronization error of the system to higher values, the Q-factor is found to be minimized in all three methods, even at low repetition bit-rates. The improvement caused by the filter usage is limited. For example, 5% synchronization error gives a Q-factor value of almost 6, in the CM method after filtering, while 30% error further limits the Q-factor value to a low value of 2.5. The above behavior proves that when the synchronization error is high, the message cannot be recovered at all. Thus, it is crucial to maintain the system well synchronized in order to achieve a fully recovered message with a high Q-factor value [8].

4 Influence of the Decoding process on the Performance of Chaos Encrypted Optical Communication Systems

Two different decoding methods of an all-optical chaotic communication system are investigated when Chaotic Modulation (CM) encoding format is employed. The transmitter consists of an external cavity semiconductor laser generating thus a chaotic carrier, modulated using an external modulator. The receiver is either a solitary semiconductor laser diode identical to that of the transmitter or a laser diode

coupled to an external cavity, forming an open or a closed-loop configuration respectively (fig. 3). The performance of the system is then evaluated by means of calculating the Q-factor extracted by the eye diagram of the recovered data when two different approaches of the decoding process for the receiver are adopted. The first decoding method relies on the normalized to receiver's amplitude output, difference of the two lasers' optical amplitude outputs, while the second one corresponds to a more realistic case by subtracting the electrical current outputs of two p-i-n photodiodes coupled to the transmitter and receiver laser correspondingly. By comparing the numerical results extracted by the two decoding methods and for various cases of interest, such as employing open- or closed-loop configuration for several message bit-rates and different laser's driving current, it was shown that, under certain circumstances, the two decoding methods result to significantly different results.

The encryption method considered in the present work is chaotic modulation (CM) (fig. 3). In this method the message encoding is carried out by modulating the chaotic carrier of the ML using an external Mach-Zehnder modulator, according to:

$$|E_t^{\text{mod}}(t)| = |E_t(t)|(1 + m(t)) \quad (6)$$

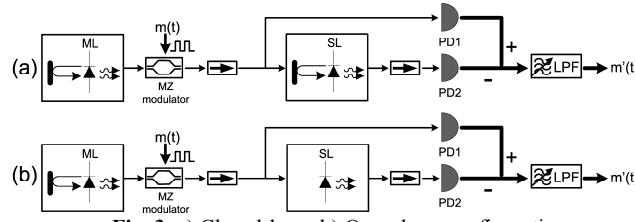


Fig. 3. a) Closed-loop, b) Open-loop configuration

where $m(t)$ is a 1-5Gb/s pseudorandom bit sequence forming a codeword with length $2^{11}-1$. The modulation depth engaged for the simulations was chosen to be no more than 5% so that the message not only was totally indistinguishable from the chaotic carrier of the ML but caused limited distortion in the synchronization process of the system as well. The proposed decoding process, in many papers exhibiting numerical simulations, is of the form:

$$m'(t) = \sqrt{\frac{|E_t^{\text{mod}}(t)|^2}{|E_r(t)|^2}} - 1 \quad (7)$$

which actually, under perfect synchronization conditions, mathematically provides a replica of $m(t)$. To approximate a possible decoding procedure in real, laboratory conditions two photodiodes were adopted to provide a subtraction between the two chaotic outputs of the ML and SL, rather than the division used in the first method. The equation describing this function is

$$I_{\text{dec}} = I_t - I_r \quad (8)$$

where $I_{t,r}$ denotes the photocurrent created by the optical injection of the output waves of the ML and SL into the photodiodes PD1 and PD2 respectively.

In fig. 4 the Q-factor value in respect to the normalized injection strength for the two different decoding methods is presented. The results correspond to the case of a

closed-loop scheme when 1 Gb/s message is embedded on the chaotic carrier. Two different regions can be easily discriminated in this figure.

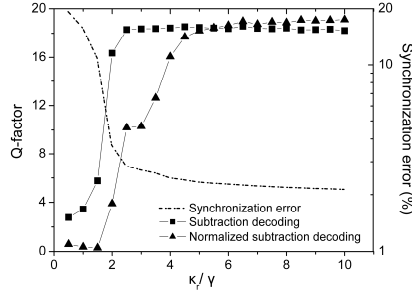


Fig. 4. Q-factor and synchronization error as a function of the normalized injection strength for 1Gb/s message bit-rate and closed-loop configuration.

In the first one, for normalized injection strength values κ_r/γ greater than 5, the two decoding methods provide similar performance results. On the other hand, for κ_r/γ values less than 5, a significant difference in the measured Q-factors appears. In detail the behavior of each decoding method can be discussed as follows:

A. Normalized difference of Electric fields

By using the normalized decoding method according to eq. 7 (which for simplicity will be further called first method), a significantly lower estimation of the system's performance is observed. The normalized difference between the two methods for low values of the injection strength exceeds 90% and gradually diminishes to zero as the injection strength increases. The synchronization error of the system (defined in [9]), also plotted in the same figure, provides a possible explanation of the observed behavior. It is evident that the Q-factor curve concerning the first method directly follows (in the opposite obviously direction) the synchronization error curve. This is somewhat mathematically expected since $m'(t)$ becomes identical to $m(t)$ under perfect synchronization conditions as indicated by equations (6), (7).

B. Subtraction of Photocurrents

For the subtraction based decoding method according to eq. 8, (which for simplicity will be further called second method) there is a major difference. Mathematically, under perfect synchronization conditions and in the absence of every added by the photodiodes noise terms, the decoded photocurrent follows the relation $I_{dec} \sim m(t)|E(t)|^2$. This means that in the second decoding method the quality of the decoded signal not only depends on a good synchronization between ML and SL but on the chaotic carrier of the transmitter itself. Therefore the spectrally wider the transmitted chaotic wave is, the worse performance is expected for the second decoding method. Moreover, it is well known that chaos spectrally grows up around the relaxation oscillation frequency of the laser. Consequently, a deterioration of the system performance is expected as more frequency components of the chaotic carrier are included in the decoded message after the filtering procedure. This could happen for example if higher message bit-rate was employed. For the 1 Gb/s message bit-rate though, a good decoding quality is expected since the electrical filter utilized has a cut-off frequency of 600MHz, eliminating the higher spectral components of the

chaotic carrier. By observing fig. 4, it can also be pointed out that for high values of synchronization error the synchronization quality is the dominant factor deteriorating the performance calculated using the second decoding method, while at moderate injection parameter values ($\kappa_i/\gamma=2-4$), the calculated Q-factor values are at a constant high level, despite the corresponding synchronization improvement. Therefore, according to the second decoding method, a floor on the best achievable performance is expected.

To further verify the dependence of system's performance mainly from the spectral characteristics of the transmitted chaotic carrier using the second decoding method, and exclusively from the degree of synchronization using the first method, the performance of the same setup is evaluated for 2.4 & 5 Gb/s message bit-rates (fig. 5). It is evident that when the first decoding method is used, a direct relation between the Q-factors and the synchronization degree is observed. In fact the Q-factor curve resembles a mirror image of the corresponding synchronization error curve with slight differences, basically due to the fact that the synchronization error was calculated in the whole spectral region while the Q-factor was calculated after filtering the remaining high frequency non-synchronized spectral components. Concerning the first decoding method, a system's performance degradation is also observed relative to the low message bit-rate. This result, depicted in fig. 5, can be attributed to the fact that when higher message bit-rates are employed, output filters with increased bandwidth are used. This allows the higher non-synchronized spectral components to be part of the filtered decoded message, deteriorating the performance. On the other hand this is not solely the case for the second decoding method. Significant performance degradation, especially when compared to the first decoding method, is observed in fig. 5 which cannot be attributed exclusively to the increased filter's bandwidth. As already mentioned, there is an enhancement of the chaotic carrier's spectral components around the relaxation oscillation frequency ($\sim 4.5\text{GHz}$) of the laser. Since the frequency components of the message bit-streams come closer (2.4 Gb/s) or even exceed (5 Gb/s) the relaxation oscillation frequency, the decoding quality, which depends on the chaotic carrier, is degraded even in the presence of a filter. Compared to the first decoding method, synchronization degree only confines its effect on the high synchronization error region since in all the other cases, although the synchronization error of the system decreases, the Q-factor stabilizes to a certain value.

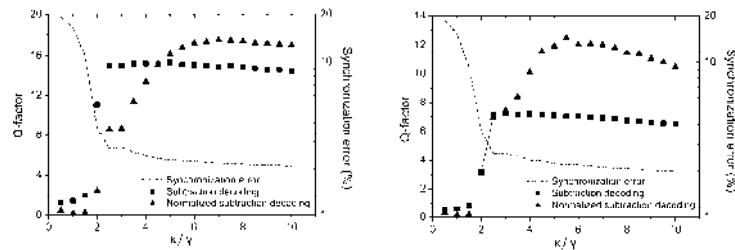


Fig. 5. Q-factor and synchronization error as a function of the normalized injection strength for 2.4Gb/s (left) and 5Gb/s (right) message bit-rate and closed-loop configuration.

To generalize the conclusions derived so far an analogous study was performed in the case of an open-loop scheme. From the study performed it was evident that the two different decoding methods exhibit similar performance behavior compared to the corresponding closed-loop cases confirming the generality of the conclusions derived for the closed-loop scheme [10].

5 Numerical Investigation of Fiber Transmission of a Chaotic Encrypted Message using Dispersion Compensation Schemes

A detailed numerical investigation of the transmission properties of all-optical chaotic communication systems is finally presented for two data-encoding techniques and for various dispersion compensation maps. A semiconductor laser subjected to optical feedback generates the chaotic carrier and the data is encoded on it by CM or CSK methods. The complete transmission module consists of different types of fiber, in-line amplifiers and Gaussian optical filters. Different dispersion maps based on either Non-zero Dispersion Shifted Fibers (NZ-DSFs) or combinations of Single Mode Fibers (SMF) along with Dispersion Compensating Fibers (DCF) were considered (fig. 6). The system's performance is numerically tested by calculating the Q -factor of the eye diagram of the received data for 1 and 2.4Gb/s. The influence of the optical power launched into fiber and the transmission distance to the quality of the decoded message has been investigated.

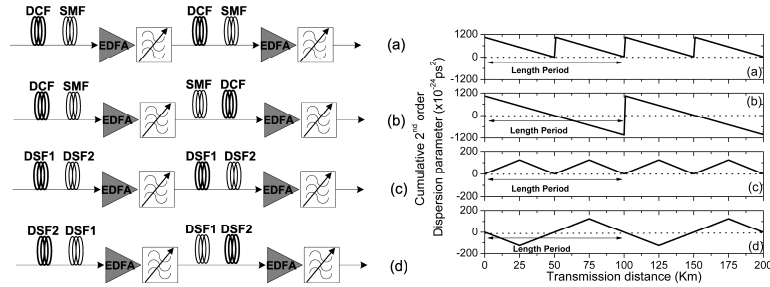


Fig. 6. Cumulative second-order chromatic dispersion parameter in respect to the transmission distance (right) for a) Pre-compensation map, b) Symmetrical map, c) One-step map, d) two-step map.

By utilizing a NZ-DSF of a small value of second order dispersion parameter ($1\text{ps}^2 \text{Km}^{-1}$) in the link between master and slave laser, it was found that for distances that exceed 10 Km no message extraction is possible. Consequently, dispersion management techniques should be appointed, so that adequate message quality is derived for longer transmission spans.

In order to identify the upper limits of the transmission maps employed, in terms of the maximum distance achievable for acceptable message quality (Q factor values higher than 6), the optimum value of the launched optical power from the transmitter to the transmission module should be determined. It is well known that Self Phase

Modulation (SPM) induced impairments are enhanced as the optical power propagating through the fiber increases. On the other hand as the optical power of the transmitted signal reduces, the gain, and in consequence the ASE noise added by the in line EDFAs, increases. Therefore, the optimum value of the optical power injected into the fiber from the transmitter results from a trade-off between the non-linear effects of the fiber, and the added by the in-line amplifiers ASE noise.

The performance of the transmission maps proposed (utilizing DCF-1 in the corresponding dispersion compensation maps), by means of the Q-factor values in respect to the launched optical power, is presented in fig. 7 for a total transmission period length of 100 Km and for both modulation formats. As expected, the curves for each map have a peak corresponding to the optimum value of the launched optical power. It is also evident that the maps utilizing the NZ-DSFs (fig. 6c and 6d) appear to be more insensitive to injected optical power, sustaining high Q-factor values for a wide range of power values. On the other hand, symmetrical and pre-compensation map show strong dependence from the optical power, due to the enhanced SPM phenomena related to the DCF fibers employed.

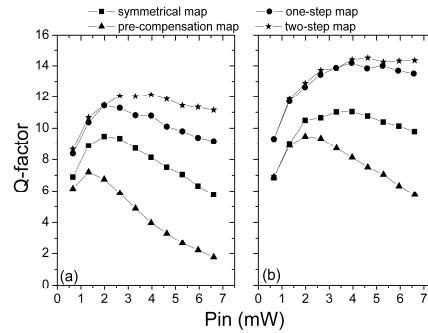


Fig. 7. Q-factor values, for 1 Gb/s bit-rate and 5% modulation depth, as a function of the launched optical power for all transmission maps proposed and for a) CM method, b) CSK method.

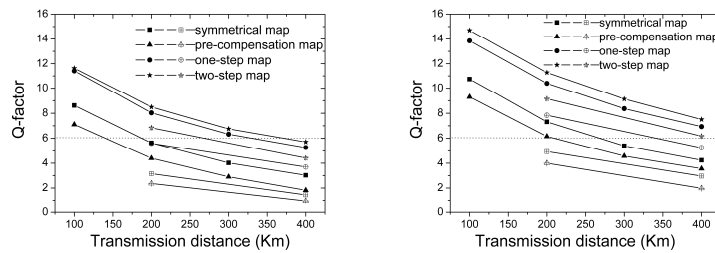


Fig. 8. Q-factor values, for CM method (left) and CSK method (right), as a function of the transmission distance for all maps proposed and for a) 100 Km map period length (solid symbols), b) 200 Km map period length (blank with crosses symbols).

For the optimum power values extracted from fig. 7, system's performance was investigated in respect to the transmission distance. The latter extends from 100 to

400 Km by simply utilizing multiple periods of the initial, 100 Km length, maps proposed. In fig. 8 the calculated Q-factor is shown versus the transmission distance, for all the maps used and for both message encoding methods. The strong degradation of the system performance, in respect to the transmission distance, is caused by the accumulated nonlinearity and ASE noise effects. The contribution of these two phenomena distorts the phase and amplitude characteristics of the chaotic carrier, degrading severely the synchronization quality in the receiver and thus the system performance. For both modulation formats, the dispersion maps based on NZ-DSFs (fig. 6c and 6d) appear to be more efficient than the other two based on DCF as dispersion compensating element. An evident performance improvement is also observed between the symmetrically constructed maps (two-step map and symmetrical map) and the corresponding (by means of the different types of fiber employed) non-symmetrical maps.

In order to investigate the performance of the system at higher bit-rates, a 2.4 Gb/s message bit-stream with a modulation depth of 5% was applied to the chaotic carrier. The curves extracted by simulations, were similar with the ones presented for the 1 Gb/s message bit-rate case (fig. 8), indicating that the inherent properties of the maps used (the type of fibers employed, the way they are placed along the whole map and the spacing between the EDFAs), mainly determine the optimum power level for each map.

In order to improve the performance of the system for the 2.4 Gb/s case, the modulation depth of the message applied to the chaotic carrier was increased. The value of 7% chosen not only improves the decoded message quality but keeps the high level security of the system as well. It was found that in CM method the upper distance limits are about 190 and 230 Km for both the one and two-step map respectively, while for the dispersion compensation maps the corresponding values, for the symmetrical and pre-compensation map, are 150 and 140 Km [11]. On the other hand, CSK method appear to have quite improved performance, since the upper distance limits determined, were 280 and 300 Km for the one and two-step map respectively, and less than 200 Km for both the symmetrical and pre-compensation map.

6 Conclusions

A performance comparison by means of Q-factor calculations between the three encoding methods was demonstrated in the above analysis. It has been shown that the CM method has an obvious advantage over the other two methods, due to the fact that the message carries the phase of the chaotic carrier. However, after filtering the chaotic high frequency oscillations, all methods result satisfactory Q-factor values for low repetition bit-rates up to 2.4 Gb/s. By increasing the bit-rate to 10 Gb/s, only CM method could be characterized as a sufficient encoding scheme, while at 20 Gb/s the best Q-factor value extracted is almost 5 and is referred to the CM method, also.

Then a performance analysis by means of Q-factor calculation was also conducted for CM codification scheme and two possible decoding methods in a back-to-back chaotic communication setup. For both closed and open loop configurations, it was

shown that the first decoding method provided performance results strictly dependent on the degree of the synchronization between ML and SL. On the other hand, the performance results extracted with the second decoding method proved to be strongly determined by the complexity of the chaotic carrier itself.

Finally it was shown that the CSK method has an obvious advantage over the CM method in all the transmission maps proposed. In order to explore the upper limits of the system, by means of the maximum distance achievable for acceptable message quality, the optimum value of the launched optical power has been determined. The maps utilizing NZ-DSFs appear to be more robust to variations in the injected optical power, compared to the dispersion compensation maps for both the message bit-rates employed. When 1 Gb/s bit-rate is used, dispersion management maps, though less practical, proved to provide better results. When a 2.4 Gb/s message is applied to the chaotic carrier the corresponding in each case, Q-factor values are decreased because of the residual spectral components of the chaotic oscillations in the higher frequency regime. However, by increasing the modulation depth of the message, improved results can be obtained.

References

1. A. Sanchez-Diaz, C. R. Mirasso, P. Colet, and P. Garcia-Fernandez, "Encoded Gbit/s digital communications with synchronized chaotic semiconductor lasers," *IEEE J. Quantum Electron.*, vol. 35, pp. 292–297, Mar. 1999.
2. S. Tang and J. M. Liu, "Message encoding/decoding at 2.5 Gb/s through synchronization of chaotic pulsing semiconductor lasers," *Opt. Lett.*, vol. 26, pp. 1843–1845, Dec. 2001.
3. T. L. Carroll and L. M. Pecora, "Synchronizing chaotic circuits," *IEEE Trans. Circuits Syst. I*, vol. 38, pp. 453–456, Apr. 1991.
4. V. Annovazzi-Lodi, S. Donati, and A. Scire, "Synchronization of chaotic lasers by optical feedback for cryptographic applications," *IEEE J. Quantum Electron.*, vol. 33, pp. 1449–1454, Sept. 1997.
5. J. M. Liu, H. F. Chen, and S. Tang, "Optical communication systems based on chaos in semiconductor lasers," *IEEE Trans. Circuits Syst. I*, vol. 48, pp. 1475–1483, Dec. 2001.
6. C. R. Mirasso, J. Mulet, and C. Masoller, "Chaos shift-keying encryption in chaotic external-cavity semiconductor lasers using a single-receiver scheme," *IEEE Photon. Technol. Lett.*, vol. 14, pp. 456–458, Apr. 2002.
7. G. P. Agrawal, *Nonlinear Fiber Optics*, San Diego, CA: Academic, 2001.
8. D. Kanakidis, A. Argyris, and D. Syvridis, "Performance Characterization of High-Bit-Rate Optical Chaotic Communication Systems in a Back-to-Back Configuration", *IEEE J. Lightwave Technol.*, vol. 21, pp. 750–758, Mar. 2003.
9. J. M. Liu, H. F. Chen, and S. Tang, "Synchronized chaotic optical communications at high bit-rates," *IEEE J. Quantum Electron.*, vol. 38, pp. 1184–1196, Sept. 2002.
10. D. Kanakidis, A. Argyris, A. Bogris and D. Syvridis, "Influence of the Decoding Process on the Performance of Chaos Encrypted Optical Communication Systems", *IEEE J. Lightwave Technol.*, vol. 24, pp. 335–341, Jan. 2006.
11. D. Kanakidis, A. Bogris, A. Argyris, and D. Syvridis, "Numerical Investigation of Fiber Transmission of a Chaotic Encrypted Message Using Dispersion Compensation Schemes", *IEEE J. Lightwave Technol.*, vol. 22, n. 10, pp. 2256–2263, Oct. 2004.

Optical Microring Devices for Optical Networks

Applications: Investigation, Design and Characterization

Alexandros Kapsalis*

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
alex@di.uoa.gr

Abstract. In this dissertation, the optical characteristics of active microring resonators are investigated. The result of this thesis is the complete investigation of the SML properties and also the optimization of the key design parameters in order to be used in future microring-based lightwave system applications. In the first part the following are presented: investigation of the influence of various structural parameters, correlation between structural and functional parameters, analytical description of a numerical simulation model based on multimode rate equations. In the second part experimental measurements are presented. These measurements are: characterization of all-active and active-passive integrated devices, examination of different ring radii, investigation of operating regimes and a strong correlation between the coupling efficiency and threshold current density was established. Mode hopping phenomena and phase and intensity noise were investigated. A complete matching of the experimental findings with the theoretical predictions was established. Tuning capabilities of SMLs are also presented for different configurations where an extra broad tuning range of 40nm is recorded. Finally some applications: direct modulation at 7Gb/s and transmission through 100km, wavelength conversion in similar data rates by exploiting cross gain modulation (XGM) and a 10 dB optical amplification at 10Gb/s modulated.

Keywords: Microring resonator waveguiding analysis, semiconductor microring lasers, tunable laser sources, direct modulation and propagation.

1 Dissertation Summary

This dissertation is structured in two parts. The first, deals with microring lasers from a theoretical point of view, whereas in the second part, experimental measurements are demonstrated from actual devices measured on the optical communications laboratory. In more detail, the summary per chapter of this dissertation is as follows.

Beginning with the first chapter of the dissertation, an overview of optical networks is presented. A synopsis of the technological evolution of optical networks is shown, alongside the requirements of second and third generation networks for all optical signal processing. These include standards for devices like laser sources, photoreceivers, optical amplifiers, wavelength converters, filters and logic gates. According to these requirements the desired functionalities of laser sources are

*Dissertation Advisor: Dimitrios Syvridis, Professor

described in order to be included in nodes of future generation networks. At the end of the chapter the contribution of the thesis in the realization of fundamental building block of WDM networks is developed.

Microring resonators are ideal devices for developing almost every fundamental block in WDM networks. They are compact and so can be integrated in very large scale. In parallel, they are wavelength selective which in combination with the ability to amplify field due to feedback makes them suitable for a broad range of applications in optical signal processing. Furthermore, the lack of costly mirrors or Bragg gratings in order to achieve optical feedback, simplifies the fabrication procedure and cost and allows them to be monolithically integrated with other structures. Here, a full description of the optical characteristics of active microrings is presented. Microring resonators are nowadays well established as the alternative way for optical resonators in general, alongside Fabry-Perot cavities and Bragg gratings. Specifically, microring lasers are offered as a means to construct sources for future optoelectronic integrated circuits (OEICs) because they combine simple fabrication and small size. Microring laser devices coupled with a straight bus waveguide that delivers light in or out of the device are explored. The study is intended to provide insight on the correlation of structural design parameters like the ring radius, the ring-bus power coupling efficiency and the bus facet residual reflectivity with operational characteristics like threshold current, phase and intensity noise, wavelength tunability and modulation capabilities. The result of this investigation is the synthesis of a general picture about the properties of microring laser diodes and the determination of the optimal design guidelines to be used in OEIC of the future.

In this context, a thorough investigation of microring resonator properties was carried out with special focus on microring lasers. This begins from a design level, passes through the development of numerical modeling and concludes with experimental results not only with characterization but also in practical applications.

The multiple InGaAsP/InP quantum well structures were fabricated at the Fraunhofer Institut in Berlin, Germany using wafer bonding techniques. The purpose of the measurements was to compare theoretical and experimental results and draw useful conclusions about the optimization of the abovementioned devices. Also, actual systems were implemented with performance comparable to existing state of the art applications, which by the time of their publication were considered to be a first. These implementations, which are a main part of this dissertation contribution, highlight a series of novel applications of microrings in WDM building blocks and emphasize on the importance of these structures on the evolution and development of optical networks in general.

The second chapter includes the theoretical background associated with the coupling of light between a bent and a straight waveguide. The influence of basic design aspects to the achievable coupling efficiency is examined which include: the ring radius, waveguide width and lateral displacement. The analysis is focused on the achievement of transversal monomode operation in both active and passive waveguides, since the existence of higher order modes leads to degradation of the performance of these devices by adding an extra loss route. Three methods are proposed and analyzed to circumvent this problem. These include: selective excitation of the fundamental mode, selective attenuation of higher order modes and monomode coupling by imposing a positive lateral displacement between the two waveguides.

Furthermore, phenomena related to polarization rotation in curved waveguides are briefly investigated in order to design polarization conversion free devices.

An extensive report on microring resonators and lasers based on microring devices is presented in chapter 3. The equations that describe the light circulation and lead to resonance are analyzed. Also, an overview of up to the date of this dissertation existing microring laser devices is given and at the end the correlation between the design parameters and operational characteristics is theoretically studied. These operational characteristics include phase and intensity noise.

The fourth chapter is dedicated to the analytical development of a complete numerical model for simulating microring lasers coupled with a straight waveguide, based on multimode rate equations. A detailed description of the related equations is demonstrated. The enclosure of linear and non-linear gain terms is analyzed as well as the symmetric and asymmetric ones that lead to the various mode hopping phenomena. Furthermore, the inclusion of external feedback due to residual reflectivity is explained which is practically incorporated as an additional injection term in the rate equations.

The previously described numerical model was used to study the microring laser behavior and the accompanying results are shown in chapter 5. The spectral characteristics of the laser are described and the various operation regimes are identified. These include unidirectional and bidirectional as well as single mode and multimode operation. Moreover, the noise of microring lasers was examined with emphasis on the study of the spectral distribution of relative intensity noise. Both all-active and active-passive devices were simulated for different ring radii and coupling efficiency values, whereas, the tuning properties of microring lasers were theoretically investigated.

In chapter 6, the experimental investigation begins. The fabrication and characterization of various devices is described. At first, results from passive ring structures are shown, which were fabricated for confirmation with theoretical trends and also for parameter extraction like facet reflectivity and waveguide loss. All active devices were also characterized, where both the ring and straight waveguide exhibit gain by injecting current. Also, active-passive integrated devices were characterized in which the bus waveguide is made of passive material in order to minimize the influence of back reflection from the cleaved facets. Characterizations were performed for a variety of ring radii and bus waveguide widths which lead to different coupling efficiency values. The investigation of the different regimes of operation was accomplished through optical and RF spectra, optical time traces and light-current (LI) curves.

Noise phenomena in microring lasers are examined in chapter 7. Initially, mode hopping phenomena are investigated when lasers are operated at the high injection multimode regime. Additionally, the behavior related to intensity fluctuations is studied through time traces, optical and RF spectra as well as relative intensity noise (RIN) measurements. Finally, the phase noise is evaluated by measuring the laser's linewidth where a matching of the experimental findings with theoretically anticipated values is confirmed.

Chapter 8 covers the tuning properties of microring lasers. In 'all-active' microring devices, tuning is achieved through the straight active bus waveguide that controls the phase of the back reflected field by adjusting the injection current. In single ring

active-passive devices the same apply with the exception that the bus waveguide is passive, therefore, non-controllable. So, the tuning properties are restricted. In double ring active-passive devices, an extensive tuning range is accomplished through the Vernier mechanism between the two active cavity spectra.

In chapter 9, applications of active microring devices are demonstrated. More specifically, microrings are investigated as directly modulated laser sources with modulation rate capabilities up to 7Gb/s. Also, the propagation of a signal, using a microring laser optically modulated at 5Gb/s, for 100km of single mode fiber (SMF) is shown. Moreover, the usage of active microring lasers as wavelength converters up to 5Gb/s was presented, by exploiting the cross gain modulation (XGM) and as an optical amplifier, capable of amplifying 10Gb/s signals by 10dB.

In the tenth and final chapter, a result and related conclusions summary is presented. Additionally, future works related to the contribution of this dissertation are discussed.

2 Results and Discussion

All results related to waveguide design for achieving monomode transversal behavior are summarized in [1-5]. All-active devices and numerical modeling related applications are described in [7-10]. Previous studies of vertically coupled microring lasers using a wafer-bonded transfer substrate have defined the bus waveguide prior to inverting the wafer. However technological challenges derived from the complexity of processing active components has lead to an inversion in the design. The laser structures are defined on the first face, prior to wafer bonding. The complete epitaxial structure is grown, the laser waveguides are defined and the p-type metal is deposited prior to wafer bonding to the transfer wafer. The substrate is removed by etching prior to the formation of bus waveguides and the opening of via holes to the p-metal at the transfer wafer. This results in a buried active microring waveguide structure, with the additional advantage that the laser is well protected during wafer dicing process. A standard laser epitaxial growth process now becomes feasible prior to wafer bonding. The buried microring structure is clad by low index benzo-cyclo-butene (BCB) material applied during the wafer-bonding process, where the high index difference ensures a good modal confinement within the waveguide.

The passive bus waveguides are formed on top of the active microring structure. Both p and n contacts are implemented on the top surface of the fabricated device. The width of both the passive bus waveguide and active microring are 1.8 μm . The active region of the microring laser comprises of six GaInAs quantum wells, with a band gap of $Q_{\text{ring}}=1.55 \mu\text{m}$, the InGaAsP bus waveguide has a band gap of $Q_{\text{bus}}=1.44 \mu\text{m}$, and a thickness of 0.35 μm . In order to minimize undesired reflections resulting from cleaving a 7° tilt is implemented on both ends of the bus waveguide. Detailed descriptions of the wafer bonding fabrication process can be found on [11-13].

2.1 Back to back measurements

Ring devices with ring radii ranging from 40 μm to 80 μm are analyzed. A copper submount is used to place the chips and temperature at 20°C is maintained throughout all measurements. A 50 Ω terminated Ground Signal Ground (GSG) probe head is used for biasing and modulation and the optical output is coupled out using tapered fibers. Single mode emission is of crucial importance for the quality of high data rate operation. In figure 1, an L-I curve is shown along with an emission spectrum for bias current 41.28 mA and modulation depth 1 V_{pp} where single mode operation with side mode suppression ratio (SMSR) above 30dB is observed.

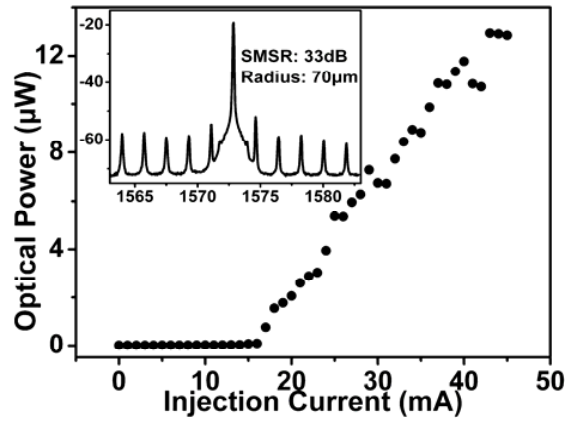


Fig. 1. L-I curve for a microring laser with radius 70 μm . The inset shows a single mode spectrum with SMSR approx. 33 dB for modulation bias at 41.28mA.

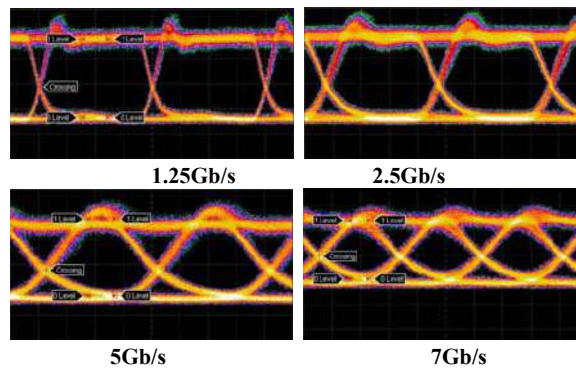


Fig. 2. Eye diagrams of ring laser modulated at 1.25, 2.5, 5 and 7 Gb/s with ER 7.6, 8.2, 8.2, and 5.8dB respectively in back to back configuration.

Eye patterns as well as Bit Error Rate (BER) measurements are carried out using various pseudo random bit sequence (PRBS) data streams for Non Return to Zero

(NRZ) pulses. Well opened eye diagrams are shown in figure 2, for 1.25, 2.5, 5 and 7 Gb/s bit rates for a 70 μm radius microring laser with extinction ratios (ER) 7.6, 8.2, 8.2 and 5.8 dB respectively. Modulation at 10Gb/s was also performed but with a poor extinction ratio of 3.7dB. Back to back BER measurements are in accordance with the eye diagrams showing error free ($\text{BER} < 10^{-12}$) operation up to 7 Gb/s where as for higher data rates BER increases significantly (e.g. 10^{-6} for 8 Gb/s).

2.2 Transmission experiment

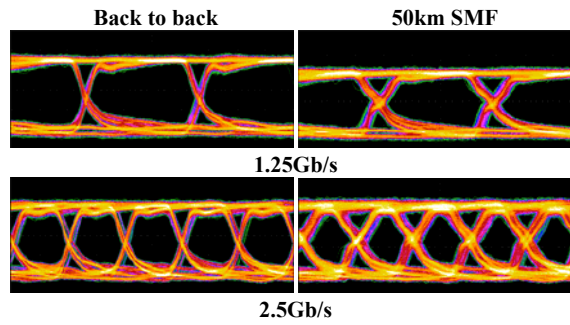


Fig. 3. Eye diagrams for a ring laser modulated at 1.25 Gb/s and 2.5 Gb/s in back to back configuration (left) and after 50 km propagation (right) in a SMF.

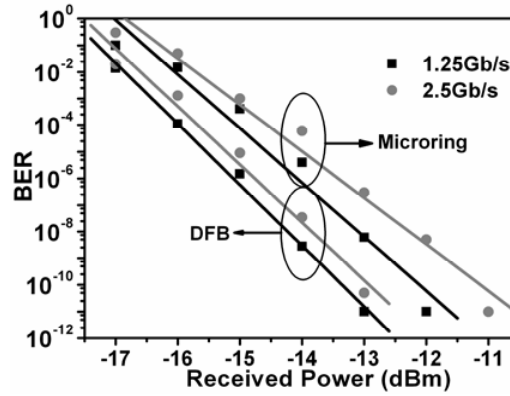


Fig. 4. BER measurements at 1.25 Gb/s and 2.5 Gb/s (50 km), for the microring and a DFB laser with the same ER.

A transmission experiment is also carried out using the microring lasers as data modulated sources for the first time to date. The signal from a 70 μm radius microring is amplified by an EDFA, filtered and launched into a 50 km single mode fiber. Error free transmission at 1.25 and 2.5 Gb/s is reported without dispersion compensation.

ER remains stable at 7.5dB for 1.25 Gb/s and reduces from 8 dB to 7.2 dB at 2.5 Gb/s and the corresponding eye diagrams are shown in figure 3.

A comparison with an edge emitting laser is made showing that microrings have comparable performance as shown in figure 4. In order to achieve a good extinction ratio a modulation depth in the order of 1V is needed which corresponds to a 5mA injection current span. Over this current span, a microring under CW operation hops through several cavity modes. Although the same mode oscillates at both '1' and '0' states, during the bit transitions, energy is transferred at different wavelengths resulting in a chirp like emission effect which is detrimental for the transmission properties of a directly modulated microring laser.

2.3 Wavelength Conversion

Next the usage of microring laser (MRL) as a wavelength conversion system is shown. The MRL was biased above threshold and emitted at a wavelength λ_o , in single mode operation. Although MRLs are multimode devices in nature, they have proven to provide enhanced spectral purity characteristics due to the absence of spatial hole burning effects. In our case the MRL has a side mode suppression ratio (SMSR) as high as 35 dB without any wavelength selection mechanism.

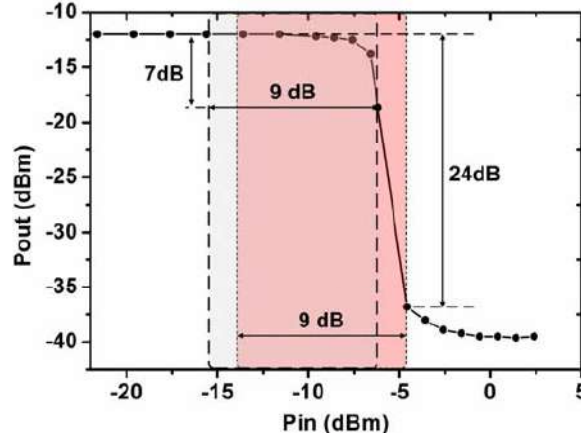


Fig. 5. Output power at 1539 nm against input at 1525 nm. The highlighted areas correspond to an input signal with extinction ratio ER = 9 dB. The darker area has sufficient logical '1' power to switch off the main mode whilst the other suppresses the main mode to produce an output signal with ER = 7dB.

An intensity modulated (IM) signal at a wavelength λ_i was injected into the laser and caused a modulation of the carrier density in the gain section resulting in IM of the output signal at λ_o . Thereby, the information on the incoming signal at λ_i was transferred to the lasing wavelength λ_o . In figure 5, the output power at $\lambda_o = 1539$ nm is plotted against input power at $\lambda_i = 1525$ nm. A complete switch-off of the mode at λ_o is obtained, for input power exceeding -5dBm, which would theoretically lead to an extinction ratio (ER) of 24 dB.

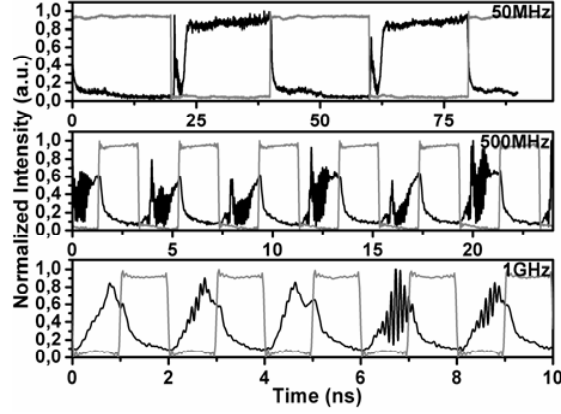


Fig. 6. Transient response (black line) for a clock input signal (gray line) that causes complete mode switching at frequencies: 50 MHz, 500 MHz and 1 GHz (top to bottom).

However, the transient response in this dynamic range is significantly limited by the required switching time of the laser modes which is in the nanosecond order. To better demonstrate this, the transient response of wavelength conversion for a clock input signal with frequencies of 50 MHz, 500 MHz and 1 GHz is shown in figure 6.

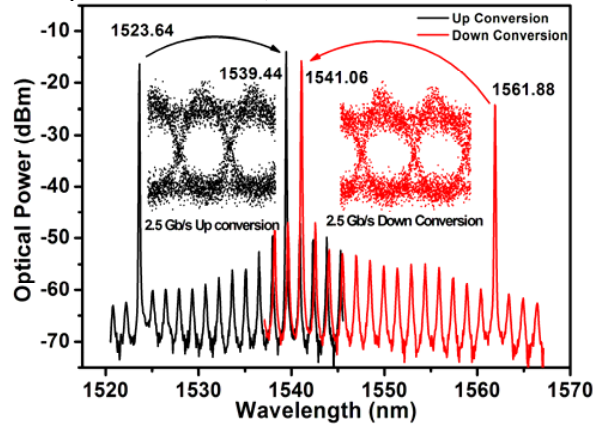


Fig. 7. Wavelength up and down conversion at 2.5 Gb/s. Spectra and eye diagrams with ER approx. 4 dB and high Q factor equal to 9.

The input optical power (P_{in} in figure 5) in these cases (varying approx. between -14 dBm and -5 dBm, as indicated in figure 5) is adequate to cause complete off-switching of the laser mode hence a high ER is measured. In this case the transient response of the output signal is affected by relaxation oscillations and mode competition phenomena occurring during the mode switch-on process. It is evident from the three time-traces that the switching time is in the nanosecond order which

means that high ER, that comes with complete mode switching, cannot be achieved for high data rate applications in the Giga scale without severe distortion.

Nevertheless, wavelength conversion up to 5 Gb/s is possible by not switching off the laser mode completely, though at the expense of lower ER. In this case, the response bandwidth is limited by the gain dynamics associated with the carrier lifetime of the laser [9]. A 2.5 Gb/s modulated input signal at 1537 nm with 9 dB ER (input optical power P_{in} varying between -16 dBm and -7 dBm, as shown in figure 5) is up-converted to one of three successive laser modes ($\lambda_1 = 1539.34$, $\lambda_2 = 1540.74$, $\lambda_3 = 1542.22$ nm). The MRL supports mainly these three modes, attainable by adjusting the current injected in the device. Although this limits the output tunability of the conversion process, various methods that utilize MRLs can be used in future implementations to greatly extend the achievable wavelength coverage. Each of these modes is achievable by properly adjusting the injection current in a repeatable and systematic manner. Nevertheless, each mode is stable within a range of 1 mA having SMSR greater than 30 dB. The switching time required to tune the MRL between modes is measured to be below 1 μ s which is well suited for λ -conversion applications.

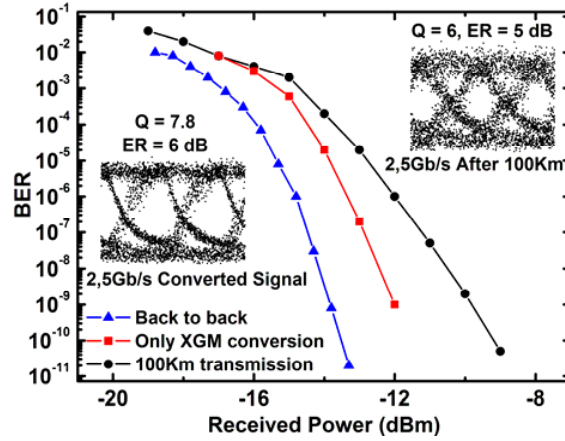


Fig. 8. BER measurements at 2.5 Gb/s for back-to-back, wavelength converted and propagated signal. The corresponding eye diagrams show an ER degradation from 6 to 5 dB and Q factor from 7.8 to 5. The input and output wavelengths in figure 8 are 1552.82 and 1541.06nm respectively.

On the other hand, input tunability is only limited by the gain bandwidth of the device which is approximately 40 nm. The input signal can be tuned to coincide with any MRL mode and the converted output resides at either one of the three laser modes depending on the injection current. Although the output wavelength tuning with injection current will slightly move the resonance peaks of the input wavelength as well, this can be addressed with proper adjustment of the temperature. In figure 7, the input tunability is shown, where successful up- and down-conversion is depicted from the output spectra and eye diagrams. The shown input values 1523.64 and 1561.88 are close to the extreme cases chosen for input tunability showing complete C-band wavelength coverage.

Next results are presented from a transmission experiment of 100 km with SMF at 2.5 Gb/s. In figure 8, the BER measurement and eye diagrams before and after transmission are demonstrated showing low power penalty of 2 dB, for an error-free conversion process. The power penalty after propagation of about 3 dB is reasonably low which indicates that the wavelength converter is suitable for metro and access applications. There is a tradeoff between high ER and Q-factor depending on the input power. This flexibility allows for a higher ER (6 dB) at the cost of Q-factor (7.8) when it comes to propagating the converted signal, assuming always ER of 9 dB for the input signal. In the insets of figure 8, the resulting eye diagrams after wavelength conversion and after the 100 km propagation in SMF are shown. The input and output wavelengths in figure 8 are 1552.82 and 1541.06nm respectively.

3 Conclusions

The modulation properties of vertically coupled microring lasers have been presented [14-15] showing successful -error free- operation up to 7 Gb/s and the potential to operate up to 10 Gb/s. A transmission experiment was also carried out showing error free transmission at 1.25 and 2.5 Gb/s over 50 km of single mode fibre without dispersion compensation proving for the first time the possibility of microring lasers serving as sources in communications systems such as access and metro optical networks.

Wavelength conversion using self pumped cross gain modulation in a microring laser vertically coupled to a passive waveguide was also presented [16]. Successful conversion at 2.5 Gb/s along with a 100 km transmission was demonstrated. Tunable conversion from almost any wavelength of the gain bandwidth to any of the three supported modes of the laser was shown. Future implementations that can accomplish tunable output, i.e. double ring lasers, could greatly enhance the tunability of the conversion process, hence providing a cost effective integrated solution for metro and access optical networks.

More results published or pending publication in other applications other than those shown here can be found in [17-20].

References

1. D. Alexandropoulos, A. Kapsalis, D. Syvridis, "Suppression of higher order modes in vertically coupled micro-ring resonators," *Microwave and Optical Technology Letters* Vol. 49, No. 12, pp. 2963-2968, (2007)
2. M. Kusko, A. Kapsalis, C. Kusko, D. Alexandropoulos, D. Cristea and D. Syvridis, "Design of single-mode vertically coupled microring resonators", *J. Opt. A: Pure Appl. Opt.* 10 (2008)
3. M. Kusko, D. Alexandropoulos, C. W. Tee, A. Kapsalis, D. Cristea, D. Syvridis, C. Kusko, "Numerical analysis of microring resonator obtained by wafer-bonding technology", *Proceedings of the SPIE*, Volume 5956, pp. 349-360 (2005)
4. D. Alexandropoulos, A. Kapsalis, D. Syvridis, U. Troppenz, M. Hamacher H. Heidrich, "Design considerations for spatial monomode operation of InP-based passive vertically

- coupled microring resonators", Proceedings of the 5th International Conference on Numerical Simulation of Optoelectronic Devices, NUSOD '05, pp. 119- 120 (2005)
5. M. Kusko, A. Kapsalis, C. Kusko, D. Alexandropoulos, D. Cristea, D. Syvridis, "Design of single-mode vertically coupled microring resonators," 2nd European Optical Society Topical Meeting: "OPTICAL MICROSYSTEMS ", 30 Sept. – 3 Oct. 2007, Capri (Napoli), Italia, post dead-line paper
 6. I. Stamataki, S. Mikroulis, A. Kapsalis, D. Syvridis, "Investigation on the Multimode Dynamics of InGaAsP/InP Microring Lasers", IEEE Journal of Quant. Electronics, vol. 42, no. 12, pp. 1266-1273(2006)
 7. A. Kapsalis, I. Stamataki, S. Mikroulis, D. Syvridis, M. Hamacher, "Widely Tunable All-Active Micro-Ring Lasers", IEEE Phot. Techn. Lett., vol. 18, no. 24, pp. 2641-2643 (2006)
 8. I. Stamataki, A. Kapsalis, S. Mikroulis, D. Syvridis, M. Hamacher, U. Troppenz and H. Heidrich, "Modal properties of all-active InGaAsP/InP microring lasers", Optics Communications, Volume 282, Issue 12, p. 2388-2393 (2009)
 9. A. Kapsalis, D. Alexandropoulos, S. Mikroulis, H. Simos, I. Stamataki, and D. Syvridis, M. Hamacher, U. Troppenz, and H. Heidrich, "Spectral properties of all-active InP-based microring resonator devices", Proceedings of the SPIE, Volume 6115, pp. 386-394 (2006)
 10. A. Kapsalis, I. Stamataki, D. Syvridis, M. Hamacher, "Widely Tunable All-Active Micro-Ring Lasers through Phase Shifted Controlled Feedback", We 3.52 ECOC 2006
 11. M. Hamacher, U. Troppenz, H. Heidrich, V. Dragoi, A. Kapsalis, D. Syvridis, C.W. Tee, K.A. Williams, M. Alexe, M. Kusko, D. Cristea, "Vertically coupled microring laser devices based on InP using BCB waferbonding," CLEO-Europe IQEC 2007, 17-22 June 2007, Munich, Germany
 12. M. Hamacher, H. Heidrich, U. Troppenz, D. Syvridis, D. Alexandropoulos, S. Mikroulis, A. Kapsalis, C.W. Tee, K.A. Williams, V. Dragoi, M. Alexe, D. Cristea, M. Kusko, "Waferbonded Active/Passive Vertically Coupled Microring Lasers," SPIE Photonics West - LASER 2008, January 19-24, 2008, San Jose, California, USA, paper no. 6896-27
 13. M. Hamacher, U. Troppenz, H. Heidrich, V. Dragoi, A. Kapsalis, D. Syvridis, C.W. Tee, K.A. Williams, M. Alexe, M. Kusko, D. Cristea, "Vertically Coupled and Waferbonded μ Ring Resonators on InP," European Semiconductor Laser Workshop 2007, 14 – 15 Sept. 2007, Berlin, Germany
 14. A. Kapsalis, U. Troppenz, M. Hamacher, H. Heidrich and D. Syvridis, "7Gb/s Direct Modulation of Vertically Coupled Microring Lasers," OFC 2008, 24-28 February 2008, San Diego, California, USA
 15. D. Syvridis, H. Simos, S. Mikroulis, and A. Kapsalis, "Microring-based devices for telecommunication applications," Proceedings of the SPIE, Vol. 7211, (2009)
 16. A. Kapsalis, H. Simos, D. Syvridis, M. Hamacher, H. Heidrich, "Tunable Wavelength Conversion using Cross Gain Modulation in a Vertically Coupled Microring Laser," IEEE Phot. Techn. Lett., vol. 21, no. 21, pp. 1618-1620 (2009)
 17. A. Kapsalis, D. Syvridis, M. Hamacher, H. Heidrich, "Broadly Tunable Laser using Double-Rings Vertically Coupled to a Passive Waveguide", IEEE Journal of Quantum Electron., vol. 46, no. 3, pp. 306-312 (2010)
 18. A. Kapsalis, C. Mesaritakis, D. Syvridis, "Design and Experimental Evaluation of Active-Passive Integrated Microring Lasers", ESWL Lauzanne (2011)
 19. A. Kapsalis, I. Stamataki, C. Mesaritakis, D. Syvridis, M. Hamacher, H. Heidrich, "Design and Experimental Evaluation of Active-Passive Integrated Microring Lasers: Threshold Current and Spectral Characteristics", submitted for publication to IEEE Journal of Quantum Electron. (2011)
 20. A. Kapsalis, I. Stamataki, C. Mesaritakis, D. Syvridis, M. Hamacher, H. Heidrich, "Design and Experimental Evaluation of Active-Passive Integrated Microring Lasers: Noise Properties", submitted for publication to IEEE Journal of Quantum Electron. Semiconductor Optoelectronic Devices special issue (2011)

BPEL scenario execution: QoS-based dynamic adaptation and exception resolution

Christos Karelitis¹

National and Kapodistrian University of Athens
Department of Informatics and Telecommunication
ckar@di.uoa.gr

Abstract. BPEL/WSBPEL is the main approach for combining individual web services into integrated business processes. A BPEL/ WSBPEL scenario allows for specifying which services will be invoked, their sequence, the control flow and how data will be exchanged between them. BPEL however does not include mechanisms for considering the invoked services' Quality of Service (QoS) parameters and thus BPEL scenarios cannot customize their execution to the individual user's needs or adapt to the highly dynamic environment of the WEB, where new services may be deployed, old ones withdrawn or existing ones change their QoS parameters. Moreover, infrastructure failures in the distributed environment of the web introduce an additional source of failures that must be considered in the context of QoS-aware service execution. In this thesis, it is proposed a framework for addressing the issues identified above; the framework allows the users to specify the QoS parameters that they require and it undertakes the task of locating and invoking suitable services. In this dissertation two strategies for selecting the most suitable service are considered: (a) a greedy strategy and (b) a partner link-level strategy. The proposed framework intercepts and resolves faults occurring during service invocation, respecting the QoS restrictions specified by the consumer. The latter also intercepts and resolves faults occurring during service invocation, respecting the QoS restrictions specified by the consumer. Finally, methods for tackling with syntactic differences between functionally equivalent services, broadening thus the pool of available services for each adaptation are considered. Finally, performance metrics for the proposed framework are presented, which validate its applicability to operational environments and present performance metrics for the proposed framework.

1 Introduction

Web services have emerged as a new standard, having as main focus to allow applications over the Internet to communicate with each other, which are independent of execution platform, programming language and implementation details. The web service paradigm has been adopted by research community and industry alike, however a number of challenges still lie ahead for fully covering the needs of both service

¹ Dissertation Advisor: Panayiotis Georgiadis, Professor

providers and consumers. [1] identifies a number of open issues in the current SOA state-of-the-art, spanning across four major categories namely service foundations (service oriented middleware backbone that realizes the runtime SOA infrastructure), service composition, service management and monitoring as well as service design and development. For service governance, in particular, [1] lists “service governance” as a major research challenge, stating that the potential composition of services into business processes across organizational boundaries can function properly and efficiently only if the services are effectively governed for compliance with QoS and policy requirements. Services must meet the functional and QoS objectives within the context of the business unit and the enterprises within which they operate.

In this context, development procedures as well as composition and execution mechanisms need to take into account the QoS dimension of web services in order to formulate successful business processes that will satisfy users’ (either business or individuals) expectations. Regarding service composition into business processes, the predominant approach used nowadays is the formulation of BPEL/WSBPEL scenarios [2], in which the BPEL designer specifies the business process logic; this includes invocation of selected web services, control flow constructs and data flow arrangements in the form of result gathering and parameter passing, while provisions for exception handling (such as service unavailability or business logic faults) also exist.

BPEL scenarios, however, do not include facilities either for specifying QoS parameters for services, or for dynamically selecting the web service to be called at runtime, therefore the BPEL scenario designer must select the concrete service implementation to be invoked in the context of the business process while creating the scenario, by examining the QoS parameters of functionally-equivalent services. This alternative, however, is not a viable one since (a) the same BPEL scenario may be used by different users with diverging or even contradictory requirements and (b) even if the “best choice” is made at some time point there is no guarantee that this choice will continue to be optimal in the future. Moreover, in the presence of failures, it would be desirable for the system to be able to locate and use “second best” choices automatically, provided that they deliver the required functionality and satisfy QoS restrictions.

2 Summary

2.1 Motivation and Challenges

The main objective of web service technology and related research [3] is to provide the means for enterprises to do business with each other and provide joint services to their customers under specified Quality of Service (QoS) levels. The collaboration of web services, possibly provided by different companies, in order to create composite and potentially highly complex business process, elevates the need of a Business Process Management (BPM) [4], [5]. Trying to model real world business processes with BPEL scenarios a series of challenging issues may emerge. Specifically:

1. processes may be long-running, in the order of hours, days or even longer. Such issues commonly arise in cases where human intervention is required for the completion of all or some of the services that comprise the process.
2. BPEL scenarios may try to model stable and established processes that remain relatively unchanged. Examples of such processes are those that represent interactions with Government-based services, spanning the range of G2x and x2G acronyms
3. as the complexity of the process and the number of cooperating services needed increase, so does the volatility of these services. New services implementing the same process may appear, existing ones may be decommissioned or the BPEL designer may not be aware of all the services that can be utilized at the time of the designing phase
4. quality requirements for the process may change during the lifetime of the BPEL scenario. This may be due to different needs of end-users (a real-world counterpart of this case is one person sending a package using courier mail to minimize delivery time, whereas another person may use ordinary surface mail to pay less), or alterations in organizational policy.

2.2 Problem Identification and Objective

In cases such as the above, the static nature of BPEL scenarios and their handling of BPEL engines fail to accommodate for the dynamic nature of real world processes. To cope with these situations, the BPEL scenario would have to be redesigned and re-deployed possibly forcing existing transactions to fail or be re-started. For accommodating different needs of end-users, the alternative approach of maintaining different versions of the BPEL scenarios could be also taken, with each version being targeted to a specific user category (e.g. “express delivery” vs. “economic delivery”); this arrangement, however, would increase development and maintenance costs and would weaken the overall system manageability.

To tackle these issues, this dissertation proposes an approach that is relying on *dynamic service selection mechanism based on functional and non-functional (quality) criteria* for selecting the most suitable service *per scenario invocation*. Furthermore, this mechanism provides for non-existent or invalidated services allowing them to be replaced with existent and valid ones, choosing the optimal candidate per service invocation based on current criteria. The criteria can be different on each run and can provide for diverse needs depending on the invoker.

So, the basic features and innovations this dissertation introduced were:

- the concept of replacement candidate for web services was formalized considering criteria related to the specific BPEL scenario execution, instead of the generic functionality or behavior of the service. Replacement candidates are used for hot-swapping failed services within a BPEL transaction, allowing thus the BPEL scenario to complete its execution. The formalization introduced allows including more services in the “replacement candidate” pool and therefore formulating execution paths with better qualitative characteristics.

- the notion of service selection affinity was introduced, which allowed for maintaining the transactional characteristics of BPEL scenarios in the presence of adaptation
- an approach to bridging the syntactic differences between functionally equivalent services was proposed, which greatly enhances the maintainability of the equivalent services repository, trading off a degradation in performance, which has been quantified to be quite small.
- a method for distinguishing between system faults and business logic faults was proposed; this distinction is important since faults in the former category can be resolved by automatically invoking a replacement candidate for the failed service, while this is not possible for faults in the second category.
- a framework that enables the automatic resolution of system faults and the dynamic adaptation of BPEL scenario execution according to QoS criteria was proposed. The framework is independent of the particular BPEL execution engine used, and methods have been proposed for setting the QoS criteria granularity (for all scenarios executing in the system; for the scenario as a whole; for each individual service within a scenario). This framework includes provisions for maintaining the transactional characteristics of BPEL scenario execution, making use of the *service selection affinity* notion.
- the feasibility of the above was proved through a complete system implementation and quantification of its performance.
- the issue of BPEL scenario adaptation in the context of secure web services invocation was identified, and a system architecture for a system that supports such an adaptation was drafted.

2.3 Related Work

In this section some related work is adduced in the following research directions:

QoS management in web services composition:

In [6] a framework is presented named AgFlow [7] as middleware platform that enables the quality-driven composition of Web services. In AgFlow, the QoS of Web services is evaluated by means of an extensible multidimensional QoS model. It presents two selection policies: the local optimization of individual tasks and a global planning. The first is similar to the one proposed in this thesis and it uses the Simple Additive Weighting [8] technique to select the optimal service for a given task. The proposed approach differentiate from this since we deal with already defined composition scenario and doesn't propose a re-planning solution method in order to change the task execution order, or replace a set of task with another set. It uses a proxy-like service that is invoked for each individual task in the business scenario in order to discover the optimal services for each one of them based on a specific consumer's quality policy at execution time.

In [9] a web service proxy is introduced in order to perform a dynamic binding of related web services under specified user's constraints. The selection of equivalent services is not only filtered by constraints but also it is measured the quality score for each equivalent service depending on a quality vector and a set of quality weights. In [10] the importance of qualify-able QoS aspect related to the issue of web services

composition and monitoring is illustrated. It describes an algorithm capable of capturing and reflecting the state of web services involved in the integration process.

In exception management web services composition:

In this research work [11] a policy-driven approach is introduced to exception management. An exception handling policy language is designed, which defines deviation situations and the associated exception handlers. The proposed approach complements the above solution by discovering an optimal alternate service task to perform the alternative action mentioned. A remarkable research in this area has been and the one introduced in [12]. It's presenting a component called BPBot (Business Process roBOT). A business process is executed by a collection of BPBots that are dynamically organized as a hierarchical structure. The proposed solution is not re-planning an execution path, but it discovers functionally and qualitatively equivalent services to perform the determined business tasks without changing the task execution sequence. Moreover, during this dissertation the author published relative papers ([20], [21], [22], [23], [24]) considering service BPEL scenario adaptation in the context of exception resolution and security issues in exception handling in [25].

Semantic Web Services:

In the past few years, the issue of exception resolution in composite web services has drawn the researchers' attention. A noteworthy approach to exception handling is the one undertaken by METEOR-S project [13], [14] in cooperation with WSMX (Web Services Execution Environment) [15]. WSMX contains the discovery component, which undertakes the role of locating the services that fulfill a specific user request. This task is based on the WSMO conceptual framework for discovery [16]. WSMO includes a Selection component that applies different techniques ranging from simple "always the first" to multi-criteria selection of variants (e.g., web services non-functional properties as reliability, security, etc.) and interactions with the service requestor. Both in the METEOR-S and other approaches, functional and non-functional properties are represented using shared ontologies, typically expressed using DAML+OIL [17] and the latter OWL-S. Such annotations enable the semantically based discovery of relevant web services and can contribute towards the goal of locating services with "same skills" [18] in order to replace a failed service in the process flow. The main difference of the research illustrated with the one referenced above is that selection of replacements to services that have failed within an execution plan is made dynamically, instead of using pre-determined exception resolution scenarios. Replacement service selection is based on both functional equivalence (performed through semantic matching) and qualitative replaceability (considering non-functional attributes). Furthermore, qualitative replaceability criteria may be defined by the composite service invoker, to more accurately specify which replacement service is the most suitable one in the context of the current execution.

2.4 Brief Description

Service Quality Vectors

In order to enable the selection of the "most suitable" operation according to some QoS specification, the QoS attributes of the operations should be represented in an unambiguous and system-processable format, while additionally means for expressing QoS-related operation selection criteria should be afforded. For brevity, in the

following we will consider only the QoS parameters cost, security, performance, response time and availability, adopting the definitions in [19]. For each such source, mappings between the domains employed by the source and numeric values are used.

Table 1. Mapping of QoS values

	QoS provider 1	QoS provider 2	Value
Cost	10 €	11 €	1
Security	6 (out of 10)	62 (out of 100)	3
Performance	High throughput	99%	5
Response time	0.0001 ms	Real-time	1
Availability	High	> 95%	4

In the approach illustrated here three vectors that define the QoS criteria for process invocation are considered; in other words it is defined a QoS specification as a triple (MAX, MIN, W), where MAX, MIN and W are quality vectors (defined below). The quality vector for the QoS attributes considered in this work can be defined as:

Table 2. Quality Vector

MAX =	(costmax, secmax, perfmax, respmax, availmax)
MIN =	(costmin, secmin, perfin, respmin, availmin)
W =	(costw, secw, perfw, respw, availw)

ASOB-Framework

Figure 2 illustrates the overall architecture of our approach to dynamic policy-driven execution of a business scenario QoS-aware and policy-adhering exception management techniques. The component undertaking this responsibility is the Alternative Service Operation Bind (ASOB).

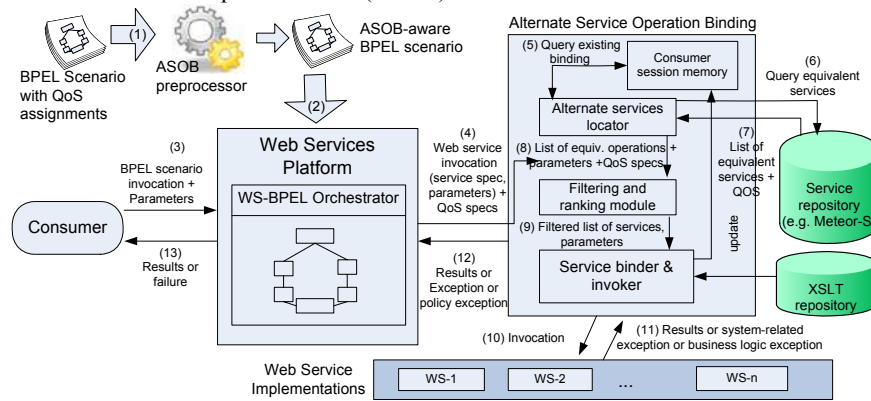


Fig. 1. Overall System Architecture

The BPEL scenario (SC) as crafted by the BPEL designer is processed by the *ASOB preprocessor*, which produces an *ASOB-aware BPEL scenario* (SC_{ASOB}) as output, so that for each service, the ASOB middleware calculates an overall score which takes into account all the operations of the service that are listed in the BPEL scenario and the respective QoS weights that the client has specified at the pre-processing phase.

$$Sc_{WS} = \sum_{op \in Ops} \sum_{attr \in \{cost, sec, \dots\}} attr_{WS, op} * QoS_w(op)_{attr} \quad (1)$$

Depending on the score of each service, in case of a failure, ASOB replaces the failed one with the service that owns the highest score Sc . The interested reader will find in more depth the main processing of the ASOB framework at the main dissertation text.

3 Results and Discussion

The contribution of the ASOB framework to the field is as follows:

1. it allows the BPEL scenario designer to specify the desired QoS parameters for each service. These parameters are specified through standard BPEL variables, thus the designer may examine scenario input parameters for setting them, tuning thus the adaptation of the particular BPEL scenario execution to the desires and needs of the scenario consumer.
2. it does not require any modification to the BPEL syntax or semantics.
3. it takes the execution flow specified by the designer as granted, and optimizes service selection within this flow, contrary to service composition approaches which define this flow dynamically. This is an important aspect in cases where execution flow is carefully crafted by the designer to reflect particularities of the business process, specialized exception handlers are used, etc.
4. it incorporates exception handling as an integral part of the adaptation process, allowing for switching to the “next best” solution when the originally selected candidate is unavailable.
5. it does not use pre-determined alternative paths, but selects services dynamically from a suitable registry.
6. It employs XSLT transformations through which the middleware bridges the syntactic differences between the service originally specified in the BPEL scenario and other services that are semantically equivalent but syntactically different. This arrangement offers to the middleware a wider range of choices, for the stage of deciding which service provider best matches the QoS specifications given in the BPEL scenario.
7. it considers service selection affinity, enabling the conducting of multi-operation transactions with providers.

8. it introduces the notion of the service replacement candidate, which relaxes the requirements for service equivalence. Service replacement candidates are computed for the context of a particular BPEL scenario and takes into account only the operations used in the scenario and not all operations offered by the services. This arrangement enables the middleware to avoid cases where some operation that is not used in a scenario breaks the equivalence of two services, and thus disallows the consideration of some alternates.
9. it elaborates on the management of consumer session memory, which supports the maintenance of service selection affinity.
10. it provides full details for the algorithms used by the middleware to process web service invocations.
11. it includes a partner link-level strategy for deciding which is the service provider that best matches the QoS profile specified in the BPEL scenario; the partner link-level strategy can significantly improve the service provider selection when a BPEL scenario uses multiple operations from the same service provider, while it may also prevent some cases where the greedy strategy is unable to find any appropriate execution path for servicing the scenario.

Algorithms in pseudo-code can be found in the main text of this dissertation.

3.1 Performance Evaluation and Results

Figure 2 illustrates the ASOB internal process time for single web service operation invocations, against the overall service repository (SR) size and the number of equivalent services present in the repository. The overhead increment, on the other hand, when the number of alternate services increases is considerable, mainly affecting the sorting of the candidate operation list (typically of complexity $O(n * \log(n))$).

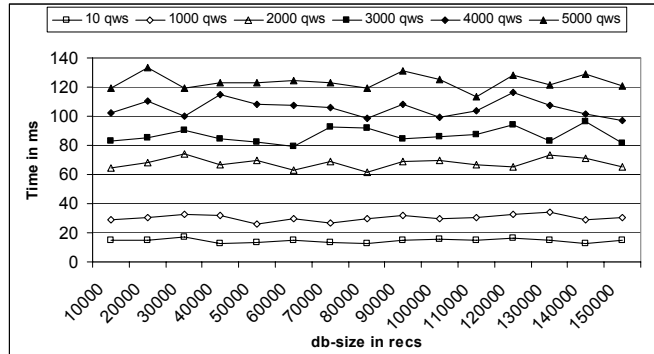


Fig. 2. ASOB internal process time

Table 3. XSLT transformation overhead

concurrent ASOB invocations	20	40	60	80	100
time in msec (average per transformation)	17.8	18.5	34.5	46.2	61.7

Table 3 shows the overhead incurred by applying XSLT transforms on request and response SOAP messages, to resolve syntactical differences between operations that are *semantically* but not *syntactically* equivalent

Figure 3 illustrates the number of operation invocations that can be served in a unit of time against the number of concurrent invocations when (a) services are directly invoked and (b) when invocations are made through the ASOB middleware.

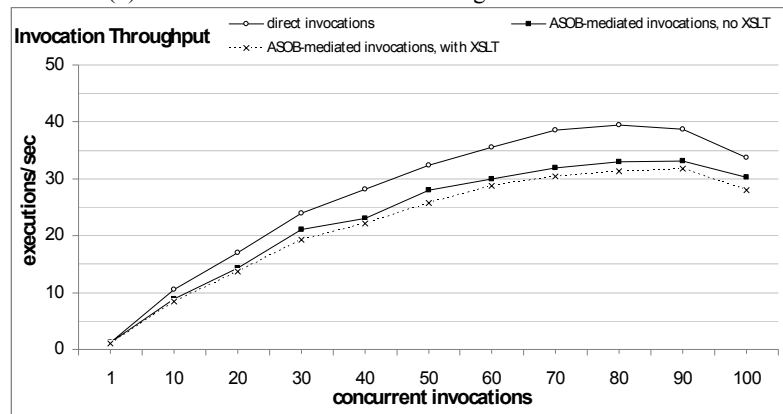


Fig. 3. Invocation throughput

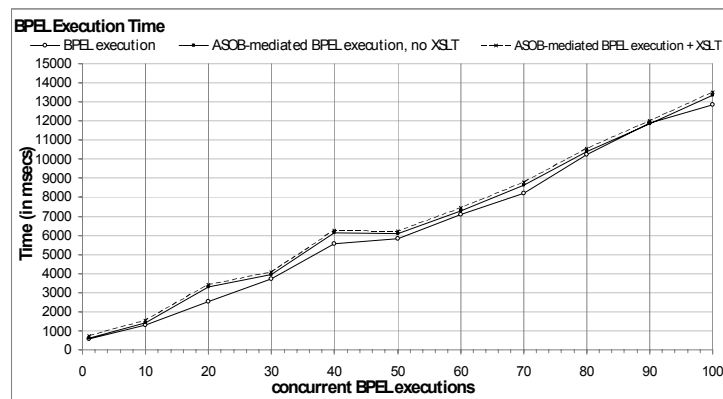


Fig. 4. BPEL scenario execution time

Figure 4 illustrates the BPEL execution time of a BPEL scenario containing two web service invocations against the number of concurrent executions. The increment is very small (4%-9% without XSLT transformations, 8-16% with XSLT transformations).

Figure 5 depicts the BPEL scenario execution throughput against the number of concurrent executions. The behavior is consistent with the previous diagrams.

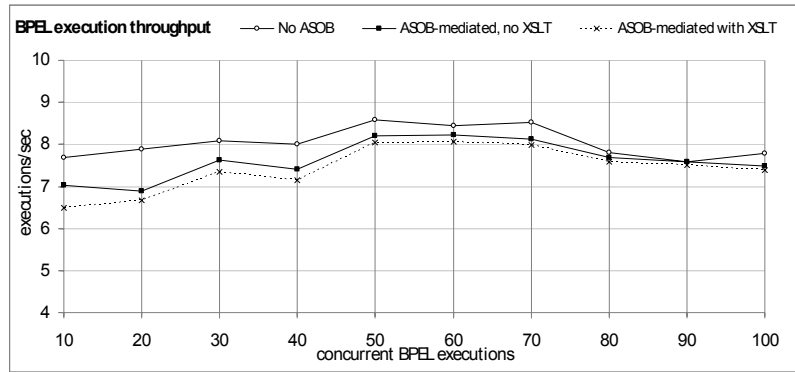


Fig. 5. ASOB-mediated vs. direct invocation BPEL scenario execution throughput

4 Conclusions

Building processes that are able to cope with the dynamics of real world requirements has always been a challenging endeavor. The adoption of BPEL in the design and execution phases of business processes has already obtained gains in speed and reliability, but has not been able insofar to successfully address issues arising from the dynamic nature of the processes themselves, the diversity in user requirements and the inherent instability of distributed environments, which leads to a number of system faults.

The framework presented in this dissertation addresses these shortcomings by employing a dynamic service selection mechanism based on QoS criteria for a BPEL process; these criteria are defined by the BPEL scenario designer and can be set to reflect the end-user requirements. Service attributes are stored in a repository that stores the services' functional and non-functional (qualitative) characteristics. Updating the repository suffices to reflect changes in the real world (service introductions or withdrawals, changing of services' QoS aspects etc). An exception resolution mechanism for faults owing to systemic reasons is also included, easing thus the work of the BPEL designer.

The strategy employed by the presented framework for binding a partner link to a specific service provider can follow either (a) a *greedy strategy*, according to which the QoS aspects of only the first operation invoked for a particular partner link are examined to determine the binding or (b) a *partner link-level strategy*, which reviews all invocations collectively, avoiding suboptimal bindings and cases where the greedy strategy leads to inability to successfully conclude the BPEL scenario.

Open issues in this field includes a detailed evaluation of the partner link-level strategy regarding (a) its performance, i.e. the time needed to determine the optimal binding for a partner link and (b) the quality of the execution plans it produces. Execution plan quality is a twofold aspect involving (i) the degree to which the bindings performed by the middleware correspond to the QoS specifications listed in the BPEL scenario and (ii) the number of cases where the partner link-level strategy bindings lead to successful execution of the BPEL scenario, contrary to the bindings of the greedy algorithm. Moreover, it could be investigates the collection and exploitation of statistics regarding the number of invocations for each particular operation in the context of a specific BPEL scenario, so as to use a more elaborate weight assignment in the phase of calculating the suitability scores of different bindings.

References

1. M. P. Papazoglou, P. Traverso, Leymann, Service-Oriented Computing: State of the Art and Research Challenges. IEEE Computer (40) 11, Nov. 2007, pp. 38-45.
2. M. Juric, Business Process Execution Language for Web Services BPEL and BPEL4WS (2nd Edition), Packt Publishing, 2006, ISBN-10: 1904811817.
3. Newcomer, E., Lomow, G.: Understanding SOA with Web Services, Addison-Wesley, (2005)
4. F. Leymann, D. Roller, and M. T. Schmidt, Web services and business process management, Available at: <http://www.research.ibm.com/journal/sj/412/leymann.html>
5. Martin Hepp, Frank Leymann, John Domingue, Alexander Wahler, and Dieter Fensel Semantic Business Process Management: A Vision Towards Using Semantic Web Services for Business Process Management, IEEE International Conference on e-Business Engineering, 2005, p:535-540
6. Liangzhao Zeng, Boualem Benatallah, Anne H.H. Ngu, Marlon Dumas, Jayant Kalagnanam, and Henry Chang. Qos-aware middleware for web services composition. IEEE Trans. Softw. Eng., 30(5):311–327, 2004.
7. L. Zeng, Dynamic Web Services Composition, PhD thesis, Univ. of New South Wales, 2003.
8. H.C.L and, K. Yoon, Multiple Criteria Decision Making,” Lecture Notes in Economics and Mathematical Systems. Springer-Verlag, 1981.
9. K. Verma, R. Akkiraju, R. Goodwin, P. Doshi, J. Lee, On Accommodating Inter Service Dependencies in Web Process Flow Composition, AAAI Spring Symposium PP: 37-43 on Semantic Web Services.
10. Hassan Issa, Chadi Assi, Mourad Debbabi, QoS-Aware Middleware for Web Services Composition - A Qualitative Approach, in Proceedings of the 11th IEEE Symposium on Computers and Communications, 2006
11. Liangzhao Zeng; Hui Lei; JunJang Jeng; Jen-Yao Chung; Benatallah, B. *Policy-driven exception-management for composite Web services*, *E-Commerce Technology*, 2005. CEC 2005. 7th IEEE International Conference on Volume , Issue , 19-22 July 2005, 355 – 363
12. Liangzhao Zeng, JunJan Jeng, Santhosh Kumaran and Jayant Kalagnanam, Reliable Execution Planning and Exception Handling for Business Process, LNCS, Springer, Technologies for E-Services, 2003. p.119-130
13. Kochut, K. J.: METEOR Model version 3. Athens, GA, Large Scale Distributed Information Systems Lab, Department of Computer Science, University of Georgia (1999)

14. K. Verma, K. Sivashanmugam, A. Sheth, A. Patil, S. Oundhakar, and J. Miller, METEOR-S WSDI: A Scalable Infrastructure of Registries for Semantic Publication and Discovery of Web services. *Journal of Information Technology and Management*, Special Issue on Universal Global Integration, Vol. 6, No. 1 (2005) 17-39
15. Cimpian, E., Moran, M., Oren, E., Vitvar, T., Zaremba, M.: Overview and Scope of WSMX. Technical report, WSMX Working Draft, <http://www.wsmo.org/TR/d13/d13.0/v0.2/>
16. D. Roman, D2v1.2 Web Service Modeling Ontology (WSMO). WSMO Final Draft April 13, 2005. Available at: <http://www.wsmo.org/TR/d2/v1.2/20050413/>
17. DAML+OIL, Available at: <http://www.daml.org/2001/03/daml+oil-index.html>
18. Dellarcas, C. and M. Klein, A knowledge-based approach for handling exceptions in business processes, *Information Technology and Management* 2000.
19. O'Sullivan J., Edmond D., and A. Ter Hofstede (2002), What is a Service?: Towards Accurate Description of Non-Functional Properties, *Distributed and Parallel Databases*, 12.
20. Kareliotis C., Vassilakis C., Rouvas E., Georgiadis P. (2008), Exception Resolution for BPEL Processes: a Middleware-based Framework and Performance Evaluation. *Procs of iiWAS 2008*, Linz, Austria.
21. Kareliotis C., Vassilakis C., Georgiadis P. (2007), Enhancing BPEL scenarios with Dynamic Relevance-Based Exception Handling, *Proceedings of the ICWS 2007*, pp.751-758.
22. Kareliotis C., Vassilakis C., Rouvas E., Georgiadis P. (2009), QoS-Driven Adaptation of BPEL Scenario Execution. *Procs of ICWS 2009*, July 6-10, 2009, Los Angeles, CA, USA.
23. Christos Kareliotis, Costas Vassilakis, Stathis Rouvas, Panagiotis Georgiadis. QoS-Aware BPEL Scenario Execution and Adaptation: A Middleware-Oriented Framework. Extended version invited from ICWS09 in *International Journal on Web Services Research. JWSR*. 2009.
24. Kareliotis Christos, Vassilakis Costas, Georgiadis Panagiotis: Towards Dynamic, Relevance-Driven Exception Resolution in Composite Web Services. *Proceedings of International Conference on Object Oriented Programming, Systems, Languages and Applications* 2006.
25. Costas Vassilakis, Kareliotis Christos: A framework for adaptation in secure web services. *4th Mediterranean Conference on Information Systems* 2009. MCIS, Athens, Greece

Study and Design of Algorithms for Information Dissemination in Unstructured Networking Environments

Dimitrios Kogias*

Department of Informatics and Telecommunications
National and Kapodistrian University of Athens
Greece
Email: dimkog@di.uoa.gr

Abstract. The focus on this thesis lies on the study of several information dissemination techniques in modern unstructured networks. Information dissemination has an important academic interest in these environments due to the special characteristics (e.g., decentralization, large-scale and dynamic nature) that they possess. One of the techniques that is studied is *probabilistic flooding* and analytic (asymptotic) bounds on the value of the forwarding probability p_f , for which a probabilistic flooding network manages to fully cover an underlying connected random graph are presented. The technique of *multiple random walkers* is, also, studied given analytical expressions regarding coverage and termination time for fully and less dense connected topologies. The observation that the network is not efficiently covered at the early stages due to the (potentially) large collection of walkers at the initiator node has led to the introduction of a new information dissemination mechanism that creates walkers (by *replicating* the existed ones) during the random walks and not from the beginning of the process. This replication technique, called *Randomly Replicated Random Walkers (RRRW)*, has been studied in various networking topologies (e.g., random geometric graphs, power-law graphs, clustered graphs) to examine whether it can fill the gap between the performance of two well-known techniques, the *Full Flooding* and *Single Random Walker (SRW)* in *stretching* the advertising information in broader networking areas.

1 Introduction

In modern networking environments, the discovery of a given piece of information plays a key role to its robustness and functionality. In particular, easy and quick access to any information source that is needed is expected, along with the possibility of sharing it with more network nodes located in further network areas.

In this thesis, a study of various information dissemination techniques (existed and newly introduced ones) will take place in unstructured networking

* Dissertation Advisor: Ioannis Stavrakakis, Professor

environments. These environments are mainly distributed topologies, where all the nodes are equal and are characterized by their large size, their scalability properties and their highly dynamic nature. Due to these inherent characteristics, it is not possible for a node to possess information regarding the global network structure in any given time. On the other hand, a node will always know the number and the identity of its one hop neighbours. This lack of knowledge, regarding the structure of the overall topology (even though it exists), by any network node is the reason for calling these networks *unstructured*. P2P and ad-hoc networks possess the described characteristics and will be, therefore, considered and studied throughout this thesis.

Because of the aforementioned characteristics, the process of disseminating information in an unstructured environment becomes very difficult and costly (both in number of messages and in termination time) but, in the same time, it is the reason for its high academic and research interest. In this thesis, the information that is disseminated is related to the knowledge of the location in the network of a node that possess a certain service. Therefore, it is considered as part of a larger process that is called *Service Discovery*.

The process of Service Discovery will be divided in two phases: *Service Advertising* and *Service Searching*, that are sequentially applied. During service advertising, a node advertises (using an information dissemination technique) his location, along with the service he possesses, to part (or even whole) of the network by constructing an *advertising network* (i.e., a network consisted by the nodes and links used for advertising). During searching, the node first checks whether he is part of the advertising network of the desired service and if he is then the searching is successfully completed. If not, then he employs an information dissemination technique to search for either a node that is part of the advertising network or the node that possesses the service itself.

It can be easily verified that this two service discovery phases are complementary in nature and that when the intensity of one is small then the intensity of the other should be large.

1.1 Related Work

Various techniques have been proposed so far for disseminating information in a network. The most popular ones will be presented here.

The simplest technique used for information dissemination (both for advertising and for searching) is *traditional flooding*, [1]. Under traditional flooding, the information messages traverse all network links and, thus, visit all the nodes in the network, producing a large number of messages, especially when the network's size (i.e., N) increases. *Termination time*, on the other hand, is significantly small, upper bounded by the network *diameter*, typically of the order of $\log(N)$.

A popular variation is the *controlled flooding* technique, which employs flooding but only for a number of K hops (i.e., K is a Time- to- Live, *TTL*, value) away from the node that initiates the process (called the *initiator node*). While the value of K remains small, the number of messages produced also remains

small and so does the size of the network covered by this technique, reducing the probability for efficient discovery of the desired node or service in a large-scale, modern environment. The controlled flooding (or *K-flooding*) technique is known to be used for searching in the Gnutella P2P system, [2].

On the other hand, techniques like *random walks*, [3], are very different than flooding. These approaches manage to reduce the total number of its messages by sending a limited number (one in the case of the *Single Random Walker* or m in the case of *Multiple Random Walkers*) of entities (as special messages) to cover the network. Each of these entities follows its own path by randomly selecting the next node to visit from one of the 1-hop neighbours of the node that the entity resides in each time slot. The termination of the algorithm takes place after some predefined time (e.g., using TTL expiration), for each entity, or after checking with the initiator node and learning that the desired information has already been discovered by another entity. A combination of the aforementioned termination conditions can also be applied.

Hybrid probabilistic techniques (e.g., local flooding process initiated after a random walk) have also been proposed and analysed, [4], as well as other schemes that adapt the employed TTL values in a probabilistic manner, [5]. Another modification, [6], allows for network nodes to forward messages to their neighbours in a random manner, thus significantly reducing the number of messages in the network. Many other works have been published proposing the selective forwarding of a certain message in the network, e.g., [7, 8].

As it has already been mentioned the main scope of this thesis is the study of several existed and newly introduced information dissemination techniques. Therefore, the rest of this work is organized as follows: in Section 2 the study of the *probabilistic flooding* technique will be presented and an asymptotic analysis following the bounds on the value of the forwarding probability p_f will take place. In Section 3 the technique of *multiple random walkers* is studied giving analytical expressions regarding coverage and termination time for fully and less dense connected topologies. For more efficient coverage of the network during the early stages, a new information dissemination mechanism that creates new entities (by *replicating* the existed ones) is introduced. In Section 4 the technique of *Randomly Replicated Random Walks (RRRW)*, is studied in various networking topologies (e.g., random geometric graphs, power-law graphs, clustered graphs) to examine whether it can fill the gap between the performance of two well-known and used techniques, the *Full Flooding* and *Single Random Walker (SRW)* in stretching the advertising information in wider network areas.

2 Probabilistic Flooding

When *probabilistic flooding* is applied every node, by receiving the message for the first time, forwards it to each of his 1 hop neighbours (apart from the one(s) that forwarded it to him) using a constant forwarding probability $p_f(N)$. By properly parametrized the value of this probability $p_f(N)$ it is possible to cover the underlying graph while producing a smaller number of messages than the

traditional flooding approach. The cost to be paid for this reduction is that the coverage is no longer deterministically guaranteed but is rather probabilistically achieved.

Each time probabilistic flooding is applied to a graph, a *probabilistic flooding network* is created. This is a connected network which contains the number of nodes and links over which the disseminated information has been forwarded. The main scope for using probabilistic flooding is the creation of a probabilistic flooding network (i.e., $P(G(N, p), p_f)$) that contains the minimum number of links/ messages, while still achieving the desired network coverage.

In this thesis, the use of probabilistic flooding is studied when the underlying graph is a random graph. In fact, it is noticed that there is a connection between the stages for the creation of a random graph $G(N, p)$ (which follows the binomial model and, thus, has an expected number of links that equals $p(N) \frac{N(N-1)}{2}$, [9]), depending on the value of $p(N)$, and the stages for the creation of the probabilistic flooding network, depending on the value of $p_f(N)$, when probabilistic flooding is applied to a connected random graph. The observation of such a connection allowed for the a further study of probabilistic flooding with the use of elements from graph theory.

2.1 Analytic bounds on p_f for global outreach of $G(N, p)$

One of the problems that has been studied was to find analytic (asymptotic) bounds for the forwarding probability $p_f(N)$, in order to achieve full coverage of an underlying random graph $G(N, p)$, while producing the minimum number of information dissemination messages. In order to achieve this, the use of two random graphs $G(N, p * p_f)$ and $G(N, p * p')$ is proposed. This two graphs have the same value of p with the one that is used to create the connected random graph $G(N, p)$.

Since a link of the underlying connected random graph $G(N, p)$ will be part of the probabilistic flooding network either with probability $p_f(N)$ (when only one of the end nodes of the link receives the message from an other link and takes a forwarding decision) or with probability $p'(N) = 2p_f(N) - p_f^2(N)$ (when both the end nodes of the link receive the message from a different than their common link and, thus, both make an (assumed independent) decision to forward it over the common link), it is expected that $P(G(N, p), p_f)$ contains on average more links than $G(N, p * p_f)$. Consequently, when $G(N, p * p_f)$ is connected with high probability (i.e., *w.h.p* from now on), then $P(G(N, p), p_f)$ is also connected *w.h.p* and, thus includes all network nodes *w.h.p*. Note also that since $p_f(N) \leq p'(N)$ (the equality holds for $p_f(N) = 1$), $G(N, p * p')$ contains (on average) more links than $G(N, p * p_f)$ and when the latter network is connected the former is also connected *w.h.p*.

Based on the two previous observations, and assuming a certain value for $p(N)$, as $p_f(N)$ increases, it is expected that there will be some probability value for $p_f(N)$ for which $G(N, p * p')$ becomes connected *w.h.p*. As $p_f(N)$ increases further, $P(G(N, p), p_f)$ becomes connected (equivalently, $C(0) = 1$) *w.h.p*. For

further increment, $G(N, p * p_f)$ also becomes connected. Consequently, the particular value of $p_f(N)$ for which probabilistic flooding disseminates information to all network nodes (thus achieve global network outreach) is “between” the values of $p_f(N)$ for which $G(N, p * p')$ and $G(N, p * p_f)$ become connected. This analysis has been also verified by simulation results in a $G(10000, 0.0008)$ random graph. The simulation results have also shown that this behaviour is also present when smaller network coverage cases are studied, something that has not been covered by the analysis.

An interesting result is that even though $G(N, p * p')$ becomes connected for smaller values of $p_f(N)$ when compared to $G(N, p * p_f)$ (as already mentioned $p' > p_f$ for $0 < p_f < 1$), these values have the same asymptotical order.

2.2 Probabilistic Versus Full Flooding

To be able to study cases of smaller coverage (than the global outreach of the underlying connected graph), a new metric was introduced. This metric is the L – coverage: it includes the number of nodes that have been informed along with those nodes that are at most L hops from (at least) one of them. The cases of $L = 0$ (i.e., global network outreach) and $L = 1, 2$ are covered here. Larger values for L were not studied due to the *small world phenomenon*.

The next problem that was studied was a comparison of the probabilistic flooding approach with the traditional flooding. A reduction on the number of messages to be achieved under probabilistic flooding at a cost of an increase in the termination time until global outreach of the underlying random graph was expected and verified both by analytic and by simulation results.

Let $R_{M,L}(N)$ denote the (asymptotic) fraction of messages under probabilistic flooding over those under traditional flooding for some L , or, $R_{M,L}(N) = \frac{\mathcal{E}(\mathbb{P}(\mathbb{G}(N,p):p_f))}{\mathcal{E}(\mathbb{G}(N,p))}$. For the case of $L = 0$,

$$R_{M,0}(N) = \frac{\ln(N)}{p(N)N}. \quad (1)$$

Obviously, $R_{M,0}(N) \rightarrow 0$, when $N \rightarrow +\infty$ ($O(p(N)) > O(p_{Q_0}(N)) = O(\frac{\ln(N)}{N})$). Note that in strict terms, $R_{M,0}(N) = \Theta\left(\frac{\ln(N)}{p(N)N}\right)$.

As already mentioned, since $\mathbb{G}(N, p)$ is a connected network w.h.p., then $O(p(N)) > O(p_{Q_0}(N)) = O(\frac{\ln(N)}{N})$. For $p(N) = \Theta\left(\frac{\ln(N)}{N}\right)$ (which means that $\mathbb{G}(N, p)$ has just become connected), it is interesting to see that $R_{M,0}(N) = \Theta(1)$, which apparently demonstrates the fact that there is no advantage under probabilistic flooding when compared to traditional flooding for this case (the number of messages is the same under both probabilistic flooding and traditional flooding). Actually, this particular case is -asymptotically- equivalent to $p_f(N) = 1$, for which probabilistic flooding reduces to traditional flooding. In order to explain further this observation, note that for the case of $p(N) = \Theta\left(\frac{\ln(N)}{N}\right)$, $\mathbb{G}(N, p)$ has just become connected w.h.p. which apparently means that the

number or “redundant” links (links over which traditional flooding forwards messages and probabilistic flooding “saves” by probabilistically “avoiding” to do so) is significantly reduced. The shape of the network – even though it contains cycles – looks mostly like a tree, and therefore, the ability of probabilistic flooding to “avoid” forwarding messages over “redundant” links is reduced.

The reduced number of messages under probabilistic flooding is achieved at the expense of larger termination delays. This is shown by comparing the network diameter of $\mathbb{G}(N, p)$ and $\mathbb{P}(\mathbb{G}(N, p), p_f)$, for $p_f(N) = \Theta\left(\frac{\ln(N)}{N}\right)$ (the upper bound of termination time corresponds to the network diameter). Let $R_{T,L}(N)$ denote the (asymptotic) fraction of the network diameter of the probabilistic flooding network over the network diameter minus L (for fairness issues), for some $L = 0, 1, 2$, or $R_{T,L}(N) = \frac{\mathcal{D}(\mathbb{P}(\mathbb{G}(N, p), p_f))}{\mathcal{D}(\mathbb{G}(N, p)) - L}$. Given that for $L = 0$, $\mathcal{D}(\mathbb{G}(N, p)) = \Theta\left(\frac{\ln(N)}{\ln(p(N)N)}\right)$, $\mathcal{D}(\mathbb{P}(\mathbb{G}(N, p), p_f)) = \Theta(\mathcal{D}(\mathbb{G}(N, p * p_f))) = \Theta\left(\frac{\ln(N)}{\ln(p(N)p_f(N)N)}\right)$, and $p(N)p_f(N) = \Theta\left(\frac{\ln(N)}{N}\right)$, it follows that,

$$R_{T,0}(N) = \frac{\ln(p(N)N)}{\ln(\ln(N))}. \quad (2)$$

So far, the global network outreach case (i.e., $L = 0$ or $C(0) = 1$) has been studied. The cases corresponding to $L = 1$ and $L = 2$ are naturally expected to yield smaller number of messages under probabilistic flooding – compared to traditional flooding – since the particular values of $p_f(N)$ are expected to be (on average) smaller than those ensuring global network outreach (i.e., $C(0) = 1$) w.h.p. The asymptotic analysis that has been previously followed for $L = 0$ applies to these particular cases as well. Note, that the resemblance is only asymptotic and savings with respect to the number of messages are greater for the case of $L = 1$ and $L = 2$ than for the case of $L = 0$ under probabilistic flooding.

Simulation results regarding the number of messages and termination time for $L = 0, 1$ and 2 have verified the aforementioned analysis.

3 Multiple Random Walker

In [10] it is shown that *multiple random walkers*, starting from the same network node, are capable of accelerating the information dissemination process and reduce termination time by a factor equal to the number of random walkers, for a wide range of topologies. On the other hand, as the number of random walkers increases, the number of messages sent increases proportionally to the number of random walkers. Moreover, since random walkers start from the same network node, it is expected for some initial movements to partially overlap (thus, not improving coverage) due to visits to already visited network nodes (i.e., *revisits*). This motivates the adoption of a *replication* approach -under which replicas of random walkers are probabilistically created after each movement- so as to avoid initiating all of them at the same time as under the multiple case. A simple replication mechanism is proposed here capable of covering larger network

areas than multiple random walkers for the same number of random walkers and allowed number of messages.

The contribution of this thesis is the study of multiple random walkers from a different perspective than the one presented in [10]. The analytical part of the work initially assumes a fully connected network topology which allows for the derivation of an analytical expression that confirms the results presented in [10], allowing also for further understanding of various aspects of information dissemination under multiple random walkers. The analysis continues capturing coverage in less dense topologies and an analytical expression is derived showing how coverage is affected by frequent random walk revisits.

3.1 Multiple Random Walkers

Having started with m random walkers from the same initiator node, each random walker moves to one of its neighbour nodes being selected randomly and independently among the set of neighbour nodes, provided that the chosen node is not the previously visited node unless this is the only neighbour node.

Let us $C_m(t)$ denote *coverage* or the fraction of the network nodes visited by any of the m random walkers at time t . $C_m(t)$ is an increasing function of t taking values between $\frac{1}{N}$ (i.e., the case when only one node is visited) and 1 (i.e., all nodes are visited).

Let us, also, define the *termination time*, denoted by T_m , as the smallest value of t such that $C_m(t) = 1$. Alternatively, it is frequently convenient to consider the *asymptotic termination time* T'_m which is defined as the smallest value of t such that $\lim_{N \rightarrow \infty} C_m(t) = 1$.

Theorem 1. *In a fully connected network topology of N nodes and m random walkers, coverage $C_m(t)$ as a function of time t is given by:*

$$C_m(t) = 1 - e^{-\frac{m}{N}t}. \quad (3)$$

3.2 A Replication Mechanism

In topologies less dense than fully connected ones, it is expected that random walkers originating from a common initiator node to frequently revisit network nodes not only due to the probabilistic nature of the random walk mechanism (as it is the case for a fully connected topology), but also due to the topology characteristics. In such a network it is expected m random walkers to cover an almost overlapping network area (frequent revisits) at the beginning, before moving to distant (and likely not previously visited) areas. Therefore, instead of m distinct movements corresponding to the m random walkers, a macroscopic observer (most likely) would observe a number of distinct movements less than m , increasing (on average) as time increases.

In order to exploit this observation and proceed with a qualitative analysis, let us assume that the underlying topology is a fully connected network (as before), in which random walkers move (and overlap) as it would have been the case

if the underlying topology was not a fully connected one. The fully connected topology assumption is useful in order to simplify the analysis reusing results derived when proving Theorem 1. Let us $f(t)$ denote the (average) fraction of random walkers seen by the macroscopic observer at time t . $mf(t)$ corresponds to the (average) number of distinct movements of random walkers in the network. In general, $f(0)$ is expected to be rather small and $f(t)$ to be close to 1 for large values of t . Let us assume that $f(t) = 1 - e^{-\alpha t}$, where α is a constant that varies depending on the characteristics of each environment (e.g., number of nodes, density, bottleneck links).

Theorem 2. *In a fully connected network topology of N nodes and $mf(t)$ random walkers, coverage as a function of time t is given by:*

$$C_m(t) = 1 - e^{-\frac{m}{N}(t - \frac{1}{\alpha}(1 - e^{-\alpha t}))}. \quad (4)$$

It is interesting to observe that $C_m(t)$ increases as time increases (as expected) but not that quickly. In particular, for small values of t , $C_m(t)$ increases slowly, then it reaches an inflection point at some point $t = t_0$.

The existence of the inflection point (confirmed by simulation results), is the basic motivation behind the introduction of a simple replication mechanism in the sequel. It is evident that due to revisits, a large number of random walkers may not always allow for significant coverage improvement, while at the same time an increased number of network resources are wasted. Under replication, a small number of $m_0 \ll m$ random walkers is initially released at the initiator node and afterwards, more random walkers are created by replicating the existing ones. Special care is taken for the total number of random walkers in the network not to exceed m . Note that the values of m_0 comparable to m eventually do not make any difference with respect to the problem of revisits since they reduce the replication mechanism to the multiple random walkers mechanism.

The replication mechanism: Having started with m_0 random walkers, for each random walker a replica is created with constant probability $\frac{1}{q}$ after each movement. All random walkers move in the network according to the multiple random walkers mechanism.

3.3 Results

The simulation results have verified the performance covered by the analysis for a fully connected topology. For the sparser topologies, a random geometric graph was considered and simulations for various values of the variable r_c (i.e., variable that defines the connectivity of the graph) took place verifying the expected inflection point, in coverage performance. Finally, in order to examine the replication performance, a simple replication policy (consisted of one initial walker but very frequent replications) was considered and the results showed that replication outperforms the multiple random walkers mechanism when the number of simultaneous walkers (i.e., m) is large.

4 Randomly Replicated Random Walks

Taking under consideration the overall good performance of the introduced replication mechanism, when it is compared to the multiple random walkers example, a further study of its efficiency in advertising the disseminated information in various networking topologies has taken place.

The efficiency of an advertising mechanism is measured not only by the number of nodes that are informed about the location of a specific service (e.g., size of the advertising network), but also by the succeeded dissemination of this information over broader areas in the network, bringing it close to as many nodes as possible so as to reduce the intensity of their search for it (i.e., searching will apply an L -controlled flooding mechanism). As a measure of the performance regarding *stretching* the information dissemination in wider areas of the network, the previously introduced metric of L -coverage is considered. In this study, the case for $L = 0$ (i.e., to measure the size of the advertising network), $L = 1, 2$ (i.e., to measure the stretching capabilities) is examined.

To this end, the performance of two widely used techniques (flooding and single random walker) is examined through simulations, along with the performance of a proposed broad class of information dissemination schemes. The introduced Randomly Replicated Random Walks (RRRWs) scheme employs random walkers that *replicate* themselves; the *replication* policy considered here creates replicas according to an exponentially decreasing probability (in contrast with the replication policy that is assumed when comparison with the multiple walkers case took place earlier on), thus creating more replicas at the beginning of the process, controls the number of walkers and increases the probability of *stretching* the information to undiscovered parts of the underlying network. When the first replication probability equals one, the number of the walkers that are used for advertising the location of the service is large and the approach closely resembles to flooding. When the first replication probability equals zero, then only one walker is used and the approach closely resembles to SRW.

Extensive simulation results over several widely employed network topologies reveal that the RRRW scheme outperforms the single random walker (although the difference is small in the random geometric graphs), while the comparison to flooding depends on the topology. The RRRW scheme outperforms flooding in the random geometric graphs and in the clustered topologies, with respect to the size of the generated advertising network, while they manage to stretch more the information dissemination in the power-law topologies and the aforementioned clustered and random geometric graph environments, even though the generated advertising networks in the former topology are smaller in size than the ones generated by flooding. From the above, it can be stated that the RRRW scheme performs better in topologies that manage to capture best the sense of geographical coverage of a network (e.g., random geometric and clustered environments).

5 Conclusions

The study of several information dissemination techniques has taken place in this thesis. The reason for this study is an effort to examine more efficiently (both regarding the number of generated messages and the termination time) the process of Service Discovery in modern unstructured network environments.

For probabilistic flooding, analytic asymptotic bounds were given for the forwarding probability p_f in order to achieve full coverage of an underlying connected random graph $(G(N, p))$. On top of it, an analysis comparing the overhead induced both by probabilistic and by full flooding was conducted, for various coverage performance of the underlying graph $G(N, p)$. The simulation results were in accordance with the analysis.

For multiple random walkers, analytical expressions for coverage and termination time when a fully connected topology is presented. The study of less dense topologies revealed that the coverage performance in the very early stages is not so effective due to the large number of walkers that is collected near the initiator node. Later this observation was confirmed by simulation results.

A replication mechanism was introduced, mainly as a solution to the aforementioned ineffective coverage performance. This mechanism is examined as an effective advertising approach, studying whether it can stretch the dissemination of information in broader networking areas and fill the performance gap between the full flooding and single random walker approach.

References

1. A. Segall, "Distributed network protocols", IEEE Trans. Inform. Theory, vol. IT-29, Jan. 1983.
2. Gnutella RFC, <http://rfc-gnutella.sourceforge.net/>, 2002.
3. C. Lv, P. Cao, E. Cohen, K. Li, and S. Shenker, "Search and Replication in Unstructured Peer-to-Peer Networks", ICS 2002, 2002.
4. C. Gkantsidis, M. Mihail and A. Saberi, "Hybrid Search Schemes for Unstructured Peer-to-Peer Networks", IEEE Infocom 2005, 2005.
5. N. B. Chang and M. Liu, "Optimal Controlled Flooding Search in a Large Wireless Network", Third International Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOpt'05), 2005, pp. 229 - 237.
6. V. Kalogeraki, D. Gunopoulos and D. Zeinalipour-Yazti, "A Local Search Mechanism for Peer-to-Peer Networks", in CIKM (International Conference on Information and Knowledge Management), 2002.
7. F. Banaei-Kashani and C. Shahabi, "Criticality-based Analysis and Design of Unstructured Peer-to-Peer Network as "Complex Systems", in Proceedings of the Third International Symposium on Cluster Computing and the Grid, 2003, pp. 51 - 358.
8. D. Tsoumakos and N. Roussopoulos, "Adaptive Probabilistic Search for Peer-to-Peer Networks," 3rd IEEE International Conference on P2P Computing, 2003.
9. B. Bollobás, "Random Graphs", Cambridge University Press, Second Edition, 1999.
10. N. Alon, C. Avin, M. Koucky, G. Kozma, Z. Lotker and M. R. Tuttle. "Many random walkers are faster than one", in SPAA '08: Proceedings of the twentieth annual symposium on Parallelism in algorithms and architectures, pages 119-128, New York, NY, USA, 2008, ACM.

Algebraic algorithms for polynomial system solving and applications

Christos Konaxis*

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
ckonaxis@di.uoa.gr

Abstract. We consider sparse elimination theory in order to describe the Newton polytope of the sparse resultant of a given overconstrained algebraic system, by enumerating equivalence classes of mixed subdivisions. In particular, we consider specializations of this resultant to a polynomial in a constant number of variables, typically up to 3. We sketch an algorithm that avoids computing the entire secondary polytope; our goal is that it examines only the silhouette of this polytope with respect to an orthogonal projection. Moreover, since determinantal formulae are not always possible, the most efficient general method for computing resultants is by rational formulae. We propose a single lifting function which yields a simple method for computing Macaulay-type formulae of sparse resultants, in the case of generalized unmixed systems, where all Newton polytopes are scaled copies of each other. As another application of sparse elimination, we consider rationally parameterized plane curves and determine the vertex representation of the implicit equation's Newton polygon.

1 Introduction

In this dissertation we study problems in sparse elimination: rational formulae for sparse resultants via a single lifting function, computation of the Newton polytope of the sparse resultant and of some of its interesting specializations, and sparse implicitization of rational parametric plane curves. The sparse (or toric) resultant captures the structure of the polynomials by combinatorial means and constitutes the cornerstone of sparse elimination theory [5, chap.7],[19,29]. It is an important tool in deriving new, tighter complexity bounds for system solving, Hilbert's Nullstellensatz, and related problems. These bounds depend on the polynomials' Newton polytopes and their mixed volumes, instead of total degree, which is the only parameter in classical elimination theory. In particular, if d bounds the total degree of each polynomial, the projective resultant has complexity roughly $d^{O(n)}$, whereas the sparse resultant is computed in time roughly proportional to the number of integer lattice points in the Minkowski sum of the Newton polytopes.

* Dissertation Advisor: Ioannis Z. Emiris, Professor

The sparse resultant is defined for an overconstrained system of $n+1$ Laurent polynomials $f_i \in K[x_1^{\pm 1}, \dots, x_n^{\pm 1}]$, in n variables over some coefficient ring K . It is the unique, up to sign, integer polynomial over K which vanishes precisely when the system has a root in the toric projective variety X defined by the supports of f_i , in which the torus $(\bar{K})^n$ is a dense subset.

2 Preliminaries

We now recall some crucial notions of sparse elimination theory. Given a polynomial f , its *support* $\mathcal{A}(f)$ is the set of the exponent vectors corresponding to monomials with nonzero coefficients. Its *Newton polytope* $\mathcal{N}(f)$ is the convex hull of $\mathcal{A}(f)$, denoted $\text{CH}(\mathcal{A}(f))$. Newton polytopes are the main tool that allows us to translate algebraic problems into the language of combinatorial geometry. The *Minkowski sum* $A+B$ of $A, B \subset \mathbb{R}^n$ is the set $A+B = \{a+b \mid a \in A, b \in B\} \subset \mathbb{R}^n$. If A, B are convex polytopes, then $A+B$ is also a convex polytope. In what follows we will also denote the support of a polynomial f_i as A_i and its Newton polytope as Q_i .

Let Q_0, \dots, Q_n be polytopes in \mathbb{R}^n with $P_i = \text{CH}(A_i)$ and Q their Minkowski sum. We assume that Q is n -dimensional. A *Minkowski cell* of Q is any full-dimensional convex polytope $B = \sum_{i=0}^n B_i$, where each B_i is a convex polytope with vertices in A_i . We say that two Minkowski cells $B = \sum_{i=0}^n B_i$ and $B' = \sum_{i=0}^n B'_i$ *intersect properly* when the intersection of the polytopes B_i and B'_i is a face of both and their Minkowski sum descriptions are compatible.

Definition 1. A *mixed subdivision* of Q is any family S of Minkowski cells which partition Q and intersect properly as Minkowski sums. A cell R is *mixed*, in particular *i-mixed* or *v_i-mixed*, if it is the Minkowski sum of n 1-dimensional segments $E_j \subset Q_j$ and one vertex $v_i \in Q_i$: $R = E_0 + \dots + v_i + \dots + E_n$.

A mixed subdivision is called *regular* if it is obtained as the projection of the lower hull of the Minkowski sum of lifted polytopes $\widehat{Q}_i := \{(p_i, \omega_i(p_i)) \mid p_i \in Q_i\}$. If the lifting function $\omega := \{\omega_0, \dots, \omega_n\}$ is sufficiently generic, then the induced mixed subdivision is called *fine* or *tight*, and $\sum_{i=0}^n \dim B_i = \dim \sum_{i=0}^n B_i$, for every cell $\sum_{i=0}^n B_i$. This construction method ensures that the lower hull facets of the Minkowski sum of the lifted polytopes \widehat{Q}_i , are projected bijectively onto Q . Thus, every cell R of the mixed subdivision can be written uniquely as the Minkowski sum $R = F_0 + \dots + F_n \subset \mathbb{R}^n$, where each F_i is a face of Q_i . Two mixed subdivisions are equivalent if they share the same mixed cells. The equivalence classes are called *mixed cell configurations* [25].

A monomial of the sparse resultant is called *extreme* if its exponent vector corresponds to a vertex of the Newton polytope $\mathcal{N}(\mathcal{R})$ of the resultant. The following corollary of [28, Thm. 2.1], allows us to compute the extreme monomials of the sparse resultant using tight regular mixed subdivisions.

Corollary 1. *There exists a surjection from the mixed cell configurations onto the set of extreme monomials of the sparse resultant.*

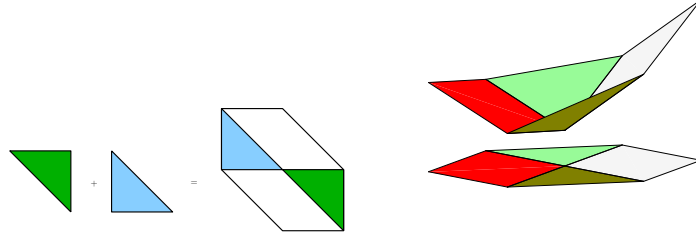


Fig. 1. Minkowski sum of two triangles (left) and construction of its regular mixed subdivision (right).

Given supports A_0, \dots, A_n , the Cayley embedding κ introduces a new point set $C := \kappa(A_0, A_1, \dots, A_n) = \bigcup_{i=0}^n (A_i \times \{e_i\}) \subset \mathbb{R}^{2n}$, where e_i are an affine basis of \mathbb{R}^n . The following proposition reduces the computation of regular tight mixed subdivisions to the computation of regular triangulations.

Proposition 1. [The Cayley Trick]. *There exists a bijection between the regular tight mixed subdivisions of the Minkowski sum P and the regular triangulations of C .*

Regular triangulations are in bijection with the vertices of the *secondary polytope* [19]. A *bistellar flip* is a local modification on a triangulation that leads to a new one. The following theorem allows us to explore the set of regular triangulations of a point set using bistellar flips.

Theorem 1. [19] *For every set C of points affinely spanning \mathbb{R}^d there is a polytope $\Sigma(C)$ in $\mathbb{R}^{|C|-d-1}$, the secondary polytope of C , such that its vertices correspond to the regular triangulations of C and there is an edge between two vertices if and only if the two corresponding triangulations are obtained one from the other by a bistellar flip.*

3 Basic results

3.1 Macaulay-type formulae for generalized unmixed sparse resultants

A resultant is most efficiently expressed by a *matrix formula*: this is a generically nonsingular matrix, whose specialized determinant is a multiple of the resultant. Its degree in the coefficients of one polynomial equals the corresponding degree of the resultant. For $n = 1$ there are matrix formulae named after Sylvester and Bézout, whose determinant equals the resultant. Unfortunately, such determinantal formulae do not generally exist for $n > 1$, except for specific cases, e.g. [7,9,22]. Macaulay's seminal result [24] expresses the extraneous factor as a minor

of the matrix formula, for projective resultants of (dense) homogeneous systems, thus yielding the most efficient general method for computing such resultants.

Matrix formulae for the sparse resultant were first constructed in [1]. The construction relies on a lifting of the given polynomial supports, which defines a mixed subdivision of their Minkowski sum into mixed and non-mixed cells, then applies a perturbation δ so as to define the integer points that index the matrix. The algorithm was extended in [3,2,28]. In the case of dense systems, the matrix coincides with Macaulay's numerator matrix.

Extending the Macaulay formula to toric resultants had been conjectured in [3,5,10,19,28]; it was a major open problem in elimination theory. D'Andrea's result [6] answers the conjecture by a *recursive* definition of a Macaulay-type formula. But this approach does not offer a global lifting, in order to address the stronger original conjecture [10, Conj. 3.1.19], [3, Conj. 13.1].

We give an affirmative answer to this stronger conjecture by presenting a single lifting which constructs Macaulay-type formulae for generalized unmixed systems, i.e. when all Newton polytopes are scaled copies of each other. We state our main result:

Theorem 2. [12] *The single lifting algorithm of Section 3.1 constructs a Macaulay-type formula for the toric resultant of an overconstrained generalized unmixed algebraic system, by means of the lifting function of Definition 5.*

Our method can be generalized to certain mixed systems: those with $n \leq 3$, as well as systems that possess sufficiently different Newton polytopes. A single lifting algorithm is conceptually simpler and also easier to implement.

D'Andrea's [6] recursive construction requires one to associate integer points with cells of every dimension from n to 1. Our method constructs the matrix formula directly, without recursion, by examining only n -dimensional cells. These are more numerous than the n -dimensional cells in [6] but our algorithm defines significantly fewer cells totally. The weakness of our method is to consider extra points besides the input supports. Related implementations have been undertaken in Maple, but cover only the original Canny-Emiris method [3], either standalone¹ or as part of library Multires². We expect that our algorithm shall lead to an efficient implementation of Macaulay-type formulae.

Let f_0, \dots, f_n be polynomials with supports $A_0, \dots, A_n \subset \mathbb{Z}^n$ and Newton polytopes

$$Q_0, \dots, Q_n \subset \mathbb{R}^n, \quad Q_i = \text{CH}(A_i),$$

where $\text{CH}(\cdot)$ denotes convex hull. A monomial with exponent $a = (a_1, \dots, a_n) \in \mathbb{Z}^n$ shall be denoted as x^a , where $x := x_1 \cdots x_n$.

Our lifting shall induce a regular and fine (or tight) mixed subdivision of the Minkowski sum $\sum_{i=0}^n Q_i$. Let Z be the integer lattice generated by $\sum_{i=0}^n A_i$. The Minkowski sum $\sum_{i=0}^n Q_i$ is perturbed by a vector $\delta \in \mathbb{Q}^n$, which is sufficiently small with respect to Z , and in sufficiently generic position with respect to the

¹ http://www.di.uoa.gr/~emiris/soft_alg.html

² <http://www-sop.inria.fr/galaad/logiciels/multires.html>

Q_i . The lattice points in $\mathcal{E} = Z \cap (\sum_{i=0}^n Q_i + \delta)$ are associated to a unique maximal cell of the subdivision, and this allows us to construct a matrix formula M whose rows and columns are indexed by these points. In particular, polynomial $x^{p-a_{ij}} f_i$ fills in the row indexed by the lattice point p in Definition 2.

Definition 2. Let $p \in \mathcal{E}$ lie in a cell $F_0 + \cdots + F_n + \delta$ of the perturbed mixed subdivision, where F_i is a face of Q_i . The row content (RC) of p is (i, j) , if $i \in \{0, \dots, n\}$ is the largest integer such that F_i equals a vertex $a_{ij} \in A_i$.

Our method is based on the matrix construction algorithm of [3,10]. For completeness, we recall the basic steps:

1. Pick (affine) liftings $H_i : \mathbb{Z}^n \rightarrow \mathbb{R} : A_i \rightarrow \mathbb{Q}, i = 0, \dots, n$.
2. Construct a regular fine mixed subdivision of the Minkowski sum $\sum_{i=0}^n Q_i$ using liftings H_i .
3. Perturb the Minkowski sum $\sum_{i=0}^n Q_i$ by a sufficiently small vector $\delta \in \mathbb{Q}^n$, so that integer points in $\sum_{i=0}^n Q_i + \delta$ belong to a unique cell of the subdivision, and assign *row content* to these points by Definition 2.
4. Construct resultant matrix M with rows and columns indexed by the previous integer points.

The main idea of both our and D'Andrea's methods is that one point, say $b_{01} \in Q_0$, is lifted significantly higher. Then, the 0-summand of all maximal cells is either b_{01} or a face not containing it. In D'Andrea's case, facets not containing b_{01} correspond to different subsystems where the algorithm recurses (each time on the integer lattice specified by that subsystem). In designing a unique lifting, the issue is that points appearing in two of these subsystems may be lifted differently in different recursions. To overcome this, we introduce several points c_{ijs} , each lying in a suitable face of Q_i indexed by s , very close (with respect to Z) to every b_{ij} , which is lifted very high at recursion i by D'Andrea's method. This captures the multiple roles b_{ij} may assume in every recursion step.

Single lifting Algorithm. Our algorithm directly generalizes the one given in [3,10], and is based on the 4 steps described above. We modify step (1) and define a new lifting function; moreover, we describe necessary adjustments to the matrix construction and extend step (4) so as to produce the denominator matrix of the Macaulay-type formula. The following three definitions suffice to specify our algorithm.

We shall use \mathcal{E} to index the rows (and columns) of the numerator matrix M , whereas the denominator shall be indexed by points lying in non-mixed cells. We focus on generalized unmixed systems, where

$$Q_i = k_i Q \subset \mathbb{R}^n,$$

for some n -dimensional lattice polytope Q and $k_i \in \mathbb{N}^*, i = 0, \dots, n$. Let the vertices of Q be $b_0, \dots, b_{|A|}$, where $Q = \text{CH}(A)$. We shall denote the vertices of each $Q_i = k_i Q$, for $i = 0, \dots, n$, as $b_{i1}, \dots, b_{i|A|}$. Obviously, $b_{ij} := k_i b_j$.

Definition 3. For $i = 0, \dots, n-2$, consider any $(n-i)$ -dimensional face $F_s^{(i)} \subset Q$, where integer s indexes all such faces. Take any vertex $b_{ij} \in k_i F_s^{(i)}$, for any valid $j \in \mathbb{N}$. Let $\delta_{ijs} \in \mathbb{Q}^n$ denote a perturbation vector such that:

1. $b_{ij} + \delta_{ijs}$ lie in the relative interior of $k_i F_s^{(i)}$,
2. It is sufficiently small compared to lattice Z , and $\|\delta_{ijs}\| \ll \|\delta\|$, where $\|\cdot\|$ is the Euclidean norm and δ as above, and
3. It is sufficiently generic to avoid all edges in the mixed subdivision of $\sum_{i=0}^n Q_i$.

For an example of Definition 3 see Figure 2, where the (appropriately translated) δ_{ijs} 's are depicted by arrows. We shall use the perturbation vectors of Definition 3 to define extra points *not* contained in the input supports. Condition (2) of Definition 3 implies that, in the mixed subdivision induced by the single lifting function β below, the cells created by the introduction of the extra points will not contain integer points after we perturb the mixed subdivision by δ . This can be checked at the end of the construction of the mixed subdivision.

Definition 4. We define points $c_{ijs} \in Q_i \cap \mathbb{Q}^n$, for $i = 0, \dots, n-2$. Firstly, set $c_{011} := b_{01} + \delta_{011} \in Q_0 \cap \mathbb{Q}^n$ where δ_{011} satisfies Definition 3. Now let $\{c_{ijs} \in k_i F_s^{(i)}\}$ be the set of points defined in Q_i , where s ranges over all $(n-i)$ -dimensional faces $F_s^{(i)} \subset Q$ and j over the set of indices of points in Q_i . Then, let $F_u^{(i+1)}$ be a facet of $F_s^{(i)}$ such that:

1. $k_i F_u^{(i+1)}$ does not contain any of the b_{ij} 's corresponding to the already defined c_{ijs} 's, and
2. $k_{i+1} F_u^{(i+1)}$ does not contain any of the already defined $c_{(i+1)l}$'s.

For each such facet choose a vertex $b_{(i+1)j} \in A_{i+1}$, for some j , and a suitable perturbation vector $\delta_{(i+1)ju}$ satisfying Definition 3, and set $c_{(i+1)ju} := b_{(i+1)j} + \delta_{(i+1)ju} \in Q_{i+1} \cap \mathbb{Q}^n$.

The previous definition implies a many-to-one mapping from the set of c_{ijs} 's to that of b_{ij} 's; it reduces to a bijection when restricted to a fixed face $k_i F_s^{(i)} \subset Q_i$ containing b_{ij} . For an application of Definition 4 for $n = 2$ see Figure 2 where Q is the unit square. In this example, for illustration purposes, we define points c_{ijs} also on edges of polytope Q_1 .

Definition 5. Let $h_0 \gg h_1 \gg \dots \gg h_{n-1} \gg 1$. The single lifting algorithm uses sufficiently random linear functions $H_i, i = 0, \dots, n$, such that:

$$1 \gg H_i(a_{ij}) > 0, \text{ and } H_i \gg H_t, \ i < t,$$

where $a_{ij} \in A_i$ and $i, t = 0, \dots, n$, $j = 1, \dots, |A_i|$. Define a global lifting β as follows:

1. $c_{ijs} \mapsto h_i$, $c_{ijs} \in k_i F_s^{(i)} \subset Q_i$, $i = 0, \dots, n-1$; this is called primary lifting.
2. $a_{ij} \mapsto H_i(a_{ij})$, $a_{ij} \in A_i$, $i = 0, \dots, n$.

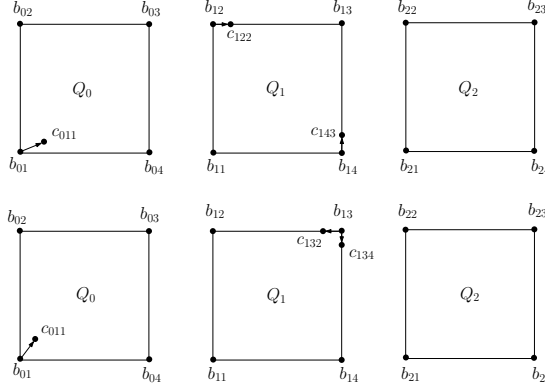


Fig. 2. Two scenarios of an application of Def. 4 for 3 unit squares. Facets are numbered clockwise starting from the left vertical edge

Let F^β denote face F lifted under β . Now c_{tjs}^β , for all valid j, s , is much higher, respectively lower, than any c_{ijs}^β , for $i > t$, respectively $i < t$. The β -induced subdivision contains edges with one or two vertices among the c_{ijs} , and edges from the Q_i . The vertex set of the upper hull of Q_i^β contains some or all of the c_{ijs}^β and the lifted vertices of Q_i .

Figure 3 shows the mixed subdivisions of three unit squares and their Minkowski sum, induced by lifting β . Here, the perturbation vectors are not sufficiently small compared to \mathbb{Z}^2 for illustration purposes.

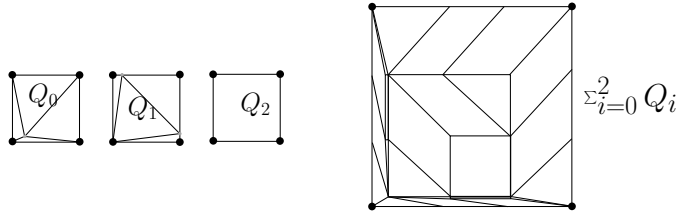


Fig. 3. The mixed subdivisions of 3 unit squares and their Minkowski sum induced by lifting β

The matrix formula M constructed by our algorithm is indexed by all lattice points in \mathcal{E} . To decide the content of each row, every point is associated to a unique (maximal) cell of the mixed subdivision according to Definition 2. The t -mixed cells contain lattice points as follows:

$$p \in k_0 E_0 + \cdots + k_{t-1} E_{t-1} + c_{tjs} + k_{t+1} E_{t+1} + \cdots + k_n E_n \cap Z,$$

for edges $E_i \subset Q$ spanning \mathbb{R}^n . This gives unique writing

$$p = p_0 + \cdots + p_{t-1} + (b_{tj} + \delta_{tjs}) + p_{t+1} + \cdots + p_n, \quad p_i \in A_i \cap E_i.$$

Hence, the row indexed by p , as with matrix constructions in [3,6], contains a multiple of $f_t(x)$:

$$x^{p_0 + \cdots + p_{t-1} + p_{t+1} + \cdots + p_n} f_t(x),$$

and the diagonal element is the coefficient of the monomial with exponent b_{tj} in $f_t(x)$. Similarly, for the rows corresponding to lattice points in non-mixed cells. The extraneous factor $\det M / \text{Res}(f_0, \dots, f_n)$ is the minor of M indexed by points in \mathcal{E} lying in non-mixed cells.

3.2 The Newton polygon of rational parametric plane curves

Implicitization is the problem of switching from a parametric representation of a hypersurface to an algebraic one. It is a fundamental question with several applications, see [20]. We consider the implicitization problem for a planar curve, where the polynomials in its parameterization have fixed Newton polytopes. We determine the vertices of the Newton polygon of the implicit equation, or *implicit polygon*, without computing the equation, under the assumption of *generic* coefficients relative to the given supports, i.e. our results hold for all coefficient vectors in some open dense subset of the coefficient space. The support of the implicit equation, or *implicit support*, is taken to be all interior points inside the implicit polygon.

This problem was posed in [32] but has received much attention lately. According to [30], “apriori knowledge of the Newton polytope would greatly facilitate the subsequent computation of recovering the coefficients of the implicit equation [...] This is a problem of numerical linear algebra ...”.

Previous work includes [15,16], where an algorithm constructs the Newton polytope of any implicit equation. That method had to compute all mixed subdivisions, then applies Cor. 1. In [19, chapter 12], the authors study the resultant of two univariate polynomials and describe the facets of its Newton polytope. In [18], the extreme monomials of the Sylvester resultant are described. The approaches in [15,19] cannot exploit the fact that the denominators in a rational parameterization may be identical.

By employing tropical geometry, [30,31] compute the implicit polytope for any hypersurface parameterized by Laurent polynomials. Their theory extends to arbitrary implicit ideals. They give a generically optimal implicit support; for curves, the support is described in [30, example 1.1].

More recently, in [17] the problem was solved in an abstract way by means of composite bodies and mixed fiber polytopes. In [8] the normal fan of the implicit polygon is determined. This is computed by the multiplicities of any parameterization of the rational plane curve. The authors reduce the problem to studying the support function of the implicit polytope and counting the number of solutions of a certain system of equations. The latter is solved by applying a refinement of the Kushnirenko-Bernstein formula for the computation of the

isolated roots of a polynomial system in the torus, given in [26]. As a corollary, they obtain the optimal implicit polygon in the case of generic coefficients.

In [13], we presented a method to compute the vertices of the implicit polygon of polynomial or rational parametric curves, when the denominators differ. We also introduced a method and gave partial results for the case when denominators are equal; both methods are described in final form in [14].

Our main contribution is to determine the vertex structure of the implicit polygon of a rational parameterized planar curve, or implicit vertices, under the assumption of *generic* coefficients. If the coefficients are not sufficiently generic, then the computed polygon contains the implicit polygon. Our approach considers the symbolic resultant which eliminates the parameters and, then, is specialized to yield an equation in the implicit variables. In the case of rationally parameterized curves with different denominators (which includes the case of Laurent polynomial parameterizations), the Cayley trick reduces the problem to computing regular triangulations of point sets in the plane. If the denominators are identical, two-dimensional mixed subdivisions are examined; we show that only subdivisions obtained by *linear* liftings are relevant. These results also apply if the two parametric expressions share the same numerator, or the numerator of one equals the denominator of the other. We prove that, in these cases, only extremal terms matter in determining the implicit polygon as well as in ensuring the genericity hypothesis on the coefficients.

The following proposition collects our main corollaries regarding the shape of the implicit polygon in terms of corner cuts on an initial polygon. A corner cut on a polygon P is a line that intersects the polygon, excluding one vertex while leaving the rest intact. ϕ is the implicit equation and $N(\phi)$ is the implicit polygon.

Proposition 2. *$N(\phi)$ is a polygon with one vertex at the origin and two edges lying on the axes. In particular, for polynomial parameterizations, $N(\phi)$ is a right triangle with at most one corner cut, which excludes the origin. For rational parameterizations with equal denominators, $N(\phi)$ is a right triangle with at most two cuts, on the same or different corners. For rational parameterizations with different denominators, $N(\phi)$ is a quadrilateral with at most two cuts, on the same or different corners.*

Example 1. Consider the plane curve parameterized by:

$$x = \frac{t^6 + 2t^2}{t^7 + 1}, y = \frac{t^4 - t^3}{t^7 + 1},$$

Our formulas yield vertices $(7, 0), (0, 7), (0, 3), (3, 1), (6, 0)$, which define the actual implicit polygon (see Figure 4, left). Changing the coefficient of t^2 to -1, leads to an implicit polygon with four cuts which is contained in the polygon predicted by our results. This shows the importance of the genericity condition on the coefficients of the parametric polynomials.

An instance where the implicit polygon has 6 vertices is:

$$x = \frac{t^3 + 2t^2 + t}{t^2 + 3t - 2}, y = \frac{t^3 - t^2}{t - 2}.$$

Our results yield implicit vertices $(0, 1)$, $(0, 3)$, $(3, 0)$, $(1, 3)$, $(2, 0)$, $(3, 2)$ which define the actual implicit polygon (see Figure 4, right).

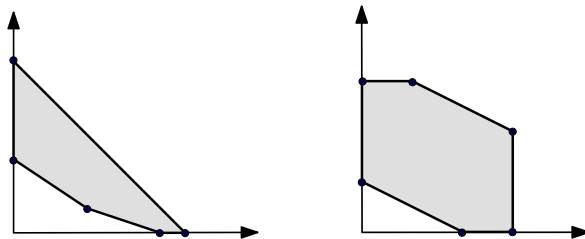


Fig. 4. The implicit polygons of the curves of Example 1

3.3 The Newton polytope of the resultant and its specializations

We describe algorithms to compute the Newton polytope of the sparse resultant, or resultant polytope, of an overconstrained system of polynomials. We rely on Corollary 1 and 1 and following [25] it suffices to enumerate a subset of the vertices of the secondary polytope associated with the input data, corresponding to mixed cell configurations. The resultant polytope allows us to compute a superset of the support of the resultant by considering all integer points contained in it; then we can reduce the computation of the resultant to linear algebra [4].

Corollary 1 establishes a surjection from the set of mixed cell configurations onto the set of vertices of the resultant polytope. Experiments³ indicate that mixed cell configurations are, depending on the input, much less numerous than mixed subdivisions, hence the computation of the resultant vertices becomes more efficient if we focus on the former.

The set of mixed cell configurations corresponds bijectively by the Cayley trick 1, to a set of equivalence classes of regular triangulations. This set can be regarded as a subset of the vertices of the secondary polytope. Thus, we can enumerate mixed cell configurations by enumerating this subset of triangulations. Several algorithms and implementations enumerate regular triangulations e.g. PUNTOS [23], TOPCOM [27], and the algorithm in [21]. We characterize the edges of the secondary polytope that connect the equivalence classes. The sub-graph of the secondary polytope with vertices, regular triangulations corresponding to mixed cell configurations, and the previous edges, is connected.

In [13], we computed the Newton polytope of specialized resultants while avoiding to compute the entire secondary polytope; our approach was to examine the silhouette of the latter with respect to an orthogonal projection. This method is revisited in [11] by studying output-sensitive methods to compute the

³ See for example the webpage <http://ergawiki.di.uoa.gr/index.php/Implicitization>

resultant polytope. Applications such as the computation of the u -resultant or implicitization of polynomial parametric curves or surfaces call for the computation of the resultant polytope after a specialization of some of its indeterminates, i.e. some of the coefficients of the input polynomials. This reduces to enumerating the vertices lying on the silhouette of the secondary polytope $\Sigma(C)$ with respect to some suitably defined projection. For example, the projection of $\Sigma(C)$ to \mathbb{R}^2 solves the problem of implicitization of polynomial curves, the projection to \mathbb{R}^3 the one of polynomial surfaces etc. The silhouette can be obtained naively by computing all the vertices of $\Sigma(C)$, then projecting them to the subspace of smaller dimension. For efficiency we want to enumerate only the vertices lying on a silhouette of $\Sigma(C)$ with respect to a projection to be defined by the problem, without computing $\Sigma(C)$.

In short, we have the following polytope theory problem: We have a high dimensional polytope $\Sigma(C)$ which we know only locally. By this we mean that from every vertex we have an oracle to find the coordinates of all of its neighbours. We describe an algorithm to compute, for a certain projection π to a space of 1, 2, or 3 dimensions, the projection $\pi(\Sigma(C))$.

References

1. J. Canny and I. Emiris. An efficient algorithm for the sparse mixed resultant. In G. Cohen, T. Mora, and O. Moreno, editors, *Proc. Intern. Symp. on Applied Algebra, Algebraic Algor. and Error-Corr. Codes (Puerto Rico)*, number 263 in Lect. Notes in Comp. Science, pages 89–104, Berlin, 1993. Springer-Verlag.
2. J. Canny and P. Pedersen. An algorithm for the Newton resultant. Technical Report 1394, Comp. Science Dept., Cornell University, 1993.
3. J.F. Canny and I.Z. Emiris. A subdivision-based algorithm for the sparse resultant. *J. ACM*, 47(3):417–451, May 2000.
4. R.M. Corless, M.W. Giesbrecht, I.S. Kotsireas, and S.M. Watt. Numerical implicitization of parametric hypersurfaces with linear algebra. In *Artificial intelligence and symbolic computation, (Madrid, 2000)*, pages 174–183. Springer, Berlin, 2001.
5. D. Cox, J. Little, and D. O’Shea. *Using Algebraic Geometry*. Number 185 in GTM. Springer, New York, 2nd edition, 2005.
6. C. D’Andrea. Macaulay-style formulas for the sparse resultant. *Trans. of the AMS*, 354:2595–2629, 2002.
7. C. D’Andrea and A. Dickenstein. Explicit formulas for the multivariate resultant. *J. Pure Appl. Algebra*, 164(1-2):59–86, 2001.
8. C. D’Andrea and M. Sombra. The Newton polygon of a rational plane curve. *Mathematics in Computer Science - MCS Special Issue on Computational Geometry and Computer-Aided Design*, 4(1):3–24, 2010.
9. A. Dickenstein and I.Z. Emiris. Multihomogeneous resultant formulae by means of complexes. *J. Symbolic Computation*, 36(3-4):317–342, 2003. Special issue on ISSAC 2002.
10. I.Z. Emiris. *Sparse Elimination and Applications in Kinematics*. PhD thesis, Computer Science Division, Univ. of California at Berkeley, December 1994.
11. I.Z. Emiris, V. Fisikopoulos, and C. Konaxis. Regular triangulations and resultant polytopes. In *Proc. 26th European Workshop on Computational Geometry*, pages 137–140, Dortmund, Germany, 2010.

12. I.Z. Emiris and C. Konaxis. Single-lifting macaulay-type formulae of generalized unmixed sparse resultants. *J. Symbolic Computation*, Elsevier, 2010. Submitted.
13. I.Z. Emiris, C. Konaxis, and L. Palios. Computing the Newton polytope of specialized resultants. In *Proc. Intern. Conf. MEGA (Effective Methods in Algebraic Geometry)*, 2007. www.ricam.oeaw.ac.at/mega2007/electronic/45.pdf.
14. I.Z. Emiris, C. Konaxis, and L. Palios. Computing the Newton polygon of the implicit equation. *Mathematics in Computer Science - MCS Special Issue on Computational Geometry and Computer-Aided Design*, 4(1):25, 2010.
15. I.Z. Emiris and I.S. Kotsireas. Implicitization with polynomial support optimized for sparseness. In V. Kumar et al., editor, *Proc. Intern. Conf. Comput. Science & Appl. 2003, Montreal, Canada (Intern. Workshop Computer Graphics & Geom. Modeling)*, volume 2669 of *LNCS*, pages 397–406. Springer, 2003.
16. I.Z. Emiris and I.S. Kotsireas. Implicitization exploiting sparseness. In R. Janardan, M. Smid, and D. Dutta, editors, *Geometric and Algorithmic Aspects of Computer-Aided Design and Manufacturing*, volume 67 of *DIMACS*, pages 281–298. AMS/DIMACS, 2005.
17. A. Esterov and A. Khovanskii. Elimination theory and Newton polytopes, 2007. arXiv.org/math/0611107v2.
18. I.M. Gelfand, M.M. Kapranov, and A.V. Zelevinsky. Discriminants of polynomials in several variables and triangulations of Newton polytopes. *Leningrad Math. J.*, 2(3):449–505, 1991. (Translated from *Algebra i Analiz* 2, 1990, pp. 1–62).
19. I.M. Gelfand, M.M. Kapranov, and A.V. Zelevinsky. *Discriminants, Resultants and Multidimensional Determinants*. Birkhäuser, Boston, 1994.
20. C.M. Hoffmann. *Geometric and Solid Modeling*. Morgan Kaufmann, 1989.
21. H. Imai, T. Masada, F. Takeuchi, and K. Imai. Enumerating triangulations in general dimensions. *Intern. J. Comput. Geom. Appl.*, 12(6):455–480, 2002.
22. A. Khetan. The resultant of an unmixed bivariate system. *J. Symbolic Computation*, 36:425–442, 2003.
23. J.De Loera. PUNTOS, 1994. http://www.math.ucdavis.edu/~deloera/RECENT_WORK/puntos2000.
24. F.S. Macaulay. Some formulae in elimination. *Proc. London Math. Soc.*, 1(33):3–27, 1902.
25. T. Michiels and J. Verschelde. Enumerating regular mixed-cell configurations. *Discr. Comput. Geometry*, 21(4):569–579, 1999.
26. P. Philippon and M. Sombra. A refinement of the Bernstein-Kushnirenko estimate. *Advances in Mathematics*, 218:1370–1418, 2008.
27. Jörg Rambau. TOPCOM: Triangulations of point configurations and oriented matroids. In Arjeh M. Cohen, Xiao-Shan Gao, and Nobuki Takayama, editors, *Mathematical Software—ICMS 2002*, pages 330–340. World Scientific, 2002.
28. B. Sturmfels. On the Newton polytope of the resultant. *J. Algebraic Combin.*, 3:207–236, 1994.
29. B. Sturmfels. *Solving Systems of Polynomial Equations*. Number 97 in CBMS Regional Conference Series in Math. AMS, Providence, RI, 2002.
30. B. Sturmfels, J. Tevelev, and J. Yu. The Newton polytope of the implicit equation. *Moscow Math. J.*, 7(2), 2007.
31. B. Sturmfels and J. Yu. Tropical implicitization and mixed fiber polytopes. In *Software for Algebraic Geometry*, volume 148 of *IMA Volumes in Math. & its Applic.*, pages 111–131. Springer, New York, 2008.
32. B. Sturmfels and J.T. Yu. Minimal polynomials and sparse resultants. In F. Orecchia and L. Chiantini, editors, *Proc. Zero-Dimensional Schemes (Ravello, June 1992)*, pages 317–324. De Gruyter, 1994.

A generic product ontology based on software agents incorporating negotiation and decision support techniques

George Kontolemakis*

National & Kapodistrian University of Athens, Department of Informatics
and Telecommunications, Information Systems Laboratory,
Panepistimiopolis, 157 84, Athens, Greece,
kontolem@di.uoa.gr

Abstract. The phenomenal growth of Internet-based information services and infrastructure in the recent years has provided a new technological basis for enabling and expanding the electronic execution of commercial transactions both on a business-to-business (B2B) and on a business-to-consumer (B2C) level. Electronic Marketplaces have increasingly played the role of an aggregator that merges potentially thousands of vendors and customers either as B2C virtual malls or as B2B electronic hubs. Virtually every working hub or marketplace focuses on either B2B or B2C business transactions. An integration of both categories would yield a generic e-hub made for all stakeholders across the process flows and covering every step of the way from production to consuming. The aim of this paper is to propose a novel architecture for the creation of economically viable e-hubs. We have argued that this can be accomplished through the ability of an agent-based electronic marketplace to transcend other taxonomical classification dimensions and, simultaneously, through the provision of an anthropocentric negotiation model and a flexible and “active” decision support. We have introduced a generic agent-based electronic marketplace architecture, comprising its three major components: ontology, negotiation and advising. According to the architecture proposed, we implemented the generic product ontology for the e-hub, focusing on its evaluation by using appropriate standards and widely accepted methodologies as well as other ontologies. We have also implemented a negotiation and a decision support system that come together and interact defining as a whole the functionality of the system.

Keywords. e-commerce, e-marketplace, e-hub, software agents, product ontology, negotiation, decision support

1 Dissertation Summary

The aim of this dissertation is the investigation of the electronic virtual communities and more specifically the application of software agents in electronic trade. Associated bibliography revealed the history of these societies, their growth, their technologies, their characteristics, their types as well as their margins of development [10], [13]. Thus general knowledge in the category of electronic marketplaces, a specialised sector of virtual societies was collected. Specifically, applications of software agents in the electronic trade and more specifically in electronic marketplaces were searched and developed [8]. These applications are generalised in three basic fragments of research: a) ontology, which is the heart of an electronic marketplace, b) negotiation, where all the activities of transaction are achieved and c) advisory services, which portray the supporting functionality of a marketplace. Combining these three technologies, a more general architecture in electronic marketplaces is proposed

*Dissertation Advisor: D.Martakos, Assoc.Professor

aiming at an essential confrontation of known problems and limitations that exist in them [11]. Then the classification of Kaplan & Sawhney [7] was studied, a classification standard for electronic marketplaces. The advantages and disadvantages of this classification were pointed out, as well as the possibility of transcending it through a generic architecture with the help of software agents [12]. By studying in detail this architecture, the three most important pieces of electronic marketplaces were developed. In the case of negotiation, a new anthropocentric system was developed and managed to put also in the game the purchaser with his rights, tactics, subjects and price limits [16]. In the advisory services, an economically viable marketplace with the help of a flexible and active department of support of decisions was proposed [3]. This system watches the actions of the user and gives advices according to his environment and demands. So through the systematic analysis, a generic product ontology that might include every possible product, every possible vendor and any prospective purchaser was developed [14]. Using techniques and already existing standards a new ontology development method was proposed and a general ontology was developed that includes negotiation and support of decisions transcending the Kaplan and Sawhney taxonomy.

1.1 Related Works

In recent years, electronic marketplaces have grown rapidly and formed marvellous giant marketplaces. For example, eBay, the most famous global electronic marketplace, has approximately 276 million registered users worldwide (http://news.ebay.com/fastfacts_ebay_marketplace.cfm). These users have posted a total number of 637 million new listings in the 4th quarter of 2007, i.e., averagely 6.7 million listings per day. In China, with an annual sale of RMB43.3 billion in the year of 2007, the dominant retail electronic marketplace—Taobao, defeated even the sum of local Carrefour and Wal-Mart and became the 2nd largest marketplace (http://forum.taobao.com/forum-14/show_thread----13526587-.htm). While the dramatically huge electronic marketplaces excite buyers by offering them abundant options, not surprisingly, they simultaneously make these options being too many to choose from. To help buyers locate their desired item, electronic marketplaces usually offer detailed item catalogues and powerful search engines. Even though, buyers still can find hundreds or thousands of items when search from electronic marketplaces.

Marketplaces can manage and present information about goods and trading status in different ways. Commodity markets typically offer a limited array of products, each described by a few parameters such as a stock name or a product grade. Other goods demand more extensive descriptions. On eBay, for example, there are usually more than four million different items for sale at any one time, and each is described by a few lines of text and perhaps a picture. Even if a user knows what item must be procured, there may be much difficulty in determining the identities of suppliers of that item and more importantly the suppliers that can supply the item with particular attributes at a particular price, for delivery before a particular date. Furthermore, even if the user is capable of determining an acceptable item-supplier combination, typically after a great deal of effort, the user is unable to determine whether a better overall deal could have been made through another supplier. As a result of any of these or other deficiencies, current procurement techniques have been inadequate for many needs. Electronic marketplaces are distinguished to controlled and uncontrolled marketplaces as shown in figure1 and discussed in [15].

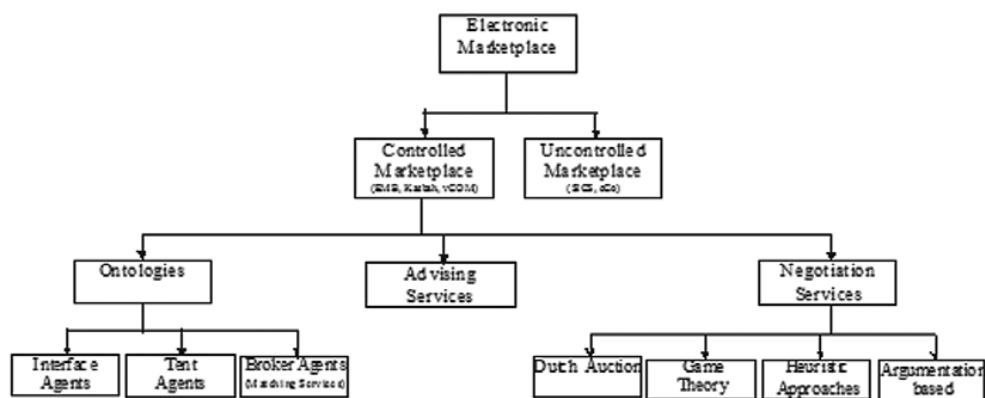


Figure 1: A categorization of the applications of agents in electronic marketplaces.

In a *controlled* marketplace the participants have to agree upon a certain set of rules concerning both what can be bought and sold and how this can be carried out. An *uncontrolled* marketplace is entirely open and decentralized; no single party for example sets the rules or controls the market. Each participant may initialize an agent that will act on its owner's interest using strategies uniquely defined for this agent. Uncontrolled electronic marketplaces are quite promising but they have to overcome an abundance of problems and are very difficult to implement. Although some interesting initiatives, such as the CommerceNet eCo System (www.commerce.net/eco/) exist, uncontrolled multi-agent marketplaces are still rather a vision than a reality. Agent-based functionality according to our research can be applied in the fields of ontologies, advising services and negotiations. Most of current commerce agents like "shopping bots" only support the Product or Merchant Brokering stages, only two of the six stages of the CBB model [4] that is used to explain the different stages of a deal. Few systems go into the negotiation part and even fewer help anticipate consumer needs (Need Identification stage) and provide paths into the subsequent CBB stages. But a flexible marketplace needs to support more stages to be successful. There is an ongoing research on the applications of agents in electronic marketplaces which was first proposed by [20]. The proposed architecture addresses four of the six stages of the CBB model. The problem with the marketplace proposed by [20] is that the agents are using a common syntax of a widely acceptable communication language. There is no ontology to strengthen the system making it rather difficult to become a widely accepted electronic marketplace architecture for agents to act.

Ontologies refer to models of a domain upon which agents rely for performing various tasks such as negotiation. Without defined ontologies, the application of agents in marketplaces and virtual communities will be severely limited and this fundamental need drives research on how ontologies can be shared and reused, how they can be revised when needed and how their consistency can be improved. Ontologies should be implemented in a way that they can be reused or even expanded with new terms but at the same time be resistant to structural revisions since they are created according to a logically consistent model [1]. For example the approach presented by [17] helps in maintaining large sets of transformation rules by providing for their decomposition into smaller and more understandable pieces and facilitating rule reuse.

When agent-based negotiation techniques were first proposed the requirement was reaching a better price for buying a product. As negotiation as a process evolved, additional information was needed so that a user could better decide for a purchase of a product and a plethora of negotiation objects appeared. These objects can be price, quality, timing, penalties, terms and

conditions or types of operation and are deemed as helpful in negotiating a product [6]. With agent-mediated negotiation, users need to be sure that the agent would achieve the best possible deal for them and that the product they are negotiating for is what they really want. Researchers emphasize on techniques that enhance trust amongst the user and his agent. This can be achieved by the continuous feedback given by the agent with additional information about the product and the negotiation phase. An agent should be able to support and advise the user for his/her actions.

A number of techniques have been proposed with each one aiming to enhance the efficiency and effectiveness of the negotiation process. The Dutch auction technique [6] is a very slow technique if none of the participants wants to buy the product. A viable solution is to provide the managing agent of the Dutch auction with additional meta-level information so as to speed up the process. Techniques borrowing principles from Game theory are generally regarded as much more efficient but suffer from one main limitation; the best possible solution is computationally intractable. With heuristic approaches, contracts that are closer to the opponent's last offer are provided but agents using these approaches often select outcomes that are sub-optimal so as to reach a deal. The best technique that has been proposed so far is the argumentation based technique where additional information over and above proposals is being exchanged.

The provision of advising services as a means for aiding the user to complete a specific task enhances the overall usability of a systems and is thus deemed critical. Traditionally, manuals and help files aided the user in the quest to find if a certain task can be performed by the system and how it can be done. Soon after, help files made their appearance enhanced with search and query capabilities. The main problem with these is that a user must know the syntax and semantics for asking the question or the answer will not be a good match to the original query. Research on agent-enabled advising services focuses on the intent, the timing and the level of intrusiveness of an advising service with researchers having proposed three styles of critic agents. These are 'Before-', 'During-' and 'After-Task' critics [19]. The main disadvantage of Before-Task critics is that as the information provided cannot be processed and filtered to match the exact user needs, redundancy and user overload is the result. During-Task critics are considered to offer the best possible advising service to the user; the drawback here being the user becoming fully dependent on the agent and the system without being able or willing to exercise any critical abilities or generate personal inferences. In contrast, After-Task critics do not distract the user, but they cannot prevent a wrong decision being made as any advice follows on the execution of the task. For agent-enabled advising services to advance in terms of usability future research should focus to a multi-style advising service using a mixture of 'During-' and After-Task style critics because the former can help in avoiding mistakes and the latter can add value in offering alternative solutions.

1.2 Innovative results of the dissertation

In this section we propose a generic and agent-mediated electronic marketplace architecture in an attempt to overcome existing impediments using latest research methods as shown in the previous section. The two main components in figure 2 are the buyer and the seller. They interact with each other via a negotiation and an active decision support system both exchanging data with the product ontology. We also show how the three basic components of an electronic marketplace (Ontology, Negotiation, and Decision Support) based on our previous research, come together and interact defining as a whole the functionality of the system. The first component is the Generic Product Ontology which is an ontology created so as to cover every possible product or input combinations which can be stored in the systems

database. The second one is the Negotiation Agent, who is responsible for managing the negotiation process between the buyer and seller using ontology attributes and for reaching a mutually acceptable promise which is then fulfilled through the logistics services. The third one refers to flexible and “active” decision support system. Flexible, in the sense, that it will accommodate all the diverse needs of the actors in the context of taxonomy classification transcendence and “active” in the sense that it will act proactively to support the decision making processes of the actors in contradiction with the traditional “passive” Decision Support Systems that required from their users to possess full knowledge of their capabilities and exercise initiative, something criticized since the late eighties [2].

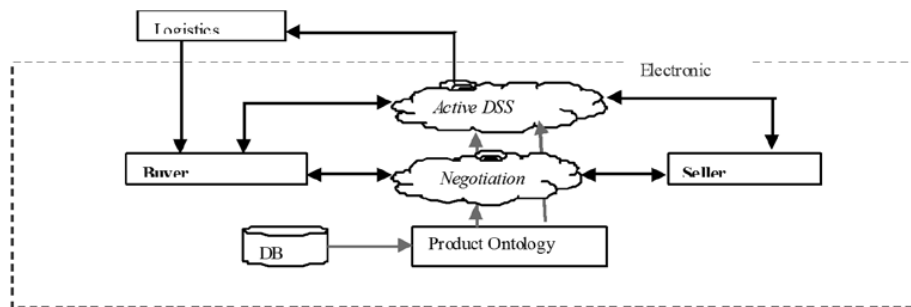


Figure 2: Generic and Agent-Mediated Electronic marketplace Architecture.

2 Results and Discussion

Ontology is a taxonomic catalogue of concept types and relation types and is a difficult, long and crucial part of the Knowledge Acquisition (KA) process (<http://ksi.cpsc.ucalgary.ca/KAW/KAW98/blazquez/#FernandezEtAl97>). A few research groups are now proposing a series of steps and methodologies for developing ontologies. However, mainly due to the fact that Ontological Engineering is still a relatively immature discipline, each work group employs its own methodology. In recent years, product ontologies have attracted both industry and academia because of their potential contribution to solving integration problems in e-commerce systems [18]. The heterogeneity of product information is a critical impediment to efficient business information exchange. There is no uniform description for each product type among vendors. In electronic commerce activities involving interactions among different vendors (business-to-business model) or between one buyer and multiple vendors (consumer-to-business model), a common ontology for the products is critical. There are countless approaches for the categorization of goods, ranging from rather coarse taxonomies, created for customs purposes and statistics of economic activities, like the North American Industry Classification System (NAICS) and its predecessor SIC (<http://www.census.gov/epcd/www/naics.html>), to expressive descriptive languages for products and services, like eCI@ss (www.eclass-online.com), eOTD (www.eccma.org), or RNTD (www.rosettanel.org), UNSPSC (www.unspsc.org) and the Epistle (www.epistle.eu).

The EPISTLE (European Process Industry STEP Technical Liaison Executive) is a brand new research effort to create a generic product ontology that stems from the need of a place to store the meaning and map between different terminologies. Using STEP - ISO 10303 and Parts Libraries - ISO 13584 developers tried to create standard instances held in external files

(class libraries) , which are also standardised by ISO, using an externally maintained registry with continuous revision. There are some drawbacks though, that need to be dealt with as the developers themselves claim. At first, there are no generic tools for access and maintenance since merging different libraries is a particular problem. Secondly, a consistent format for all levels is needed. Thirdly, there is no sharp distinction between data and meta-data. Lastly, the domain class libraries are not yet published as an ontology on the web. Nevertheless, it is a worthwhile research effort that may become a global standard for product ontologies.

The UNSPSC_v8 library, widely cited as an example of a product ontology, provides an industry neutral taxonomy of products and services categories, but no standardized properties for the detailed description of products [5]. The UNSPSC_v8 library was originally developed in 1998 during a collaboration program between the United Nations (www.un.org) and the Dun & Bradstreet (www.dnb.com) company. This ontology is today a worldwide standard that provides a wide range accurate categorization of products and services with a growing ratio of about 230 new classes per 30 days, showing significant maintenance of existing entries [5]. UNSPSC is used widely in business, especially in electronic commerce system. For example, Commerce One's Commerce Chain Solution (<http://www.commerceone.com/solutions>) and Ariba.com's Network have adopted it in their work on product content management.

For the development of our ontology, since we are interested in implementing software, we adopted the software development process standard (IEEE 1074 1996) which helped us through the procedure of developing and testing our ontology. In figure 3 we show the transformation of the IEEE directions into an ontology task analysis.

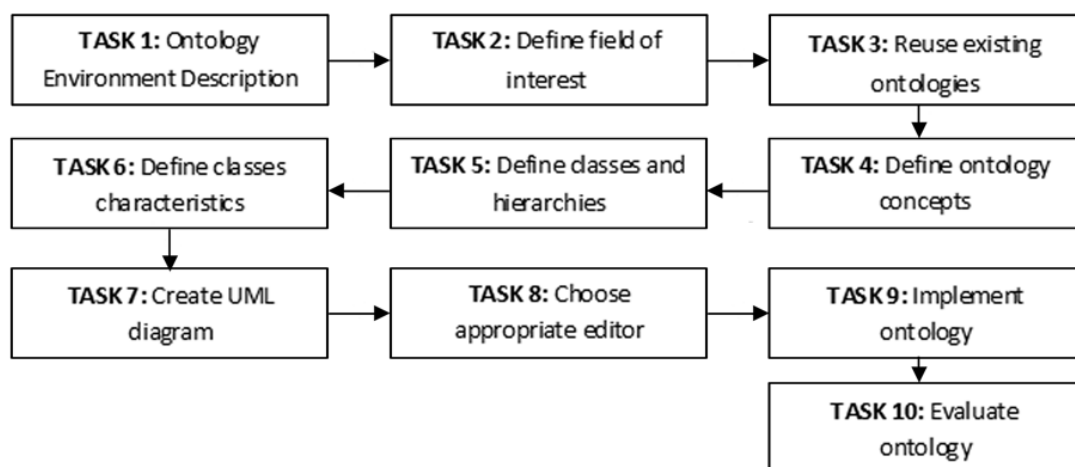


Figure 3: Ontology Task Analysis

According to the IEEE standard the first task belongs to the Pre-Development Processes. Accordingly, tasks 2-4, belong to Requirements Process, tasks 5-7 to the design process whilst tasks 9-10 refer to the implementation process. The last task is the one that was used for evaluation of the ontology and belongs to the Integral Processes. In the following paragraphs we deal with each task shown in figure 3 and gradually implement the system accordingly.

2.1 Environment description and field of interest

The generic product ontology aims to cover all products. This means that the ontology should describe the products, as well as aspects that concern their management. It might also provide information on the products that concern their natural characteristics. This information is extended by the product's materials, by its functionalism or even its presentation in the web. Moreover, through the ontology a product might participate in advisory as well as in negotiation processes. In this way, a complete knowledge of the product is defined, that gives the ability on the suitable users to handle with a dynamic and flexible way. All these dynamics of the proposed ontology will be used by all the likely purchasers and salesmen who will participate in one electronic marketplace. These will be also the users that will be responsible for the maintenance of the ontology. We will finally use the ontology as means of transcending the taxonomy of [7].

2.2 Reusing existing ontologies

The UNSPSC_v8 library is an invaluable tool for doing business globally although it has not addressed product attribute issues. Its hierarchical structure ensures that a company finds a meaningful level of product analysis conveniently. Its unique coding scheme makes it suitable for multi-language uses.

According to this ontology, products are taxonomized according to a general category (segment), a sub-category (family), a class of products and finally the product which is the general category of our implementation (Figure 4). This means that we are using an existent ontology as metadata for our own generic product ontology. This also means that the user of the ontology should have in his disposal nearly every product or service exists. This ontology is continuously updating which means that our ontology has no fear of not including every new product available, introducing standardized properties for the detailed description of them. Finally, we address the problem of the UN ontology by designing and implementing product attributes.

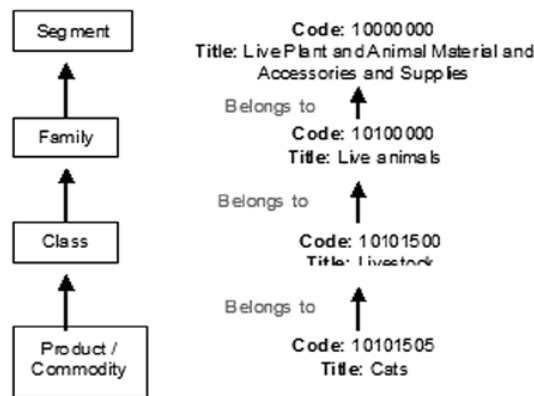


Figure 4: UNSPSC_v8 Ontology

2.3 Defining Ontology concepts

While the structure and properties of the standards such as the UNSPSC are known in advance and can be used for the comparison of alternatives, the actual coverage and level of

detail provided in a given category of products is not obvious. This leads to a situation where the decision for a standard is based mainly on its skeleton (e.g. whether it in general provides properties for a more detailed description of a category) and not on the degree to which such properties are actually defined for the product range of interest [5]. We therefore, after several interviews with stakeholders in the field and research co operations in negotiation and advising, developed a generic product ontology providing detailed description of products that are using the UNSPSC Ontology as meta-ontology and is shown in figure 5. The main class, common to the product class of the UNSPSC, is the **Product**. This class defines two subclasses which distinguish products to raw materials and final products. These are **Manufacturing Inputs** and **Operational Inputs**. With this simple categorization every product can be included in the ontology. The **Identifier** contains the product id along with some recognition details. These recognition details are dealing with the Name of the product, the Color, the Weight, the Size in all possible measurements, a description of the product, its packaging and of course detailed description of the Manufacturer of the product. The **Physical** property corresponds to a single material when we talk about manufacturing inputs or to a collection of raw materials or other products so that when synthesized an operating input is created. It contains a Code, a Name, the Origin which means the manufacturer of the product and the Type of product that is stating its input type. These inputs are also supported by the two isA relationships to the product. Simply stated, a company that in the past transcended these categories only in the physical world, it can now do it in the virtual. So both manufacturing and operating products are supported by the ontology. The Functional property refers to the possible applications of the product. This is crucial to the proper advising of a best fit product. The **Presentational** property is related to the way in which the product is represented to the user. It contains Type, Path, and Size of the file presenting the product. This means that the product can be seen from any known viewer on the internet. If any additions are required then the “Add_On” field is responsible for downloading the proper format or viewer to the system.

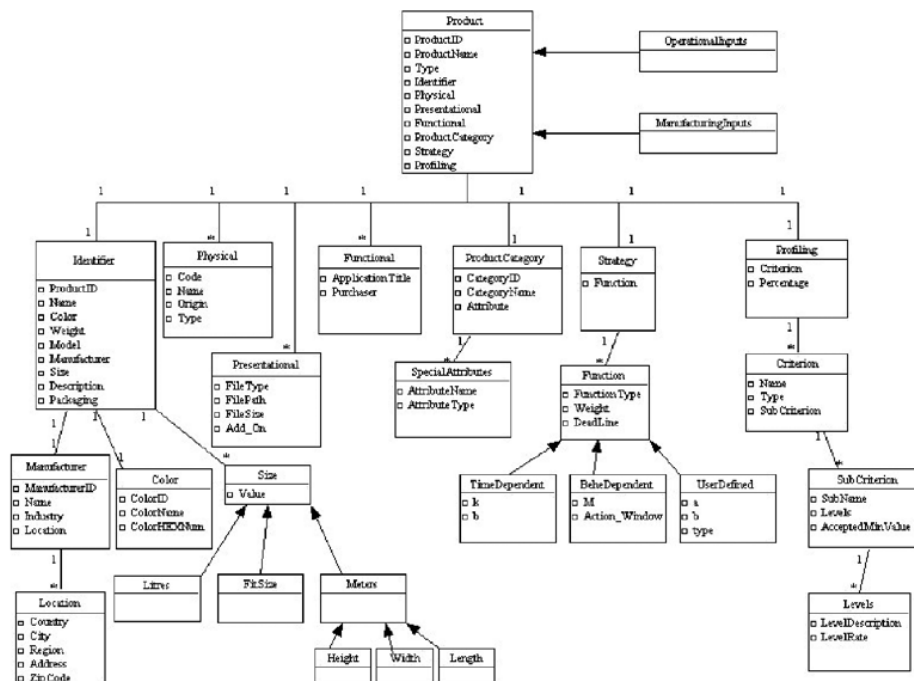


Figure 5: Generic Product Ontology (UML)

The **Product Category** property provides the vendor with the ability to classify his product into a broader category. Each category is assigned with specific properties called Special Attributes. The **Special Attributes** property includes alternate characteristics or meta-attributes of a product. This property contributes to producing a flexible system since additional product attributes are not predefined by the ontology, but can be created at run-time by appropriately configuring the Product Category. The seller is responsible to apply side-categories and product profiles, as well as the management and classification of the products he wants to sell. According to this ontology, the Product Category property provides the seller with the ability to classify his product into a broader category. Each category is assigned with specific properties called Special Attributes. The **Special Attributes** property includes alternate characteristics or meta-attributes of the product used for negotiation. For every characteristic a name and a permissible value type are inserted by the user. These include Fix Number for numeric values, String Sequence for one or more strings and Ranged Space for attributes within a specified set. This property produces a flexible system since additional product attributes aren't predefined by the ontology, but can be created at run-time by appropriately configuring the Product Category. In this way multiple issues can be negotiable facing the weakness of existing systems which use multi-issued negotiation with constant issues. Furthermore, taking into account that in manufacturing inputs quantity determines price, the ontology offering the Special Attributes property can accept ranged space attributes other than price, here quantity. In this way the E-hub supports both vertical and horizontal business purchases. In conjunction with the physical property it provides the flexibility to the user to promote his product or service in any way that he sees fit. It is worth mentioning that there is a predefined meta-knowledge given to the ontology automatically by embodying the UNSPSC_v8 library of products and services. So, the seller's product is probably included in the United Nations categorization. However, the company can incorporate into the hub products that differentiate her in comparison with other possible competitors.

Strategy is a property that helps the user to define his deal-making tactics based on the products' negotiable attributes. It comprises functions that are widely used during any negotiation nowadays. For every product inserted in the system, the seller can define tactics and strategies that will be used in the negotiation procedure. The ability of combining tactics such as time or behavior dependent is provided with the addition of a new tactic category we named "user-defined". This new tactic provides the ability to define a chosen function offers user the potential to actively participate in the negotiation process. This property could also be a property of the vendor as a whole but in real life and according to a shops stock, a vendor can have different strategies for example on a new model of a vacuum cleaner than on an older one.

Profiling is a property that allows one to define the characters of the people for whom the product will most likely have greater appeal, according to criteria defined by the users themselves. It consists of several criteria that can hold sub-criteria accordingly. These criteria can include cost reduction or discounts that can be handled at run-time favoring any side of the transaction. So the proposed architecture could prompt us to classify it as neutral but offering at the same time the flexibility to become either forward, reverse and biased.

Our goal is the generic product ontology to be used so that it can transcend the categorizations that [7] taxonomy imposes. Furthermore, it can unify all B2B and B2C mechanisms under one ontology and therefore under one E-hub.

3 Conclusions and further research

The aim of this dissertation has been to propose a novel architecture for the creation of economically viable e-hubs. We have argued that this can be accomplished through the ability of an electronic marketplace to transcend other taxonomical classification dimensions and, simultaneously, through the provision of an anthropocentric negotiation model and a flexible and “active” decision support.

We have introduced a generic agent-based electronic marketplace architecture, comprising its three major components: ontology, negotiation and advising. According to the architecture proposed we implemented the generic product ontology for the e-hub, focusing on its evaluation by using appropriate standards and widely accepted methodologies as well as other ontologies. We have also implemented a negotiation and a decision support system that come together and interact defining as a whole the functionality of the system.

The proposed electronic marketplace may provide enhancement of procurement experience amongst several sellers. It may also enable buyers to achieve optimal or near optimal procurement results with relatively little user interaction. This not only reduces errors and their associated transaction costs, but may significantly increase the speed, efficiency, and overall effectiveness of the procurement process. Buyers may also develop confidence that the procurement decisions made using this DSS and negotiation system result in the best overall deal. It may also provide the buyers with a much more detailed description of product, in terms of text, presentational and special attributes.

In terms of theory implications our research overcomes the most significant e-commerce issue: it can support both B2B and B2C e-commerce. In addition, it transcends a well-known taxonomy of e-hubs making the e-marketplace flexible to all possible market dimensions. It offers both systematic and spot sourcing, it supports both manufacturing and operating inputs, and it supports both neutral and biased behaviour. The ontology is by nature (UN taxonomy) continuously expanding and can be shared and reused. The negotiation system provides additional information without adding a considerable overhead to the system. The advising system matches exactly the user needs and present it in the most efficient and effective manner.

In terms of practice implications our research combines B2B and B2C marketplaces with the aid of a single ontology and not by trying to integrate differently developed ontologies (e.g. Alibaba). It also helps buyers locate their desired item effortlessly and with accuracy making the user able to determine the best overall deal. It can be used by any e-commerce vendors, whether small or large.

This work is just the beginning of a series of steps that must be done before this marketplace becomes a final product. The most important drawback is that it has not been tested in real time markets. To be tested in real markets this must include a huge effort of data entry that must be done by a series of businesses around the globe.

It is surely needed to use the ontology developed in a variety of case studies; typically these case studies should span multiple industries, since the ontology is meant as a generic product ontology. We should also check to what extent existing product ontologies (or reference models for product data) fit our own ontology.

Finally, it is of far most importance that the issue of trust in the context of this e-hub to be examined since trust is one of the most important factors for any successful software.

References

1. Albers, M. Jonker, C.M., Karami, M. and J. Treur. An Electronic Market Place: Generic Agent Models, Ontologies and Knowledge, 2000, *Available at: <http://www.cs.umbc.edu/kbem/kim.pdf>*, last accessed March 30 2004.
2. Angehrn, A., and Jelassi, M.T. DSS Research and Practice in Perspective, *Decision Support Systems* 12(4), 1994, 267-275.
3. Chamodrakas, Ioannis, Kontolemakis, George, Kanellis, Panagiotis, & Martakos, Drakoulis. "Liquid" Electronic Marketplaces. *Project E-Society, Building Bricks*, Springer Boston, 366-379, 2007.
4. Guttman, R., Moukas, A., Maes, P. Agent-mediated Electronic Commerce: A Survey. *The Knowledge Engineering Review*. Cambridge University Press 13 (2), 1998, 147 – 159.
5. Hepp, M., Leukel, J., Schmitz, V. A. Quantitative Analysis of eCl@ss, UNSPSC, eOTD, and RNTD: Content, Coverage, and Maintenance, 2007.
6. Jennings N.R., Faratin, P., Lomuscio, A.R., Parsons, S., Sierra, C. and Wooldridge. M. Automated Negotiation: Prospects, Methods and Challenges, 2001, *Available at: <http://www.csc.liv.ac.uk/~mjw/pubs/gdn2001.pdf>*.
7. Kaplan S., Sawhney S., E-hubs: The New B2B Marketplaces, *Harvard Business Review* 78(3), 2000, 97-103.
8. Kontolemakis G., Kanellis, P. and Martakos, D. Software Agents for Electronic Marketplaces: Current and Future Research Directions. *Journal of Information Technology Theory and Applications*. Vol.6, No.1, 2004, 43-61.
9. Kontolemakis G., Kanellis, P. and Martakos, D. Software Agents: An Overview of Research Directions, *The Knowledge Engineering Review*, 2010, Cambridge Press, to be reviewed.
10. Kontolemakis G., Kanellis, P. and Martakos, D. (2005). *Virtual Communities*. Encyclopedia of Multimedia Technology and Networking, Vol.2. Edited By: Margherita Pagani. Published by Information Science Reference (an imprint of IGI Global). ISBN: 1-59140-561-0.
11. Kontolemakis, George, Masvoura, Marisa, Kanellis, Panagiotis, & Martakos, Drakoulis. (2005) *Software Agents and the Quest for a Generic Virtual Marketplace Architecture*. 2005 Information Resources Management Association International Conference, vol.2, 15-18 May, San Diego, California, USA.
12. Kontolemakis, George, Masvoura, Marisa, Kanellis, Panagiotis, & Martakos, Drakoulis. (2005). *Transcending Taxonomies with Generic And Agent-Based E-Hub Architectures*. In Proc. 7th International Conference on Enterprise Information Systems, ICEIS 2005, 24-28 May, Miami, USA, 297-300.
13. Kontolemakis G., Kanellis, P. and Martakos, D. (2009) *Virtual Communities*. Encyclopedia of Multimedia Technology and Networking, Second Edition, Vol.3, pp.1512-1519. Edited By: Margherita Pagani. Published by: Information Science Reference(an imprint of IGI Global). ISBN: 978-1-60566-014-1. Last viewed in: <http://www.igi-global.com/downloads/pdf/Pagani1512.pdf>
14. Kontolemakis G., Kanellis, P. and Papadopoulou, P. (2010). Towards a "liquid" Electronic Marketplace, *Electronic Commerce Research and Applications*, Elsevier Press, to be reviewed.

15. Kurbel, K. and Loutchko, I. A Framework for Multi-agent Electronic Marketplaces: Analysis and Classification of Existing Systems, 2002, Available at: http://www.bi.euv-frankfurt-o.de/de/research/project/fp_softcomp_veroeff_pdf/Anlage%203.pdf.
16. Masvoura, M., Kontolemakis, G., Kanellis, P., and Martakos, D. Design and Development of an Anthropocentric Negotiation Model. In *Proceedings of the Seventh IEEE International Conference on E-commerce*, CEC 2005, 19-22 July 2005, Munich, Germany, 383 – 386.
17. Omelayenko, B. and D. Fensel. Scalable document Integration for B2B Electronic Commerce, 2001, Available at: http://informatik.uibk.ac.at/users/c70385/ftp/paper/OF_SII4.pdf.
18. Shim, J. and Shim, S. S. Y. Ontology-based e-Catalog in e-commerce: Special Section, *Electronic Commerce Research and Applications* (5:1), 2006, 1.
19. Silverman, B., Critiquing Human Error: A Knowledge-Based Human-Computer Collaboration Approach, *Academic Press*, New York, 1992, 521-528
20. Viamonte, M. J, Ramos, C. A Model for an Electronic MarketPlace, *Agent Mediated Electronic Commerce*, European AgentLink Perspective. Lecture Notes in Artificial Intelligence 1991, Frank Dignum and Carlos Sierra, Springer, 2000, 115-125.

A multimedia content modeling and classification methodology using visual information for the protection of sensitive user groups.

Alexandros Makris *

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
amakris@di.uoa.gr

Abstract. The thesis concerns the problems of visual tracking and violence detection in video sequences. For the visual tracking problem, two feature fusion frameworks are presented. For violence detection, a system that classifies movie segments as violent or non-violent is proposed. The first tracking framework called '*Model Fusion via Proposal*' (*MFP*) framework, provides a way to efficiently fuse visual cues using independent trackers to construct an improved proposal distribution for the main tracker. The fusion method results in reduced computational requirements due to the better proposal and the gradual exploitation of the state space. The '*Hierarchical Model Fusion*' (*HMF*) framework, extends the MFP framework by integrating the multiple models into a single tracker which exploits all the visual cues. This way the robustness of the approach is further increased. To this end, we extended the Bayesian framework to allow the integration of multiple models and we derived a particle filtering based approximation algorithm which allows the efficient integration of complementary models of different complexity with redundancy of information. Both tracking frameworks use multiple object models to describe the target. This feature enables the development of an adaptation strategy which adds or deletes models to cope with target appearance changes. For violence detection, a system that classifies movie segments as violent or non-violent is proposed. The system fuses audio and visual information. The audio module uses state-of-the-art methods. The visual features concern the general motion in the scene, the detection of gunshots, and the motion of the detected people.

1 Introduction

1.1 Automated Video Understanding

A very active research area, automated video understanding, concerns the analysis of video to extract semantic information. Most vision systems adopt a bottom-up approach. First low level tasks such as motion segmentation, object recognition, and tracking are performed followed by scene analysis or event recognition

* Dissertation Advisor: Professor Sergios Theodoridis

to extract high level concepts about the scene. Its applications lie in the fields of surveillance, control, and analysis. *Motion segmentation* is usually the first step in surveillance systems and consists of detecting the objects of interest and segmenting them from the static background. *Object recognition* is an important step required for video understanding. The recognition may concern specific objects (e.g. a face of a specific person [15], [14], [3]) or object classes (e.g. vehicles, people [2], [4]). The objects can be described by different cues. The next crucial step in vision systems consists of *tracking* the objects of interest extracted by motion segmentation or by object recognition. This is essential to establish the correspondences between the detected objects from frame to frame and possibly reduce the computational cost by avoiding detecting the objects in every frame. Usually the detection process is much more complex than tracking. *Event recognition* is the step that naturally follows the tracking of the objects of interest [7], [12].

1.2 Contribution

The contribution of this thesis concerns the development of novel tracking algorithms based on the particle filtering framework. Furthermore, a method for violence detection in movies is presented, using a classification approach, with features stemming from the tracking of objects using the proposed algorithm and from other computer vision techniques. The work resulted in the following publications: [9], [11], [1].

Proposed Tracking Frameworks The most important issues of the PFs are the efficient and information-rich *target representation* and the selection of the *proposal distribution*. We tackle these issues by proposing two generic frameworks (Model Fusion via Proposal(MFP), Hierarchical Model Fusion(HMF)) for fusing visual cues within the Bayesian framework. The MFP framework, provides a way to efficiently fuse visual cues using independent trackers to construct an improved proposal distribution for the main tracker. The fusion method results in reduced computational requirements due to the better proposal and the gradual exploitation of the state space. The HMF framework, extends the MFP framework by integrating the multiple models into a single tracker which exploits all the visual cues. This way the robustness of the approach is further increased. To this end, we extended the Bayesian framework to allow the integration of multiple models and we derived a particle filtering based approximation algorithm which allows the efficient integration of complementary models of different complexity with redundancy of information. Additionally, we developed an adaptation technique, to automatically delete and re-initialize the auxiliary models.

Proposed Violence Detection Method A system that classifies movie segments as violent or non-violent is proposed. The system fuses audio and visual information to increase the robustness. Two independent modules for audio and video based classification were developed. The audio module has been developed

in [5] where more details can be found. The video module uses features which describe the amount and the direction of motion in the scene, the motion of the detected people, and the illumination variations caused by gunshots. The features are used to classify the segments in three activity classes according to the amount of human activity of the segment (no, normal, and high activity) and to two classes according to the existence or not of gunshots in the scene. Similarly, The audio module uses several features to classify the movie segments in one of several audio classes (e.g. music, speech, gunshots, fights). The output of the video and audio classifiers is fed to a meta-classifier which decides for the presence of violence in the segment. The system as well as the independent modules were tested in real movie dataset. Both the modules when used independently reach a satisfactory level of performance. The fusion methods boosts that performance resulting in a system that detects about 4 out of 5 violent events (80% recall) and about 1 out of 2 events classified as violent are indeed violent (50% precision).

2 Hierarchical Model Fusion Framework

2.1 Algorithm Description

In the HMF framework the target is represented by several models of increasing dimension, which are probabilistically linked. The parameter update for each model takes place hierarchically so that the simpler models, which are updated first, guide the search in the state space of the more complex models to relevant regions. The most complicated model (in terms of state dimension) and the last in hierarchy, is called main model and its parameters fully describe the target. The rest of the models are referred as auxiliary as the estimation of their state is not required by the application.

A simple example (see Figure 1) will clarify the proposed concept. Let us consider a case of a target of which we want to estimate the bounding box. We will use two models, an auxiliary that tracks a feature point in the target and the main model, the bounding rectangle. The state of the first model has two parameters, the point's coordinates $x_s = [i_{s_x}; i_{s_y}]$, while the rectangle model has three, the coordinates of its center and a scale parameter, $x_b = [i_{b_x}; i_{b_y}; s_b]$. When the tracking is initialized the relative position of the rectangle's center and the point's is measured. If the tracked object is rigid this relative position should be almost constant between two consecutive frames. Thus if the location of the feature point is found on the next frame we can infer the coordinates of the center of the rectangle. The advantage of this strategy is that we first search in a two-dimensional space for the feature point and then we search in an one-dimensional space for the scale instead of searching in a three-dimensional state space to locate the rectangle directly. One should argue here that the coordinates obtained from the feature point model might not be very accurate or that this strategy will fail for non rigid objects. These issues are addressed by relaxing the link between the two models which will be discussed in detail in the following.

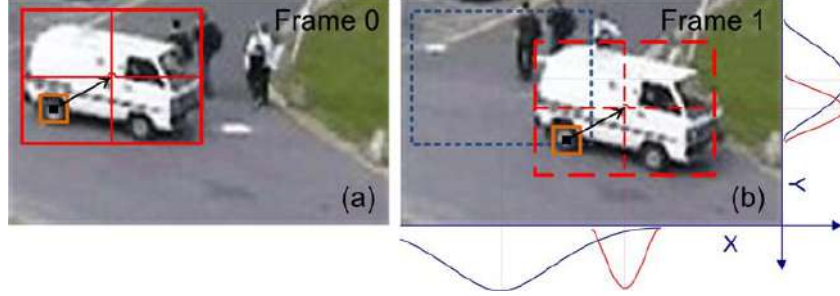


Fig. 1. In (a) the tracking is initialized with two models describing the target, a bounding rectangle and a salient point. The arrow shows the relative position of the salient point and the center of the rectangle. In (b) the position of the point is updated and using the stored relative distance the proposal for the rectangle given this position is shown in the x and y axis (red). This proposal is much closer to the target than the proposal derived by the state evolution model of the rectangle (blue).

In the general case M object models are used for target representation. The state can be written as ¹:

$$\mathbf{x} = [\mathbf{x}_{[1]}; \mathbf{x}_{[2]}; \dots; \mathbf{x}_{[M]}] \quad (1)$$

Where $\mathbf{x}_{[i]}$ are the state vectors of each object model. To each object model corresponds a measurement model with parameters $\mathbf{z}_{[i]}$. The graphical model of Figure 2c encodes the architecture of our framework. It depicts the following assumptions which we make to derive an algorithm for the recursive calculation of the posterior:

- The total likelihood is given by multiplying the likelihoods of individual models: $p(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^M p(\mathbf{z}_{[i]}|\mathbf{x}_{[i]})$.
- The state evolution is decomposed as:
 $p(\mathbf{x}_t|\mathbf{x}_{t-1}) = \prod_{i=1}^M p(\mathbf{x}_{[i]t}|\mathbf{Pa}(\mathbf{x}_{[i]t}))$.

where $\mathbf{Pa}(\mathbf{x}_{[i]t})$ denotes the parent nodes of $\mathbf{x}_{[i]t}$.

To construct algorithms that will be able to update the posterior of each model sequentially we derive the following equation which is an extension to the classical Bayesian tracking equation, but takes place in M steps. Each step updates the state of the corresponding model. Using simple probability rules and the assumptions mentioned above the filtering equation for the i – th step is given by:

$$\begin{aligned} p(\mathbf{x}_{[1:i]t}, \mathbf{x}_{0:t-1} | \mathbf{z}_{[1:i]t}, \mathbf{z}_{1:t-1}) &= \\ &= p(\mathbf{x}_{[1:i-1]t}, \mathbf{x}_{0:t-1} | \mathbf{z}_{[1:i-1]t}, \mathbf{z}_{1:t-1}) \frac{p(\mathbf{z}_{[i]t} | \mathbf{x}_{[i]t}) p(\mathbf{x}_{[i]t} | \mathbf{Pa}(\mathbf{x}_{[i]t}))}{p(\mathbf{z}_{[i]t} | \mathbf{z}_{[1:i-1]t}, \mathbf{z}_{1:t-1})} \end{aligned} \quad (2)$$

¹ For notational clarity we avoid using the T superscript, instead we use the Matlab's $'$ notation to concatenate vectors.

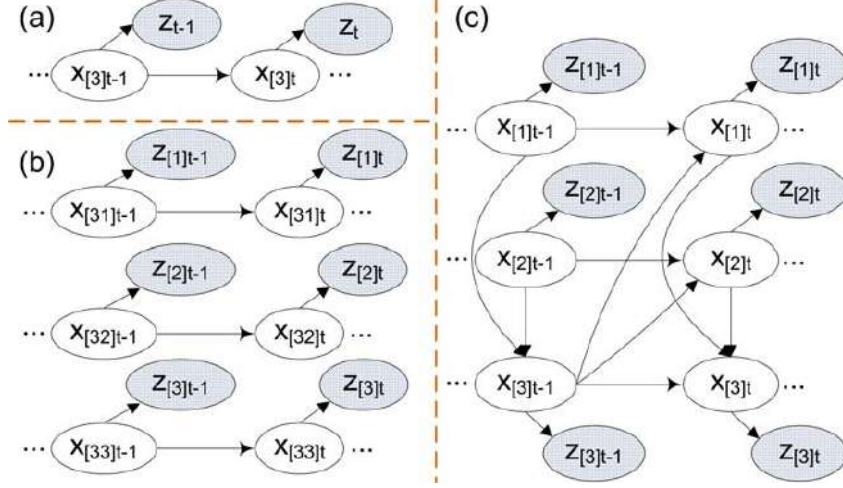


Fig. 2. Graphical Models depicting only slices $t - 1$ and t of the temporal dimension. The state to be approximated is denoted as $x_{[3]}$. (a) Standard Particle Filter [8]. (b) Partitioned Sampling [10],[13], the state is partitioned in 3 parts ($x_{[31]}$, $x_{[32]}$, $x_{[33]}$), which are updated independently, each one depending on different measurements. (c) Proposed Graphical Model, using 2 auxiliary models ($x_{[1]}$, $x_{[2]}$) which are linked to the main model ($x_{[3]}$).

As mentioned above, equation (2) can be used to construct Bayesian tracking algorithms that use multiple object models. Here, we use it to construct a PF based algorithm to iteratively update the posterior of each model.

We assume that at time $t - 1$ the posterior $p(\mathbf{x}_{[1:M]0:t-1} | \mathbf{z}_{[1:M]1:t-1})$ is approximated by a weighted particle set comprised of N weighted sample trajectories: $\{\mathbf{x}_{[1:M]0:t-1}^{(n)}, w_{[M]t-1}^{(n)}\}_{n=1}^N$. To update the particle set that approximate the posterior at time t we proceed in a sequential fashion. Each model is updated using the information from the already updated models at time t . The proposal distribution is selected to factorize as:

$$q(\mathbf{x}_{[1:i]t}, \mathbf{x}_{0:t-1} | \mathbf{z}_{[1:i]t}, \mathbf{z}_{1:t-1}) = q(\mathbf{x}_{[i]t} | Pa(\mathbf{x}_{[i]t}, \mathbf{z}_{[i]t})) q(\mathbf{x}_{[1:i-1]t}, \mathbf{x}_{0:t-1} | \mathbf{z}_{[1:i-1]t}, \mathbf{z}_{1:t-1}) \quad (3)$$

As in standard PF the samples are drawn from the first factor of Eq. (3).

The weights are given by:

$$w_{[i]t}^{(n)} = \frac{p(\mathbf{x}_{[1:i]t}^{(n)}, \mathbf{x}_{0:t-1}^{(n)} | \mathbf{z}_{[1:i]t}, \mathbf{z}_{1:t-1})}{q(\mathbf{x}_{[1:i]t}^{(n)}, \mathbf{x}_{0:t-1}^{(n)} | \mathbf{z}_{[1:i]t}, \mathbf{z}_{1:t-1})} \quad (4)$$

By substituting equations (2) and (3) into (4) we get the following weight update equation:

$$w_{[i]t}^{(n)} \propto w_{[i]t-1}^{(n)} \frac{p(\mathbf{z}_{[i]t} | \mathbf{x}_{[i]t}^{(n)}) p(\mathbf{x}_{[i]t}^{(n)} | Pa^{(n)}(\mathbf{x}_{[i]t}))}{q(\mathbf{x}_{[i]t}^{(n)} | Pa^{(n)}(\mathbf{x}_{[i]t}), \mathbf{z}_{[i]t})} \quad (5)$$

The steps of the iteration for the update of model i at time t of the proposed algorithm are the following (see Figure 3):

Given the particle set:

$$\{\mathbf{x}_{[1:i-1]t}^{(n)}, \mathbf{x}_{[0:t-1]}^{(n)}, w_{[i]t}^{(n)}\}_{n=1}^N:$$

1. **Sample:** For $n = 1$ to N draw $\mathbf{x}_{[i]t}^{(n)}$ from $q(\mathbf{x}_{[i]t}^{(n)} | Pa(\mathbf{x}_{[i]t})^{(n)}, \mathbf{z}_{[i]t})$.
2. **Update** the weights of each particle using Eq. (5).
3. **Normalize** the weights.
4. **Resample** the particle set according to its weights, so that the resulting particle set will be un-weighted and with the same number of particles.

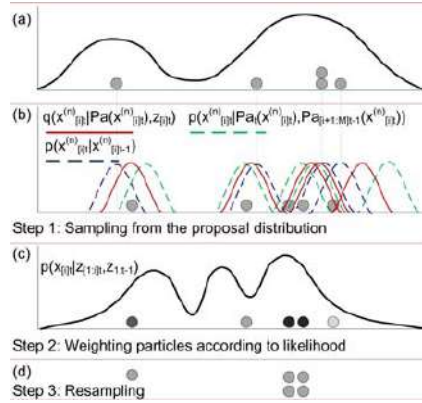


Fig. 3. Update for model i at time t . (a) The pdf and particles at time $t-1$. (b) The proposal is formed by fusing information from the current model's previous state and from the rest of the models. $Pa_t(\mathbf{x}_{[i]t})$ denotes the parent nodes of $\mathbf{x}_{[i]t}$ from time t , $Pa_{[i+1:M]t-1}(\mathbf{x}_{[i]t})$ denotes the parent nodes of $\mathbf{x}_{[i]t}$ from time $t-1$ excluding $\mathbf{x}_{[i]t-1}$. (c) The new particles are weighted. Darker particles have higher weight. (d) Resampling.

2.2 Tracker Implementation

In this section we use the framework to build a tracker, which we applied in tracking various targets in challenging situations. We combined three different object models to represent the target which are in the order which are updated:

- (i) A salient point tracking model. This model has only 2 position parameters.
- (ii) A blob tracking model. The blob represents a rectangular region of the target with homogeneous color with 3 parameters which describe its position and scale.
- (iii) The target's contour. This is the main object model. It is represented as a b-spline curve and contains 5 parameters which allow several geometric transformations.

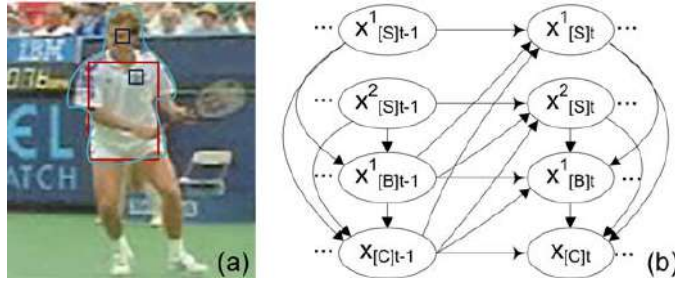


Fig. 4. (a) Sample image showing two salient points (dark), one blob (red) and one curve model (blue). (b) Graphical model of the implemented tracker, depicting slices $t-1$ and t of the temporal dimension for the aforementioned models. The evidence nodes are omitted for clarity.

The combined state vector is: $\mathbf{x} = [\mathbf{x}_{[S]}; \mathbf{x}_{[B]}; \mathbf{x}_{[C]}]$, where $\mathbf{x}_{[S]}$ represents the salient points $\mathbf{x}_{[B]}$ represents the blobs and $\mathbf{x}_{[C]}$ represents the contour curve. For a single target more than one salient points or blobs models can be used.

Adaptation Procedure So far we proposed a framework to fuse information from different cues using several object models which are initialized at the first frame. Here, we will expand the proposed framework to adapt the auxiliary models during tracking using information from the main model. The adaptation of the auxiliary models is integrated in the update equation. Each auxiliary model has a parameter that encodes the confidence of the object to belong to the target (target confidence). The adaptation consists of deleting an object if the target confidence is below a threshold and detecting and initializing new objects. For the k_s salient point model the target confidence parameter $\mathbf{x}_{[s_t]}^{k_s}$ is initialized when the point is detected and updated during tracking by the following filtering equation:

$$\mathbf{x}_{[s_t]}^{k_s} = a_{s_{tc}} \mathbf{x}_{[s_t]}^{k_s} + (1 - a_{s_{tc}}) f_{stc}(\mathbf{x}_{[S]}^{k_s}, \mathbf{x}_{[C]}) \quad (6)$$

Where $a_{s_{tc}}$ is the filtering parameter and t_{stc} is the threshold for deleting an auxiliary model. $f_{stc}(\cdot)$ is a metric measuring the compatibility between the

points and the contour model on the previous frame:

$$f_{stc}(\mathbf{x}_{[S]t-1}^{k_s}, \mathbf{x}_{[C]t-1}) = \exp \left\{ -\frac{d_{ht}^2(L_{[S]}^{k_s}, L_{[C]})}{2\sigma_{stc}^2} \right\} \quad (7)$$

where σ_{stc} is the deviation, $L_{[S]}^{k_s}$, $L_{[C]}$ are the likelihood vectors for the k_s -th salient point and the curve model respectively defined as

$L_{[S]}^{k_s} = [p(\mathbf{z}_{[S]}^{k_s} | \mathbf{x}_{[S]}^{k_s(1)}); \dots; p(\mathbf{z}_{[S]}^{k_s} | \mathbf{x}_{[S]}^{k_s(N)})]$ and similarly for the curve model. This equation models the similarity between the likelihood vectors which is high when the two models describe the same target. In that case the link from one model to the other is meaningful. In contrast when one of the models is distracted by clutter then the similarity between the two vectors is expected to be lower. The same equations hold for the initialization and update of the target confidence parameter of the blob model $\mathbf{x}_{[b_t]t}^{k_b}$.

When $x_{[s_t]t-1}^{k_s(n)} < t_{stc}$ or $x_{[b_t]t-1}^{k_b(n)} < t_{stc}$ for the k_s -th salient point and k_b -th blob models respectively then the auxiliary model is deleted and a detection procedure searches for new salient points or blobs to re-initialize in the target region as defined by the main model.

3 Tracking Experiments

The experiments have been executed using several challenging video sequences and various objects have been tracked to verify our methods. More specifically, we experimented with sequences containing deformable objects, abrupt motion, heavy clutter, partial occlusions, and short full occlusions. For the experiments we implemented the following trackers which we compare: **SIR** - the original SIR algorithm, **MFP** and **aMFP** - the MFP tracker with and without adaptation, **HMT aHMT** - the HMF tracker with and without adaptation.

To compare the trackers we annotated several sequences by hand and we calculated the ‘Tracker Detection Rate’ (TDR) and ‘False Alarm Rate’ (FAR) measures [6]:

$$TDR = \frac{TP}{TP + FN}, FAR = \frac{FP}{FP + TP} \quad (8)$$

where TP , FN and FP denote the true positive, false negative and false positive area respectively.

In the experiment displayed in Figure 5 we compare our HMT tracker to the SIR in a PETS 2006 surveillance sequence using the blob and contour models. The contour hypotheses of our method are much more concentrated near the actual target than those of the SIR because of the strong prior provided by the blob model.

The experiment displayed in Figure 6, illustrates the concept of the model adaptation using the aHMT tracker with two model types, salient points and contour. The salient point models are deleted and re-initialized as the initial

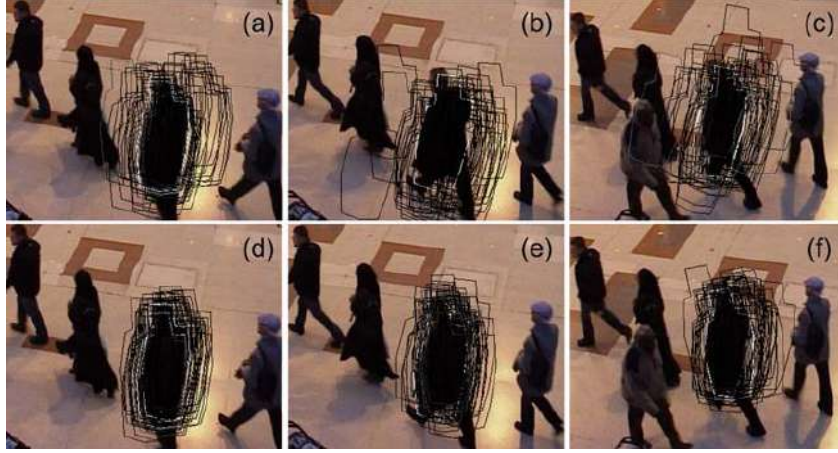


Fig. 5. Tracking Results - Surveillance Sequence:(a),(b),(c)-SIR, 50 particles (d),(e),(f)-HMT, 50 particles. Frames 1,35,50. Both trackers use the blob and the contour models, for image clarity we only show the contour particles. The particles of the HMT tracker are more concentrated near the target due to the better proposal distribution.



Fig. 6. Tracking Results - 7up Sequence:aHMT, 60 particles, frames 1,180,300

points are occluded due to the object rotation. The main model (contour), defines the target area and the search for new points is performed there.

In several cases some type of auxiliary models do not provide valid information for the target. In such cases the adaptation mechanism discards these models without replacing them. One such situation is observed in Figure 7 where the blob models are misled by similar background colors and are quickly discarded. The spots are not helpful throughout the whole sequence and therefore are discarded for several frames as well.

In the HMT framework the object models are connected and form a single graphical tracking model whereas in the MFP framework several independent trackers are used with each one having a single object model and they exchange information only through the proposal distribution. This type of connection

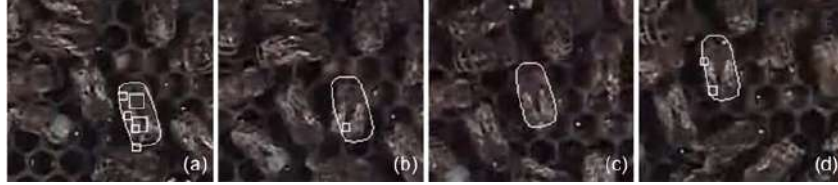


Fig. 7. Tracking Results - Sequence without blobs:aHMT, 50 particles, frames 1,15,30,50. The blob models, although initialized in (a), are quickly discarded because the background contains similar colors that distract it. In (c), the spot models are also discarded and only the contour is used. In (d), several new spots are detected.

between the models does not guarantee that the trackers will remain locked on the same target, especially when the adaptation method is not used. In Figure 8, this point is highlighted. The contour model is distracted by the player of the other team and destabilizes the tracker.

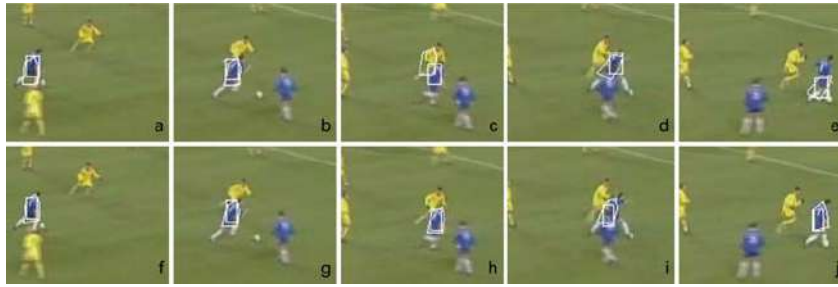


Fig. 8. Tracking Results - Soccer Sequence:(a)-(e) MFP, 50 particles, (f)-(j) HMF, 50 particles, frames 1,21,24,29,40. In (c), the contour model is distracted by the player of the other team as opposed to the HMF tracker which is not distracted as seen in (h).

4 Conclusions

Two feature fusion frameworks for visual tracking were presented. The implemented trackers were used in the following to create features for the developed violence detection system, which classifies movie segments as violent or non-violent. The trackers are based on the particle filtering methods. Their main goal was to create better hypotheses thus reducing the computational cost and enabling the use of high dimensional object models in real time applications. A violence detection system was also developed. The proposed system fuses audio and visual information using features that do not restrict the scope of its application and was tested on a real film dataset.

References

1. Anagnostopoulos, V., Kosmopoulos, D., Doulamis, A., Makris, A., Lalos, C., Varvarigou, T.: Automated production of personalized video content for visitors of thematic parks. In: IE06. pp. 173–181 (2006)
2. Belongie, S., Malik, J., Puzicha, J.: Shape matching and object recognition using shape contexts. *IEEE Trans. Pattern Anal. Mach. Intell.* 24(4), 509–522 (2002)
3. Bowyer, K.W., Chang, K., Flynn, P.: A survey of approaches and challenges in 3d and multi-modal 3d + 2d face recognition. *Comput. Vis. Image Underst.* 101(1), 1–15 (2006)
4. Fergus, R., Perona, P., Zisserman, A.: Object class recognition by unsupervised scale-invariant learning. In: CVPR (2). pp. 264–271 (2003)
5. Giannakopoulos, T.: Study and application of acoustic information for the detection of harmful content, and fusion with visual information. PhD Dissertation, NKUA (2009)
6. Hall, D., Nascimento, J., Ribeiro, P., Andrade, E., Moreno, P., Pesnel, S., List, T., Emonet, R., Fisher, R.B., Victor, J.S., Crowley, J.L.: Comparison of target detection algorithms using adaptive background models. In: Visual Surveillance and Performance Evaluation of Tracking and Surveillance, 2005. 2nd Joint IEEE International Workshop on. pp. 113–120 (15–16 Oct 2005)
7. Hu, W., Tan, T., Wang, L., Maybank, S.: A survey on visual surveillance of object motion and behaviors. *Systems, Man and Cybernetics, Part C, IEEE Transactions on* 34(3), 334–352 (2004), <http://dx.doi.org/10.1109/TSMCC.2004.829274>
8. Isard, M., Blake, A.: Condensation – conditional density propagation for visual tracking. *International Journal of Computer Vision* 29(1), 5–28 (1998), [cite-seer.ist.psu.edu/isard98condensation.html](http://citeseer.ist.psu.edu/isard98condensation.html)
9. Kosmopoulos, D.I., Doulamis, A., Makris, A., Doulamis, N., Chatzis, S., Middleton, S.E.: Vision-based production of personalized video. *Image Commun.* 24(3), 158–176 (2009)
10. MacCormick, J., Isard, M.: Partitioned sampling, articulated objects, and interface-quality hand tracking. In: ECCV '00: Proceedings of the 6th European Conference on Computer Vision-Part II. pp. 3–19. Springer-Verlag, London, UK (2000)
11. Makris, A., Kosmopoulos, D.I., Perantonis, S.J., Theodoridis, S.: Hierarchical feature fusion for visual tracking. In: ICIP (6). pp. 289–292 (2007)
12. Pantic, M., Pentland, A., Nijholt, A., Huang, T.: Human computing and machine understanding of human behavior: a survey. In: ICMI '06: Proceedings of the 8th international conference on Multimodal interfaces. pp. 239–248. ACM, New York, NY, USA (2006)
13. Perez, P., Vermaak, J., Blake, A.: Data fusion for visual tracking with particles. *Proceedings of the IEEE* 92(3), 495–513 (2004)
14. Tan, X., Chen, S., Zhou, Z.H., Zhang, F.: Face recognition from a single image per person: A survey. *Pattern Recogn.* 39(9), 1725–1745 (2006)
15. Zhao, W., Chellappa, R., Phillips, P.J., Rosenfeld, A.: Face recognition: A literature survey. *ACM Comput. Surv.* 35(4), 399–458 (2003)

Nonlinear Signal Processing and its Applications to Telecommunications

Gerasimos Mileounis*

Department of Informatics and Telecommunications
National and Kapodistrian University of Athens
TYPA Buildings, University Campus, 15784 Athens, Greece
gmil@di.uoa.gr

Abstract. This dissertation is primarily concerned with the estimation of nonlinear communication systems that are modeled by Volterra series. The major methods used for estimating the unknown channel parameters can be classified into two main categories: training-based and blind. First, orthobasis representation and training-based identification through the respective Fourier series are investigated for most modulated signals of interest. Next, higher order cumulants are used for the blind identification of nonlinear channels. The proposed algorithms for blind nonlinear channel estimation take advantage of the inherent sparseness of the higher order cumulants of common communication signals. Then, sparse Volterra channels are employed to mitigate the enormous computational complexity of the full Volterra channels. Sparse Volterra channels are approached by two newly developed sparse adaptive (greedy and ℓ_1 -regularized) algorithms. Last, the problem of blind sparse channel estimation is formulated by modifying the Expectation-Maximization framework to accommodate channel sparsity.

Keywords: Volterra series, Higher-Order-Statistics, Adaptive filters, Blind identification, Nonlinear compressed channel sensing.

1 Introduction

Nonlinear behavior is observed in almost all digital communication systems including satellite, telephone channels, mobile cellular communications, wireless LAN devices, radio and TV stations, digital magnetic systems and so forth. In those cases, possible remedies based on linear approximations degrade system performance. Significant benefits in the performance of a digital communication system are expected when appropriate nonlinear models, methods and algorithms are developed, taking into account nonlinear effects.

Nonlinear systems have been systematically studied in the past [1], but they have not been widely used in communications due to their computational complexity. Computational complexity depends on several factors, including: (1) dimension of the unknown parameter vector, (2) sparseness characteristics of the parameter vector, (3) degree of nonlinearity, (4) number of available measurements, and (5) the probability density function of the input. The motivation

* Dissertation Advisor: Prof. Nicholas Kalouptsidis.

of this dissertation is to develop methods and algorithms that considerably reduce the computational complexity and increase the performance of existing algorithms for nonlinear channel estimation, by using the above factors in a beneficial manner.

We study nonlinearities in communication systems using polynomial filters, a special class of which is Volterra series. More precisely, we propose new estimation techniques and apply them to linear and nonlinear communication channels. Our research efforts focus on two areas: (1) nonlinear channel estimation using higher order statistics, and (2) adaptive and blind algorithms for sparse channel estimation.

Initially, we develop training based algorithms for the identification of (passband and baseband) Volterra channels modulated by QAM, PSK and OFDM inputs [2]. When the Volterra channel is excited by QAM or PSK inputs, multivariate orthogonal polynomials are used to estimate the unknown parameters. Closed form expressions are established for baseband Volterra channels driven by i.i.d complex Gaussian (OFDM) signals.

Blind methods identify the unknown channel merely based on the received signal, without consuming any of the available channel capacity. However, blind nonlinear channel estimation is a hard problem and the development of blind methods remains at a very preliminary stage dealing with special model subclasses. We investigate sparseness of the higher-order output cumulants in order to simplify the blind identification problem of two different nonlinear models: (1) passband and baseband Hammerstein channels excited by common communication signals [3], and (2) linear-quadratic Volterra with complex random inputs [4].

Volterra models employ a large number of parameters, to adequately represent many real-world systems, which increases exponentially with the order of nonlinearity and memory length. For this reason, their applicability is limited to weak nonlinearities, e.g. only up to third order. Therefore, there is a strong need to decrease the parameter space by only considering those parameters that actually contribute to the output. This observation lead us to the exploitation of sparse Volterra models which constitute a major component of this dissertation.

The use of adaptive filtering is crucial in applications like communications where channel measurements arrive sequentially and in many cases the channel response is time-varying. All adaptive algorithms in the literature for nonlinear channel estimation treat each parameter equally and identify the complete set of parameters. The major drawback of estimating the complete set of parameters is the large computational/implementation cost. In this dissertation, we investigate the performance gains that can be achieved if insignificant parameters are ignored. Two Novel adaptive algorithms are developed that recursively update the parameters of interest. The first adaptive algorithm combines the Expectation-Maximization and Kalman filtering [5], whereas the second one relies on greedy methods [6]. Finally, using the Expectation-Maximization framework, we address the problem of blind identification of sparse linear and nonlinear channels [7].

This summary is organized as follows. Firstly, Chapter 2 establishes the necessary background needed in the sequential chapters. Section 3 deals with training-based methods for the identification of Volterra channels. Sections 4 and 5 tackle

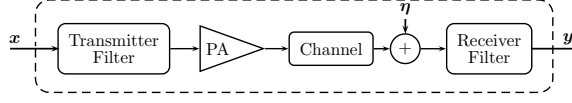


Fig. 1. Nonlinear communication system

the problem of blind identification in Hammerstein channels and second order Volterra systems, respectively, by making use of Higher-Order output Statistics. Two different adaptive training-based algorithms for the estimation of sparse channels are presented in Sections 6 and 7. Section 8 proposes a blind identification algorithm for the estimation of sparse channels. Finally, Section 9 presents the overall conclusion of this summary.

2 Background

Modern high-speed communication systems are frequently operated over nonlinear channels with memory. Most transmitters are equipped with Power Amplifiers (PAs) operating close to saturation to achieve power efficiency [8]. To properly analyze a communication system, like the one of Fig. 1, the nonlinear effects caused by the presence of PAs must be combined with the transmitting, receiving and channel filters.

One of the most popular models that are applied for the description of nonlinear phenomena are Volterra series [1] and allows us to capture these combined effects. Causal discrete Volterra series of finite-order have the following form (also referred to as *passband Volterra*):

$$y(n) = \sum_{p=1}^P \sum_{\tau_1=0}^{M_p} \cdots \sum_{\tau_p=0}^{M_p} h_p(\tau_1, \dots, \tau_p) \left[\prod_{i=1}^p x(n - \tau_i) \right]. \quad (1)$$

where $x(n)$ and $y(n)$ are the system input and output respectively. The function $h_p(\tau_1, \dots, \tau_p)$ is called the *pth-order Volterra kernel* of the system. P is the highest order of nonlinearity while M_p is the *pth-order system memory*, note that the Volterra model of Eq. (1) becomes linear when $P = 1$.

In many communication systems the signal bandwidth is very carefully defined depending on the application. The receiver filter is used to eliminate signal components outside the desired bandwidth. Therefore the output signal only contains spectral components near the carrier frequency ω_c . This leads to the *baseband Volterra* system [8, Ch. 14], given by

$$y(n) = \sum_{p=0}^{\lfloor \frac{P-1}{2} \rfloor} \sum_{\tau_1=0}^{M_{2p+1}} \cdots \sum_{\tau_{2p+1}=0}^{M_{2p+1}} h_{2p+1}(\tau_1, \dots, \tau_{2p+1}) \prod_{i=1}^{p+1} x(n - \tau_i) \prod_{j=p+2}^{2p+1} x^*(n - \tau_j) \quad (2)$$

where $\lfloor \cdot \rfloor$ denote the floor operation. The above representation only considers odd-order powers with one more unconjugated input than conjugated input.

This way the output does not create spectral components outside the frequency band of interest.

The key feature of Volterra series is that the nonlinearity is due to multiple products of delayed input values, while the kernel coefficients appear linearly in the output. This allows us to rewrite them as a linear regression model using Kronecker products. Indeed, consider the passband case of Eq. (1) and let $\mathbf{x}_{M_1}^{(1)}(n) = [x(n), x(n-1), \dots, x(n-M_1)]^T$ and the p th-order Kronecker power $\mathbf{x}_{M_p}^{(p)}(n) = \mathbf{x}_{M_1}^{\otimes p}(n)$. The Kronecker power contains all p th-order products of the input. Likewise $\mathbf{h} = [\mathbf{h}_1(\cdot), \dots, \mathbf{h}_p(\cdot)]^T$ is obtained by treating the p -dimensional kernel as a $(M_p)^p$ column vector. We now rewrite the output of Eq. (1) as follows

$$y(n) = \left[\mathbf{x}_{M_1}^T(n) \cdots \mathbf{x}_{M_p}^{(p)T}(n) \right] \begin{bmatrix} \mathbf{h}_1 \\ \vdots \\ \mathbf{h}_p \end{bmatrix} = \mathbf{x}^T(n) \mathbf{h}. \quad (3)$$

Because of this property, linear estimation techniques can be exploited in the identification of Volterra coefficients.

3 Nonlinear system identification using orthogonal bases and cumulants

In this section baseband and passband nonlinear channels featuring PSK, QAM and OFDM modulation are considered. Channel estimation is performed using multivariate orthogonal polynomials of complex variables.

An i.i.d complex valued signal is orthogonalizable [9] and there are various ways to construct associated orthogonal bases. A common construction relies on one dimensional orthogonal base and its separable extension to higher dimensions. One dimensional polynomials constructed by the Gram-Schmidt procedure is a notable case. Multidimensional orthogonal polynomials are formed as products of one dimensional orthogonal polynomials ($P_k(x_i)$, where k is the degree of the polynomial and $x_i \equiv x(n-i)$) [9]. For the monomials (in one variable) $\{x_n^i\}_{i=0}^p$ associated with passband Volterra models, we introduce a degree ordering $1 < x_n < \dots < x_n^p$. For monomials in two variables ($x_n x_n^*$) related to baseband models we apply the graded lexicographic ordering [2].

By definition the multivariate orthogonal polynomials, $Q_{\mathbf{i}_{1:p}}^{(p)}(\mathbf{x}(n))$ of degree p , are orthogonal to all lower orders and to the same order. The passage to the original Volterra kernels from the orthogonal coefficients is effected by the following expression

$$\begin{aligned} \mathbb{E}[y(n) P_{\tau_1}^*(x_{i_1}) \cdots P_{\tau_k}^*(x_{i_k})] &= \pi(\mathbf{i}_{1:k}) h_k(i_1, \dots, i_k) \|P_{\tau_1}(x_{i_1})\|_{\ell_2}^2 \cdots \|P_{\tau_k}(x_{i_k})\|_{\ell_2}^2 \\ &+ \sum_{v=1}^{\lfloor \frac{p-k}{2} \rfloor} E \left\{ \mathcal{H}_{k+2v}(\underline{z}_n) Q_{i_1 \dots i_k}^{(*p)}(\mathbf{x}(n)) \right\} \end{aligned}$$

where $\mathcal{H}_v(\cdot)$ is the homogeneous term of order v . Similarly for the baseband case. The identification process starts by estimating the highest order kernel which has no contribution from other kernels and moving downwards.

If the input is approximately complex white Gaussian (OFDM case) the relevant orthogonal polynomials are the Hermite polynomials. The method described above is applicable. Alternately cumulant operators can be used. Expressions invoking cumulants are much simpler because cumulants are equivalent to multiples of Hermite moments.

Theorem 1 *Consider the baseband Volterra model (2). The cross-cumulant of $y(n)$ with $(p+1)$ conjugated copies of the input and p unconjugated copies of the input is given by*

$$c_{y, \mathbf{x}^*}^{(p+1)}(\boldsymbol{\tau}_{1:p+1}, \boldsymbol{\tau}_{p+2:2p+1}) = p!(p+1)! \gamma_{1,1}^{2p+1} h_{2p+1}(\boldsymbol{\tau}_{1:p+1}, \boldsymbol{\tau}_{p+2:2p+1}) \\ + \sum_{v=1}^{\lfloor \frac{P-2p-1}{2} \rfloor} \frac{(p+1+v)!(p+v)!}{v!} \gamma_{1,1}^{2p+1+v} \sum_{k_1} \cdots \sum_{k_v} h_{2p+1+2v}(\boldsymbol{\tau}_{1:p+1}, \mathbf{k}_{1:v}, \boldsymbol{\tau}_{p+2:2p+1}, \mathbf{k}_{1:v})$$

$\gamma_{1,1}$ represents the variance of $x(n)$ and $\mathbf{k}_{1:v} = (k_1, \dots, k_v)$.

Detailed proof of the above theorem is given in [2]. The algorithm identifies the highest order kernel first. Then the lower order kernels are identified recursively using the previous estimated kernels and cross-cumulants information.

4 Blind identification of Hammerstein channels

This section considers the blind identification problem of Hammerstein channels. The Hammerstein model corresponds to a diagonal Volterra model, since all off-diagonal elements are zero.

Baseband: The proposed method starts by expressing the baseband Hammerstein model (similarly for passband) as a linear multichannel system of the form

$$y(n) = \sum_{i=0}^{q_1} \mathbf{b}(i) \mathbf{w}(n-i) + \eta(n)$$

where

$$\mathbf{w}(n) = (x(n) \cdots |x(n)|^{2p} x(n))^T, \quad \mathbf{b}(i) = (h_1(i) h_3(i) \cdots h_{2p+1}(i))$$

and T denotes matrix transpose. We assume that the linear kernel has the largest memory $q_1 > q_{2l+1} \forall 1 \leq l \leq p$. Moreover, we find it convenient to impose the following normalization $h_1(q_1) = 1$.

The output cumulant of order $(k+l)$, with k unconjugate and l conjugate output lags, is given by:

$$c_{y(k)}^{(l)}(\tau_1, \dots, \tau_{k+l-1}) = \sum_{i=0}^{q_1} \mathbf{b}(i) \otimes \mathbf{b}(i + \tau_1) \otimes \cdots \otimes \mathbf{b}^*(i + \tau_{k+l-1}) \boldsymbol{\Gamma}_{\mathbf{w}(k)}^{(l)},$$

where $\mathbf{\Gamma}_{\mathbf{w}(k)}^{(l)}$ is the input intensity vector (zero lag cumulant) of order $k + l$ of $\mathbf{w}(n)$. Then, parameter estimation relies on the following equation and the solution of a system of linear equations

$$\tilde{c}_y^{(l)}(q_1, \tau) = c_y^{(l)}(q_1, \dots, q_1, \tau) = \mathbf{b}^*(\tau) \left(\bar{\mathbf{\Gamma}}_{\mathbf{w}(k)}^{(l)} \right)_{s \times s} \mathbf{b}^T(0), \quad s = p + 1.$$

For PSK inputs we always consider cumulants with an equal number of conjugate/unconjugate copies of the output. Whereas, for QAM inputs we may employ output cumulants with unequal number of unconjugate/conjugate entries of the output. In this manner we reduce significantly the order of the output cumulants.

Passband: Prakriya et al. [10] have proved that if the order of nonlinearity is p , then the only non-zero multilinear function of $c_y^{(1)}(\tau_1, \dots, \tau_p)$ will be the one which includes the linear part p times and the p th-order term one time. Based on this remark, we can estimate the linear and the p th-order kernel as follows:

$$\begin{aligned} h_1(\tau) &= \tilde{c}_y^{(1)}(\tau, q_p) / \tilde{c}_y^{(1)}(q_1, q_p) \\ h_p(\tau) &= \tilde{c}_y^{(1)}(q_1, \tau) / \left(h_1(0) \text{cum}\{\underbrace{x(n), \dots, x(n)}_{p \text{ copies of } x(n)}, x^{*p}(n)\} \right). \end{aligned}$$

So far we have estimated the first and last kernel of the passband Hammerstein channel by combining the techniques in [11, 10]. The kernels sandwiched between the linear kernel and the p th-order term can be obtained through the following recursion.

Theorem 2 *Consider a passband Hammerstein model. For $2 \leq k \leq p$, the following equation holds:*

$$\tilde{c}_y^{(1)}(q_1, \tau) = \sum_{\mu=0}^{p-k} \text{cum}\{\underbrace{x^{1+\mu}(n), x(n), \dots, x(n)}_{k \text{ copies of } x(n)}, x^{*(k+\mu)}(n)\} h_{1+\mu}(0) h_{k+\mu}(\tau).$$

The proof is supplied in [3]. Theorem 2 is based on the fact that the linear kernel and the p th-order kernel are identified first then the kernel of order $k = p - 1$ is calculated. This process is iterated until $k = 2$. The above technique is applicable to Hammerstein channels excited by PSK inputs of arbitrary order. However when the channel is excited by QAM inputs, the procedure is limited to quadratic Hammerstein channels.

5 Blind identification of second order Volterra systems with complex random inputs

In this section blind identification methods for second order Volterra systems excited by complex valued random variables are developed. The proposed blind identification method relies on output cumulants of order up to 4. The computation of these cumulants and the resulting expressions are provided in [4]. If

Table 1. Algorithms for blind Volterra identification

Algorithm 1 ($q_1 > q_2$)	Algorithm 2 ($q_1 = q_2 = q$)
Require : $h_1(0) = 1$	Require : $h_1(0) = 1, h_2(0, 0) = 0$
1: $\gamma_{4,0} = \frac{c_{y(4)}^{(0)}(q_1, 0, 0)^2}{c_{y(4)}^{(0)}(q_1, q_1, 0)}$	1: $\gamma_{4,0} = \frac{c_{y(4)}^{(0)}(q, q, 0)^3}{c_{y(4)}^{(0)}(q, q, q)^2}$
2: $\gamma_{1,1} = \frac{c_{y(1)}^{(1)}(q_1) c_{y(4)}^{(0)}(q_1, 0, 0)}{c_{y(4)}^{(0)}(q_1, q_1, 0)}$	2: $\gamma_{1,1} = \frac{c_{y(1)}^{(1)}(q) c_{y(4)}^{(0)}(q, q, 0)}{c_{y(4)}^{(0)}(q, q, q)}$
3: $h_1(\tau) = \frac{c_{y(4)}^{(0)}(q_1, \tau, 0)}{c_{y(4)}^{(0)}(q_1, 0, 0)},$ $c_{y(3)}^{(0)}(q_1, q_1) = h_2(0, 0) \gamma_{4,0} h_1^2(q_1)$	3: $h_1(\tau) = \frac{c_{y(4)}^{(0)}(q, q, \tau)}{c_{y(4)}^{(0)}(q, q, q)}, h_2(q, q) = \frac{c_{y(3)}^{(0)}(q, q)}{2\gamma_{4,0} h_1(q)}$
4: $h_2(\tau, \tau) = \frac{c_{y(3)}^{(0)}(q_1, \tau) - \gamma_{4,0} h_2(0, 0) h_1(q_1) h_1(\tau)}{\gamma_{4,0} h_1(q_1)}$	4: $h_2(\tau, \tau) = \frac{c_{y(3)}^{(0)}(q, \tau) - \gamma_{4,0} h_2(q, q) h_1(\tau)}{\gamma_{4,0} h_1(q)}$
5: for $h_2(\tau_1, \tau_2)$ use Eq. below with $q = q_1$	5: for $h_2(\tau_1, \tau_2)$ use Eq. below
$h_2^*(\tau_1, \tau_2) = \frac{1}{2\gamma_{1,1}^2 h_1^2(q)} \left[c_{y(2)}^{(1)}(q - \tau_1, q - \tau_2) - \sum_{i=0}^{\tau_1 - \lfloor \frac{\tau_2}{\tau_2} \rfloor} \sum_{j=0}^{\tau_2} \tilde{h}_2(i, j) h_1(i + q - \tau_1) h_1(j + q - \tau_2) \right]$	

these expressions are evaluated at suitably chosen lags, sparse equations with respect to the Volterra kernels result. To proceed with the Volterra kernel identification algorithms, we distinguish two different cases: $q_1 > q_2$ and $q_1 = q_2 = q$. Appropriate normalization constraints are imposed for each case.

The proposed methods involve four steps carried out in the following sequence:

1. Compute the linear kernel using fourth order cumulants and a q-slice formula
2. Compute the input intensities $\gamma_{4,0}$ and $\gamma_{1,1}$
3. Compute the diagonal elements of the second order kernel using third order cumulants and a q-slice formula
4. Compute the off-diagonal elements of the second order kernel using third order cumulants and linear system solvers.

The above four steps are implemented by the algorithms of Table 1. Note that all relevant expressions are exact and hence the Volterra system is uniquely identifiable.

6 Sparse adaptive ℓ_1 -regularized algorithm

In this section, we propose a family of sparse adaptive ℓ_1 -regularized algorithms that can be used for sparse parameter estimation. The derived family of sparse adaptive algorithms is based on the Expectation Maximization (EM) framework. Let us start by considering a model that captures the dynamics of the unknown parameter vector $\mathbf{h}(n)$ (at time n). A popular technique in the adaptive filtering

Table 2. EM-KALMAN filter for sparse adaptive tracking

Algorithm description
Initialization : $\mathbf{h}_0 = \bar{\mathbf{h}}_0$, $\mathbf{P}_0 = \delta^{-1} \mathbf{I}$ with $\delta = \text{const.}$
For $n := 1, 2, \dots$ do
1: $\mathbf{k}(n) = \frac{\mathbf{P}(n-1)\mathbf{x}^*(n)}{\sigma_\eta^2 + \mathbf{x}^T(n)\mathbf{P}(n-1)\mathbf{x}^*(n)}$
2: $\boldsymbol{\psi}(n) = \mathbf{h}(n-1) + \mathbf{k}(n)\varepsilon(n)$
3: $\mathbf{P}(n) = \mathbf{P}(n-1) + r_n \mathbf{I} - \mathbf{k}(n)\mathbf{x}^T(n)\mathbf{P}(n-1)$
4: $\mathbf{h}(n) = \text{sgn}(\boldsymbol{\psi}(n)) \left[\boldsymbol{\psi}(n) - \gamma(\sigma_{\boldsymbol{\psi}_{n-1}}^2 + r_n) \mathbf{I} \right]_+$
end For

literature is to describe parameter dynamics by the first-order model [12]

$$\mathbf{h}(n) = \mathbf{h}(n-1) + \mathbf{q}_{|\Lambda_0}(n) = \mathbf{h}_0 + \sum_{i=1}^n \mathbf{q}_{|\Lambda_0}(i); \quad \mathbf{h}_0 \sim \mathcal{N}(\bar{\mathbf{h}}_0, \sigma_0^2 \mathbf{I}_{|\Lambda_0}) \quad (4)$$

where Λ_0 denotes the true support set of \mathbf{h}_0 , i.e. the set of the non-zero coefficients. The noise term $\mathbf{q}(n)$ is zero outside $|\Lambda_0$ and zero-mean Gaussian inside $|\Lambda_0$ with diagonal covariance matrix $\mathbf{R}_{|\Lambda_0}(n) = \text{diag}(\sigma_{q_1}^2(n), \dots, \sigma_{q_d}^2(n))$, where d is the ℓ_0 -norm of \mathbf{h}_0 . The variances $\{\sigma_{q_i}^2(n)\}_{i=1}^d$ are in general allowed to vary with time. The stochastic processes $\boldsymbol{\eta}(n)$, $\mathbf{q}(n)$ and the random variable \mathbf{h}_0 are mutually independent.

To apply the Expectation-Maximization method we have to specify the complete and incomplete data. The vector $\mathbf{h}(n)$ at time n is taken to represent the complete data vector, whereas $\mathbf{y}(n-1)$ accounts for the incomplete data [13]. In this context the conditional density $p(\mathbf{h}(n)|\mathbf{y}(n-1))$ plays a major role. This density is Gaussian with mean $\boldsymbol{\psi}(n) = \mathbb{E}[\mathbf{h}(n)|\mathbf{y}(n-1)]$. Under broad conditions the maximizer of the incomplete likelihood is obtained by maximizing the complete likelihood function through successive application of the following two steps:

E-step : computes the conditional expectation

$$Q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}(n-1)) = \mathbb{E}_{p(\mathbf{h}(n)|\mathbf{y}(n-1); \hat{\boldsymbol{\theta}}(n-1))} [\log p(\mathbf{h}(n); \boldsymbol{\theta})]$$

M-step : maximizes the Q -function minus the ℓ_1 -penalty with respect to $\boldsymbol{\theta}$

$$\hat{\boldsymbol{\theta}}(n) = \arg \max_{\boldsymbol{\theta}} \left\{ Q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}(n-1)) - \gamma \|\boldsymbol{\theta}\|_{\ell_1} \right\}$$

Note that $p(\mathbf{h}(n); \boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\psi}_n(\boldsymbol{\theta}), \boldsymbol{\Sigma}(n))$ and hence the Q -function takes the form

$$Q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}_{n-1}) = \text{const.} + \boldsymbol{\theta}^H \boldsymbol{\Sigma}^{-1}(n) \boldsymbol{\psi}(n) - \frac{1}{2} \boldsymbol{\theta}^H \boldsymbol{\Sigma}^{-1}(n) \boldsymbol{\theta} \quad (5)$$

where the constant incorporates all terms that do not involve $\boldsymbol{\theta}$ and hence do not affect the maximization.

The parameter $\psi(n)$ is recursively computed by the Kalman filter [12], see Table 2 steps 1 – 3, which in the special case of the time-varying random walk model Eq. (4) takes an RLS type appearance. Note that $\varepsilon(n)$, in Table 2, denotes the prediction error given by $\varepsilon(n) = y(n) - \mathbf{x}^T(n)\mathbf{h}(n-1)$.

Maximization of the Q function leads to the *soft thresholding* function, see Table 2 step 4. This operation shrinks coefficients above the threshold in magnitude value. The complete algorithm is presented in Table 2.

7 Sparse Adaptive Orthogonal Matching Pursuit algorithm

This section converts a powerful greedy scheme developed in [14] into an adaptive algorithm. Greedy algorithms form an essential tool for sparse parameter estimation. However, their inherent batch mode discourages their use in time-varying environments due to significant complexity and storage requirements.

The proposed algorithm relies on three modifications to the CoSaMP structure [14]: the proxy identification, estimation, and error residual update. The error residual is now evaluated by

$$v(n) = y(n) - \mathbf{x}^T(n)\mathbf{h}(n). \quad (6)$$

The above formula involves the current sample only, in contrast to the CoSaMP scheme which requires all the previous samples. Eq. (6) requires s complex multiplications, whereas the cost of the sample update in the CoSaMP is sn multiplications. A new proxy signal that is more suitable for the adaptive mode, can be defined as:

$$\mathbf{p}(n) = \sum_{i=1}^{n-1} \mathbf{x}^*(i)v(i) \quad (7)$$

and is updated by $\mathbf{p}(n) = \mathbf{p}(n-1) + \mathbf{x}^*(n-1)v(n-1)$. The last modification attacks the estimation step. The vector $\mathbf{h}(n)$ is updated by standard adaptive algorithms such as the LMS and RLS.

LMS is one of the most widely used algorithm in adaptive filtering due to its simplicity, robustness and low complexity. Hence, for reasons of simplicity and complexity we focus on the LMS algorithm. At each iteration the current regressor $\mathbf{h}(n)$ and the previous estimate $\mathbf{w}(n-1)$ are restricted to the instantaneous support originated from the support merging step. The resulting algorithm is presented in Table 3, where $\mathbf{h}_{|A}$ and $\mathbf{w}_{|A}$ denote the sub-vectors corresponding to the index set A , $\max(|a|, s)$ returns s indices of the largest elements of a and A^c represents the complement of set A . The following Theorem establishes the steady state Mean Square Error (MSE) error performance of the SpAdOMP algorithm:

Table 3. SpAdOMP Algorithm

Algorithm description		Complexity
$\mathbf{h}(0) = 0, \mathbf{w}(0) = 0, \mathbf{p}(0) = 0$	{Initilization}	
$v(0) = y(0)$	{Initial residual}	
$0 < \lambda \leq 1$	{Forgetting factor}	
$0 < \mu < 2\lambda_{\max}^{-1}$	{Step size}	
For $n := 1, 2, \dots$ do		
1: $\mathbf{p}(n) = \lambda \mathbf{p}(n-1) + \mathbf{x}^*(n-1)v(n-1)$	{Form signal proxy}	M
2: $\Omega = \text{supp}(\mathbf{p}_{2s}(n))$	{Identify large components}	M
3: $\Lambda = \Omega \cup \text{supp}(\mathbf{h}(n-1))$	{Merge supports}	s
4: $\varepsilon(n) = y(n) - \mathbf{x}_{ \Lambda}^T(n)\mathbf{w}_{ \Lambda}(n-1)$	{Prediction error}	s
5: $\mathbf{w}_{ \Lambda}(n) = \mathbf{w}_{ \Lambda}(n-1) + \mu \mathbf{x}_{ \Lambda}^*(n)\varepsilon(n)$	{LMS iteration}	s
6: $\Lambda_s = \max(\mathbf{w}_{ \Lambda}(n) , s)$	{Obtain the pruned support}	s
7: $\mathbf{h}_{ \Lambda_s}(n) = \mathbf{w}_{ \Lambda_s}(n), \mathbf{h}_{ \Lambda_s^c}(n) = \mathbf{0}$	{Prune the LMS estimates}	
8: $v(n) = y(n) - \mathbf{x}^T(n)\mathbf{h}(n)$	{Update error residual}	s
end For		$\mathcal{O}(M)$

Theorem 3 (*SpAdOMP*)¹. The proposed algorithm, for large n , produces an s -sparse approximation $\mathbf{h}(n)$ that satisfies the following steady-state error bound

$$\|\mathbf{h} - \mathbf{h}(n)\|_{\ell_2} \lesssim C_1(n)\|\boldsymbol{\eta}(n)\|_{\ell_2} + C_2(n)\|\mathbf{x}_{|\Lambda}(n)\|_{\ell_2}|e_o(n)|,$$

where $e_o(n)$ is the estimation error of the optimum Wiener filter and $C_1(n)$, $C_2(n)$ are constants independent of \mathbf{h} (given explicitly in [6]) and are only functions of the restricted isometry constants, λ_{\min} (the minimum eigenvalue of the input covariance matrix) and the step-size μ .

8 Blind Identification of sparse channels via the EM algorithm

The purpose of this section is to develop a blind identification algorithm for the estimation of sparse channels, under the assumption that the transmitted symbols are i.i.d. and take values in a finite alphabet set. A batch algorithm for blind channel estimation is reported in [15] using the iterative nature of the Expectation Maximization (EM) algorithm. We propose exploitation of the sparse nature of the channel by regularizing the cost function of the blind EM algorithm via the use of the ℓ_1 norm constraint.

In blind identification, the EM algorithm can be used to iteratively maximize $\log p(\mathbf{y}(n); \boldsymbol{\theta})$ (where $\boldsymbol{\theta} = \bar{\mathbf{h}}$), without explicitly computing it. To use the EM algorithm, we consider the observations $\mathbf{y}(n)$ as the incomplete data and

¹ Proof is omitted due to space limitations.

Table 4. The Sparse BW Algorithm

Algorithm description	
$\alpha_1(i) = \pi_i b_i(\mathbf{y}_1), i := 1, \dots, M^L, \quad \beta_N(i) = 1, i := 1, \dots, M^L, \quad \sigma_\eta^2 = 1$	{Initiliazation}
For $\ell := 0, 1, \dots$, do	
1: $\alpha_{n+1}(j) = \sum_{i=1}^{M^L} \alpha_n(i) p_{ij} b_j(\mathbf{y}_{n+1}), \quad n := 1, \dots, N-1, \quad j := 1, \dots, M^L$	{Forward Recursion}
2: $\beta_n(i) = \sum_{j=1}^{M^L} \beta_{n+1}(j) p_{ij} b_j(\mathbf{y}_{n+1}), \quad n := N-1, \dots, 1, \quad i := 1, \dots, M^L$	{Backward Recursion}
3: $\gamma_n(i \boldsymbol{\theta}^{(\ell)}) = \frac{\alpha_n(i)\beta_n(i)}{\sum_{j=1}^{M^L} \alpha_n(j)\beta_n(j)}, \quad n := 1, \dots, N, \quad i := 1, \dots, M^L$	{Posterior probabilities}
4: $\mathbf{h}_i^{(\ell+1)} = \frac{\text{sgn}(\mathbf{r}_i^{(\ell)})}{\mathbf{R}_{i,i}^{(\ell)}} \left[\mathbf{r}_i^{(\ell)} - \tau \right]_+$	{Channel estimation}
5: $\sigma_\eta^{(\ell+1)2} = (N+1)^{-1} \sum_{n=1}^N y_n - \hat{\mathbf{x}}_n^{(\ell)T} \mathbf{h}^{(\ell)} ^2$	{Noise Variance est.}
end For	

$(\mathbf{y}(n), \mathbf{X}(n))$ as the complete data. The EM algorithm is a two-step iterative procedure which under mild conditions converges to a local maximum. The proposed variant of the EM algorithm for blind sparse channel estimation iterates between the following two steps until convergence is reached:

- 1) **E-step:** Compute $Q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(\ell)}) = \mathbb{E}\{\log p(\mathbf{y}(n), \mathbf{X}(n)|\boldsymbol{\theta})|\mathbf{y}(n); \hat{\boldsymbol{\theta}}^{(\ell)}\}$
- 2) **M-step:** Solve $\hat{\boldsymbol{\theta}}^{(\ell+1)} = \arg \max_{\boldsymbol{\theta}} \{Q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(\ell)}) - 2\tau \|\boldsymbol{\theta}\|_{\ell_1}\}$.

The E-step is a symbol detector and is carried out by the forward-backward recursions of Table 4 steps 1-2. Maximization of the penalized Q -function with respect to \mathbf{h} at the M-step, has a closed form expression to each component of $\mathbf{h}^{(\ell+1)}$ and is given by the *soft-thresholding* function, see Table 4 step 4. The method outlined above is summarized in Table 4.

9 Conclusions

This dissertation aimed to develop new methods and algorithms for the estimation of nonlinear communications systems that are modeled by Volterra series. The estimation is achieved either with the help of a training signal or by blind identification methods. Moreover, the estimation performance depends on the pattern of channel (sparse or dense). Thus, the developed methods take into account the pattern of the channel, and simply estimate those parameters that actually contribute to the output

References

1. M. Schetzen, *The Volterra and Wiener Theories of Nonlinear Systems*, Wiley, 1980.
2. G. Mileounis, P. Koukoulas, N. Kalouptsidis, Input-output identification of nonlinear channels using PSK, QAM and OFDM inputs, *Signal Process.* (Elsevier) 89 (2009) 1359–1369.
3. G. Mileounis, N. Kalouptsidis, P. Koukoulas, Blind identification of Hammerstein channels using QAM, PSK and OFDM input using higher order cumulants, *IEEE Trans. Commun.* 57 (12) (2009) 3652–3661.
4. G. Mileounis, N. Kalouptsidis, Blind identification of second order Volterra systems with complex random inputs using higher order cumulants, *IEEE Trans. Signal Process.* 57 (10) (2009) 4129–4135.
5. N. Kalouptsidis, G. Mileounis, B. Babadi, V. Tarokh, Adaptive algorithms for sparse system identification, accepted in *Signal Process.* (Elsevier), subject to revisions.
6. G. Mileounis, B. Babadi, N. Kalouptsidis, V. Tarokh, An adaptive greedy algorithm with application to nonlinear communications, *IEEE Trans. Signal Process.* 58 (6) (2010) 2998–3007.
7. G. Mileounis, N. Kalouptsidis, B. Babadi, V. Tarokh, Blind identification of sparse channels and symbol detection via the EM algorithm, submitted in *IEEE Signal Process. Lett.*
8. S. Benedetto, E. Biglieri, *Principles of Digital Transmission: with wireless applications*, Springer, 1999.
9. S. Yasui, Stochastic functional fourier series, Volterra series, and nonlinear system analysis, *IEEE Trans. Autom. Control* 24 (2) (1979) 230–242.
10. S. Prakriya, D. Hatzinakos, Blind identification of LTI-ZMNL-LTI nonlinear channel models, *IEEE Trans. Signal Process.* 43 (12) (1995) 3007–3013.
11. N. Kalouptsidis, P. Koukoulas, Blind identification of Volterra-Hammerstein systems, *IEEE Trans. Signal Process.* 53 (8) (2005) 2777–2787.
12. L. Ljung, General structure of adaptive algorithms: Adaptation and tracking, in *Adaptive System Identification and Signal Processing Algorithms*, Eds. N. Kalouptsidis and S. Theodoridis, 1993.
13. M. Feder, *Statistical signal processing using a class of iterative estimation algorithms*, Ph.D. thesis, M.I.T., Cambridge, MA (1987).
14. D. Needell, J. Tropp, CoSaMP: Iterative signal recovery from incomplete and inaccurate samples, *Appl. Comput. Harmon. Anal.* 26 (2009) 301–321.
15. G. Kaleh, R. Vallet, Joint parameter estimation and symbol detection for linear or nonlinear unknown channels, *IEEE Trans. Commun.* 42 (7) (1994) 2606–2413.

Adaptive Educational Hypermedia Systems on the Web for the Didactics of Science and Technology

Alexandros Papadimitriou*

Department of Informatics and Telecommunications
National and Kapodistrian University of Athens
alexandr@di.uoa.gr

Abstract: This doctoral dissertation deals with the design and implementation issues of an Adaptive Educational Hypermedia System (AEHS).

More specifically, it deals with: (i) the design of an educational framework based on constructivism that it will guide the AEHS's educational decisions and the design of domain model. Also, it will set the goals and the functionality of adaptation, the feedback, the assessment, the participation of student in collaborative problem solving activities, and it also will determine the combination of adaptation techniques, taking into consideration the content and making the best of the learner's specific characteristics, (ii) the selection of the learning styles model, (iii) issues of sharing the control between the learner and the system, (iv) the selection of the appropriate adaptive navigation techniques, (v) a learner-centered method of meta-adaptive navigation, (vi) the interactive problem solving support through the activities, (vii) a learner-centered method of adaptive group formation, and (viii) issues of adaptive presentation.

In the frame of this dissertation, we designed and implemented the AEHS MATHEMA that combines the constructivist, the socio-cultural and the meta- cognitive didactic models, and it also supports the individual and collaborative learning. The didactic strategies supported by the MATHEMA based on Kolb's learning styles model. The general aim of the MATHEMA is to support senior high school learners or novices of higher education, through an interactive and constructivist educational material, in learning science individually and/or collaboratively, and overcoming their possible misconceptions and learning difficulties. At the current version, appropriate educational material has been developed for the learning of electromagnetism.

The adaptive and intelligent techniques supported by the MATHEMA are the following: *curriculum sequencing, adaptive presentation, adaptive and meta-adaptive navigation, interactive problem solving support, and adaptive group formation.*

At the end of this dissertation, researches for the development and the formative assessment of the MATHEMA, and the conclusions are presented.

Researches indicated that the senior high school students increase their performance by studying the MATHEMA. Also, the formative assessment of the MATHEMA indicated that almost all its functions are useful and user- friendly.

Dissertation Advisor: P. Georgiadis, Professor.

Dissertation Summary

This dissertation is referred to:

- (i) the design of an educational framework based on constructivism that it will guide the AEHS's educational decisions and the design of the domain model. Moreover, it will set the goals and the functionality of adaptation, the feedback, the assessment, the participation of student in collaborative problem solving activities, and it also will determine the combination of adaptation techniques, taking into consideration the content and making the best of the learner's specific characteristics,
- (ii) the selection of the learning styles model, proper for the specific application, among of all those that have been proposed by various psychologists as well as the design of adaptation on the basis of learning style model,
- (iii) issues concerning the effective design of the learners' involvement in the educational process and issues of sharing the control between the learner and the system, in a clear way, taking into consideration the learner's needs and current state of his/her model,
- (iv) the selection of the appropriate adaptive navigation techniques, so that the learner to be assisted during his/her navigation according to his/her Web experience and knowledge level of the current goal,
- (v) a learner-centered method of meta-adaptive navigation, so that the learner to be assisted in his/her selection of the most appropriate navigation technique suited to his/her profile,
- (vi) the interactive problem solving support through the activities based on modern approaches in teaching combining individual and collaborative learning,
- (vii) a learner-controlled method of adaptive group formation so that the learner to have the ability in selection of the appropriate collaborator among them that the system suggests to him/her by a priority list taking into account his/her learning style and learning style and knowledge level of his/her collaborators,
- (viii) issues of adaptive presentation.

b) Innovations of dissertation

According to the framework we mentioned above, we designed the Web- based AEHS MATHEMA which supports appropriate adaptive and intelligent techniques (curriculum sequencing, adaptive presentation, adaptive and meta-adaptive navigation, interactive problem solving, and adaptive group formation).

From these techniques, the innovative ones are the following:

- the adaptive navigation techniques (direct guidance, adaptive links hiding, adaptive link annotation, and adaptive link sorting) which assist the student in his/her navigation according to his/her Web experience and knowledge level of the current goal,
- the meta-adaptive navigation technique in which the student is assisted in his/her selection of the most appropriate navigation technique suited to his/her profile,
- the interactive problem solving support through the activities based on modern approaches in teaching and combining individual and collaborative learning,
- a learner-controlled method of adaptive group formation in which the learner have the ability in selection of the appropriate collaborator among them that the system suggests to him/her by a priority list taking into account the learner

learning style and learning style and knowledge level of his/her collaborators.

Related Work

i) Interactive problem solving support

ActiveMath is an intelligent learning environment on the Web. It provides high-quality Web presentations of mathematical documents, intelligent selection of content-items to achieve learning goals, search for text and mathematical objects, copy and paste of formula, and interactive exercises with learner inputs evaluated by classical computer algebra systems. *ActiveMath* design aims at supporting truly interactive, exploratory learning and assumes the student to be responsible for her learning to some extent. Therefore, a relative freedom for navigating through a course and for learning choices is given to the students.

ELM-ART II was designed for learning programming in LISP and integrates a LISP compiler. *ELM-ART II* provides a unique example of *example-based problem solving support*. *ELM-ART II* contains “live examples” and short programming problems. In *ELM-ART II* if learners fail to solve a LISP programming problem, they can ask the system to diagnose the code of their solution and give detailed explanation of error. It also helps learners find the relevant examples from their previous experience by presenting an ordered list of examples based on their relevancy.

ii) Adaptive group formation

Students' personal features are taken into account by many researchers for forming student groups. Gogoulou et al. (2007) presented the *OmadoGenesis* tool that accommodates learners' characteristics (gender, ethnic background, motivations, attitudes, interests, etc.) in the formation of pure homogeneous, pure heterogeneous or mixed groups; that is, groups that satisfy heterogeneity for specific learners' characteristic and homogeneity for another characteristic. The *OmadoGenesis* tool implements three algorithms: one for pure homogeneous groups, one for pure heterogeneous groups and one based on the concept genetic algorithms for homogeneous, heterogeneous and mixed groups. Muehlenbrock (2006) combines information from learner profiles and information on the learner context. This combination has a potential of improving the quality of the grouping. It allows for the ad-hoc creation of learning groups, which is especially useful for peer help for immediate problems, by reducing the risk of disruptions. It also leverages the forming of face-to-face learning groups based on the presence information. The context sensing has been tested with a set of experiments, and a distributed application has been developed that helps teachers form groups. Potentially, other context information can be used to improve the group formation, such as agenda information from personal calendars, or the availability of preferred communication channels. The building of learning groups could also be enriched by information available on the experience from past collaborations, which could be provided by peers but also from a teacher if available. Furthermore, in addition to the topic of the collaboration, the group formation could include information on the type of support needed, among others.

Some researchers form students' groups according to students' learning styles. Papanikolaou et al. (2006) form students' groups based on their learning style by using the Honey and Mumford's learning style categorization and the visual/verbal dimension of Felder-Silverman model aiming at the collaboration in the formation of concept maps.

Examples in adaptive group formation and/or peer help include forming a group for collaborative problem solving (Ikeda et al., 1997), for finding the most competent peer to answer a question (McCalla et al., 1997) or for self-, peer- and collaborative-assessment process (Gouli et al., 2006). Ikeda et al suggest the *Opportunistic Group Formation* to form collaborative learning group dynamically and context-dependently. When the system detects the situation for a learner to shift from individual learning mode to collaborative learning mode, it forms a learning group each of whose members is assigned a reasonable learning goal and a social role which are consistent with the goal for the whole group. McCalla et al presented a practical approach for just-in- time workplace training that uses artificial intelligence techniques to extend informal peer help networks to a broader scale. They have developed a system called *PHelpS* that is a situated, peer-supported, AI-based approach to training in procedural, task-oriented domains. PHelpS supports workers as they perform their tasks, offers assistance in finding peer helpers when required, and mediates communication on task-related topics. Gouli et al developed the *PECASSE* system, which implements self-, peer- and collaborative-assessment in a Web-based educational setting by offering facilities for group formation, collaboration of learners, activity submission, review process, assignment of assessors, revision of the activity, and evaluation of assessors. Peer assessment refers to these activities of learners in which they judge and evaluate the work and/or the performance of their peers.

Until recently, most support for group formation in CSCL systems was based on the learner profile information such as gender, class and other features of the learner. In traditional classrooms, the teachers group students in work teams, but in CSCL systems, group formation can be performed either by the teacher (in classroom or using the information stored in the system) or automatically by the system (Carro et al., 2003). Carro et al use adaptation techniques to dynamically generate adaptive collaborative Web-based courses in their *TANGOW* system. These courses are generated at runtime by selecting, at every step and for each student, the most suitable collaborative tasks to be proposed, the time at which they are presented, the specific problems to be solved, the most suitable partners to cooperate with and the collaborative tools to support the group cooperation. This selection is based on the users' personal features, preferences, knowledge and behaviour while interacting with the course.

If the group formation is done by the system, it can be done randomly or by taking into account personal features included in the user and group models (Read et al., 2006). Read et al designed and implemented the *COPPER* system, where individual and collaborative learning are combined within a constructivist approach to facilitate second language learning. The adaptive group formation algorithm dynamically generates communicative groups based on the linguistic capabilities of available students, and a collection of collaborative activity templates. Students initially work individually on certain linguistic concepts, and subsequently participate in authentic collaborative communicative activities.

In some AEH systems students are grouped according to their learning styles. In Martin & Paredes (2004) system, the default criteria for group formation consist of combining active students with reflective ones of Felder-Soloman model in similar percentages. They studied the impact of learning styles and group homogeneity/heterogeneity on the results obtained by students in collaborative tasks. Thus, they concluded that some dimensions of the learning style model, seem to affect the quality of the resulting work.

Results and Discussion

In order to implement the AEHS MATHEMA we adopted the constructivist, socio-cultural and meta-cognitive learning model. The didactic design of the MATHEMA supports the students:

- in construction of their knowledge,
- in recognition of their misconceptions and correction of their errors through reflection,
- in selection and achievement of their learning goals, recognizing what they have already learned and what they be able to do by evaluating their progress of learning by themselves,
- by providing them the appropriate didactic strategies suited to their learning style,
- in development of their critical thinking,
- in their self-regulation,
- in collaboration through collaborative activities,
- by giving them motivations for their participation, and
- by giving them multiple representations for the improvement of their learning.

The didactic approaches adopted by the MATHEMA are the following: questioning-visualization, exercises solving, theory and examples, and problem solving activities. The adaptive presentation is done according to students' learning style.

For the design and implementation of the innovative techniques in AEHS MATHEMA, we studied the literature and we concluded to the following:

a) adaptive navigation techniques

There are a lot of AEHS. Some of them use one only adaptive technique (e.g., Knowledge Sea), while some other use more than one adaptive technique (e.g., ISIS-Tutor, ELM-ART). Nothing of them suggests to learner any technique to begin his/her navigation according to his/her Web experience and knowledge level of the current goal. After an extended review of the literature, we found interesting researches concerning the navigation difficulties that the user deal with on the Web. Also, a lot of researches suggest the most appropriate navigation techniques for students who have nothing, little or enough Web experience or nothing, little or enough knowledge experience on the current goal.

Taking into consideration researches suggesting that not all of navigation techniques are appropriate for all learners and that the most significant role for the selection of navigation technique plays the Web experience and the knowledge level of the current goal of the learner, we decided to design the AEHS MATHEMA so that to offer four navigation techniques (direct guidance, adaptive link annotation, adaptive link hiding, and link sorting) to assist the learner during his/her navigation according to his/her Web experience and knowledge level of his/her current goal.

In AEHS MATHEMA, the first time that the learner logs in the system, is called to declare his/her Web experience and his/her knowledge level of the selected current goal. After the declaration, the system suggests for him/her the most appropriate navigation technique, according to his/her Web experience and his/her knowledge level of the selected current goal. The system informs the learner about the navigation technique that it suggests for him/her and it also explains the reasons why it suggests the particular navigation technique (e.g., to protect him from navigation problems). In addition, the system suggests to the learner to change the suggested navigation technique when he/she has fulfilled the terms of meta-adaptation.

b) meta-adaptive navigation

All the suggestions for meta-adaptation in literature are intended to the meta-adaptation based on the system (system-controlled) where it decides the most appropriate navigation technique, for each learner and environment, and acts respectively.

So far, nothing AEH system has developed using a kind of meta-adaptation. Meta-adaptation in the MATHEMA is based on the learner (learner-controlled), where the system assists the learner make decision for the navigation technique suited better to him/her by presenting the advantages and disadvantages of navigation techniques supported by it, when the learner has fulfilled the terms of meta-adaptation. This proposal enhances the self-regulating learning offered by hypermedia.

The meta-adaptation engine of the MATHEMA, after the learner's succeeded assessments on n main concepts and considering that he/she has obtained the adequate Web experience, appears a window on the screen with information about the advantages and disadvantages of each offered navigation technique. Then, the learner decides if he/she will select another navigation technique.

c) Interactive problem solving support

After an extended study of literature, we found AEHS supporting this intelligent technique as are the ISIS-Tutor, ELM-ART and ActiveMath. From the study of these systems arises that they support the learner in:

- problem solving, errors correction and conceptual change,
- the development of his/her critical thinking,
- his/her reflection.

In addition, from the study of these systems arises that they do not support the learner in:

- constructivist type problem solving activities making the most of the modern didactic approaches,
- both individual and collaborative learning.

The interactive problem solving activities in the MATHEMA make use of the following didactic approaches: experimentation through simulations, explorations, guided discovery and collaboration.

d) adaptive group formation

The most of the systems in literature use several characteristics of learners to support adaptive group formation, and they implement it based on a system-controlled design. That is, the system decides the group formation and the learners are informed the group that the system includes them without having the possibility to change it. An exception is the tool of Christodouloupoulos & Papanikolaou (2007) in which the group formation is done as follows: The system categorizes the learners into groups and then it allows them to communicate with the educator for the purpose of negotiating of their group.

In the same research line with the systems that we study, the MATHEMA offers the adaptive group formation technique by supporting the learner to select the most appropriate learner by a priority list. For the adaptation the system makes use of the Kolb's learning styles model and knowledge level of learners on the current goal.

In conclusion, the above-mentioned techniques are innovative and upgrade the science.

Conclusions

In the frame of the design of the MATHEMA, we conducted researches with high school students and the results of these researches helped us to implement the innovative techniques of the MATHEMA.

The first research had as main purpose the investigation if the suggested problem solving method facilitates the learners to deal with their misconceptions and learning misunderstanding.

The results of this research indicated that the suggested problem solving method assists the participants:

- to compute physical quantities and to correct their errors,
- to comprehend their misconceptions or misunderstandings and revise their points of view,
- to be able to explain their choices, and
- 1 to accept the restrictions of the formula usage.

According to the above results, we designed the next stage of our research through an interactive educational material taking into account the learning styles of the learners.

The second research had as main purpose the investigation whether the MATHEMA assists the learners in improvement of their performance when they collaborate in solving problems of electromagnetism.

The research questions were the following:

- Do the learners improve their performance when they study through the MATHEMA?
- Do the learners improve their performance when they carry out the problem solving activity?
- Are there differences among concrete-concrete, abstract-abstract and concrete-abstract groups when they carry out problem solving activities?

The conclusions arose from the second research are the following:

- the performance of participants is significantly improved when they study through the MATHEMA ($F_{1,22} = 49.120, p = 0.000$),
- the concrete-concrete and abstract-abstract groups almost equivalently performed in problem solving activities but both of them performed significantly better than

the concrete-abstract groups in problem solving activities,

Thus, in problem solving activities:

- (1). A Diverging (concrete) collaborates better with a Diverging (concrete) rather than an Assimilating (abstract).
- (2). An Assimilating (abstract) collaborates better with an Assimilating (abstract) rather than a Diverging (concrete).
- (3). A Converging (abstract) collaborates better with a Converging (abstract) rather than an Accommodating (concrete).
- (4). An Accommodating (concrete) collaborates better with an Accommodating (concrete) rather than a Converging (abstract).
- (5). A Diverging (concrete) collaborates better with an Accommodating (concrete) rather than a Converging (abstract).
- (6). An Assimilating (abstract) collaborates better with a Converging (abstract) rather than an Accommodating (concrete).
- (7). A Converging (abstract) collaborates better with an Assimilating (abstract) rather than a Diverging (concrete).
- (8). An Accommodating (concrete) collaborates better with a Diverging (concrete) rather than an Assimilating (abstract).

The results of the second research helped us in the implementation of the adaptive group formation technique offered by the MATHEMA.

In the frame of the formative assessment of the MATHEMA, we conducted a research in the Department of Informatics and Telecommunications of University of Athens in order to assess the functionality, usefulness and usability of the MATHEMA functions. The opinions of the participants (students) are that almost all of the functions offered by the MATHEMA are useful and user-friendly. Also, the same students consider that the MATHEMA is an enjoyable environment (88,4%), facilitates the users' attention (90,7%) and gives opportunities in learning science (90,7%).

References

1. **Papadimitriou, A., Grigoriadou, M., & Gyftodimos, G.** (2009). Interactive Problem Solving Support in the Adaptive Educational Hypermedia System MATHEMA, *IEEE Transactions on Learning Technologies*, vol. 2, no. 2, pp. 93-106.
2. **Papadimitriou, A., Grigoriadou, M., & Gyftodimos, G.** (2009). Learning Electromagnetism Through the MATHEMA. *The International Journal of Learning*, vol. 16, issue 6, pp. 371-390.
3. **Papadimitriou, A., Gyftodimos, G., & Grigoriadou, M.** (2010). The Learning Facilities and Adaptation Techniques of the MATHEMA. *The International Journal of Learning*, vol. 17, issue 1, pp. 155-172.
4. **Papadimitriou, A., Gyftodimos, G.** (2007). Use of Kolb's learning cycle through an adaptive educational hypermedia system for a constructivist approach of electromagnetism, In proceedings of 4th WSEAS/ IASME International Conference on Engineering Education, Crete Island, Greece, pp. 226-231.
5. **Papadimitriou, A., Grigoriadou, M., Gyftodimos, G.** (2008). Adaptive Group Formation and Interactive Problem Solving Support in the Adaptive Educational hypermedia System MATHEMA, In proceedings of 20th World Conference on Educational Multimedia, Hypermedia and

- Telecommunications (ED-MEDIA 2008), Vienna, Austria, pp. 2182-2191.
6. **Papadimitriou, A., Gyftodimos, G., Grigoriadou, M.** (2009). "MATHEMA: A Constructivist Environment for Electromagnetism Learning." Proceedings of 8th IEEE International Conference on Advanced Learning Technologies (ICALT '09), I. Aedo, *et al.*(Eds.), New Jersey: IEEE, 2009, pp. 453-454.
 7. **Papadimitriou, A., Gyftodimos, G., Grigoriadou, M.** (2009). "Adaptive Navigation Support and the Learner-Centered Meta-Adaptation in the Adaptive Educational Hypermedia System MATHEMA." Proceedings of 21th World Conference on Educational Multimedia, Hypermedia and Telecommunications (ED-MEDIA 2009), Hawaii, Honolulu. Chesapeake, VA: AACE. pp. 1453-1462.
 8. **Papadimitriou, A., Grigoriadou, M. & Gyftodimos, G.** (2010). MATHEMA: A Learner-Centered Design for Electromagnetism Learning. *In Proceedings of 22th World Conference on Educational Multimedia, Hypermedia and Telecommunications (ED-MEDIA 2010)*, Toronto, Canada, Chesapeake, VA: AACE. pp. 1249-1254.
 9. **Παπαδημητρίου, Α., Γυφτοδήμος, Γ., Γρηγοριάδου, Μ., Γλέζου, Κ.,** (2007), *Αντιμετώπιση των παρανοήσεων και των δυσκολιών που προκύπτουν από το φορμαλισμό του ηλεκτρομαγνητισμού, με τη χρήση προσομοιώσεων*, Ζ' Πανελλήνιο Συνέδριο ΠΤΔΕ Πανεπιστημίου Αθηνών με θέμα: «Έρευνα και εκπαίδευση στις Φυσικές Επιστήμες στα Παιδαγωγικά Τμήματα Δημοτικής Εκπαίδευσης» (υπό έκδοση).
 10. **Παπαδημητρίου, Α., Γυφτοδήμος, Γ., Γρηγοριάδου, Μ.** (2009). *Μια δραστηριότητα βασισμένη σε σύγχρονες διδακτικές προσεγγίσεις για μια εποικοδομιστική προσέγγιση του ηλεκτρομαγνητισμού*. Πρακτικά 6^{ου} Πανελλήνιου Συνεδρίου Διδακτικής Φυσικών Επιστημών και Νέων Τεχνολογιών στην Εκπαίδευση με θέμα «Οι πολλαπλές προσεγγίσεις της διδασκαλίας των Φυσικών Επιστημών», σ.σ. 648-655.
 11. **Παπαδημητρίου, Α., Γρηγοριάδου, Μ., Γυφτοδήμος, Γ.** (2009). *Διδακτικές στρατηγικές και εκπαιδευτικό υλικό του προσαρμοστικού εκπαιδευτικού συστήματος υπερμέσων MATHEMA*. Πρακτικά 5^{ου} Πανελλήνιου Συνεδρίου των Εκπαιδευτικών για τις ΤΠΕ στην Εκπαίδευση "Αξιοποίηση των Τεχνολογιών της Πληροφορίας και της Επικοινωνίας στη Διδακτική Πράξη", σ.σ. 732-741.
 12. **Παπαδημητρίου, Α., Γυφτοδήμος, Γ., Γρηγοριάδου, Μ.** (2009). *Μια μαθητοκεντρική σχεδίαση εποικοδομιστικού τύπου για την οικοδόμηση της γνώσης στον Ηλεκτρομαγνητισμό*. 1^ο Επιστημονικό συνέδριο Σ.Ε.Π. ΑΣΠΑΙΤΕ με θέμα: Η εκπαίδευση των εκπαιδευτικών της δευτεροβάθμιας επαγγελματικής και τεχνολογικής εκπαίδευσης στην Ελλάδα, Αθήνα, σ.σ. 536-545.

Performance Analysis of Wireless Single Input Multiple Output Systems (SIMO) in Correlated Weibull Fading Channels

Zafeiro G. Papadimitriou *

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications

`zefipap@hotmail.com`

November 29, 2010

Abstract

In this dissertation the statistical characteristics of the trivariate and quadrivariate Weibull fading distribution with arbitrary correlation, non-identical fading parameters and average powers are analytically studied. Novel expressions for important joint statistics are derived using the Weibull power transformation. These expressions are used to evaluate the performance of selection combining (SC) and maximal ratio combining (MRC) diversity receivers in the presence of such fading channels.

Multi-branch diversity, arbitrary correlation, Weibull fading.

I. Introduction

In recent years, the use of various telecommunication systems, their applications and usefulness to real - life has become significantly important. In

*Dissertation Advisor: Lazaros Merakos, Professor

today's telecommunications era, new technologies are constantly being developed since nowadays there is a continuous need for instant access to information from a geographical point to another.

In telecommunications systems, the term "communication channel", or simply "channel", refers to either a physical transmission medium such as a wire, or to wireless connection over a physical medium such as a radio channel. As far as wireless technologies are concerned, they are being used to meet many needs such as span a distance beyond the capabilities of typical cabling, link portable or temporary workstations, overcome situations where normal cabling is difficult and/or financially impractical to remotely connect mobile users or networks. As a consequence, there is a continuing demand for increased capacity and integration of the provided services over radio channels. The radio channel though is subject to multiple phenomena such as fading that degrade the telecommunications system performance. Hence, the determination and combat of fading effects constitute an important R&D topic which is also the subject of this PhD thesis.

One of the most important effect of the fading interference is the small scale fading which results in fluctuation of the received signal's amplitude, phase and angle of arrival. In order to combat this destructive effect, in this PhD thesis diversity reception techniques are being employed. According to this method, the receiver employs more than one antennae, in order to receive multiple copies of the transmitted signal. These copies are being appropriately combined in order to satisfy network administrator demands [1]. There are several diversity schemes, classified according to the combining technique employed at the receiver, the most well-known being selection combining (SC), maximal ratio combining (MRC) and equal gain combining (EGC). In diversity reception studies, it is frequently assumed that the different replicas of the same information signal are received over independent fading channels. However, in many practical wireless system applications, e.g. for small-size mobile units or indoor base stations, the receiving antennas are not sufficiently wide separated and thus the received and combined signals are correlated with each other. In order to model and analyze such realistic wireless channels with correlated fading it is mathematically convenient to use multivariate statistics [1], [2].

In the open technical literature there have been many papers published concerning multivariate distributions in relation to performance analysis of digital communication systems in the presence of correlated fading channels [3–9]. Most of these papers deal specifically with the so-called "constant"

and “exponential” correlation model. For the first one, correlation depends on the distance among the combining antennas and thus this model is more suitable for equidistant antennas [1, pp. 392]. The second one, corresponds to the scenario of multichannel reception from equispaced diversity antennas. This model has been widely used for performance analysis of space diversity techniques [3] or multiple-input multiple-output (MIMO) systems. The arbitrary correlation model [4], used in our paper, is the most generic correlation model available, since it allows for arbitrary correlation values between the receiving branches. Clearly it includes the constant and exponential correlation models as special cases.

Most of the published works concerning multivariate distributions with arbitrary correlation deal with Rayleigh and Nakagami- m fading channels [2, 4–6]. In [4], new infinite series representations for the joint probability density function (PDF) and the joint cumulative distribution function (CDF) of three and four arbitrarily correlated Rayleigh random variables have been presented. In [5], expressions for multivariate Rayleigh and exponential PDFs generated from correlated Gaussian random variables have been derived, as well as a general expression in terms of determinants for the multivariate exponential characteristic function (CF). In [2] useful closed-form expressions for the joint Nakagami- m multivariate PDF and CDF with arbitrary correlation, were derived and the correlation matrix was approximated by a Green’s matrix. In a recent paper [6], infinite series representations for the PDF, CDF and CF for the trivariate and quadrivariate Nakagami- m distribution have been presented.

The Weibull distribution [10], although originally used in reliability and failure data analysis, it has been recently considered as an appropriate distribution for modeling wireless communication channels [7–9]. The main motivation of this choice is its very good fit to experimental fading channel measurements for both indoor and outdoor terrestrial radio propagation environments. In [7] it was argued that the Weibull distribution could also been considered as a generic channel model for land-mobile satellite systems. Recently, expressions for the joint PDF, CDF and the moment-generating function (MGF) for the bivariate Weibull distribution have been presented [8]. In the same reference the multivariate Weibull distribution has also been studied for the exponential and constant correlation case considering equal average fading powers. In [9] a Green’s matrix approximation for the multivariate Weibull distribution with arbitrary correlation has been presented and an analytical expression for the joint CDF has been derived. However,

the performance analysis presented in [9] is restricted to SC receivers and is applicable only to the evaluation of outage probability (OP).

Motivated by the above, in this dissertation, we present a detailed and thorough analytical study of the statistical characteristics of the arbitrary correlated trivariate and quadrivariate Weibull fading distributions and their applications to various diversity receivers. For both distributions we consider the most general correlation model available, namely the arbitrary correlation model, with non-identical fading parameters or average powers and without making any approximation for the covariance matrix. In particular, novel expressions utilizing infinity series representations for the joint PDF, CDF, MGF and moments of the arbitrary trivariate and quadrivariate Weibull distributions will be presented. These analytical expressions are being conveniently used to evaluate the OP, the average bit error probability (ABEP) and other significant performance metrics for the case of SC and MRC diversity reception.

II. Results and Discussion

To investigate the trivariate and quadrivariate Weibull distributions, it is convenient to consider the multivariate Weibull distribution, $\mathbf{Z}_L = \{Z_1, Z_2, \dots, Z_L\}$. \mathbf{Z}_L is assumed to be arbitrarily correlated according to a positive definite covariance matrix $\mathbf{\Psi}_L$, with elements $\psi_{i\kappa} = \mathbb{E} \langle G_i G_\kappa^* \rangle$, where $\mathbb{E} \langle \cdot \rangle$ denotes expectation, $*$ complex conjugate, $i, \kappa \in \{1, 2, \dots, L\}$ and $\mathbf{G}_L = \{G_1, G_2, \dots, G_L\}$ being joint complex zero mean Gaussian L RVs. Since $\psi_{i\kappa}$ can take arbitrary values, the analysis presented in this section refers to the most general correlation case.

A. Trivariate Weibull Distribution

For the case of the trivariate (i.e. $L = 3$) arbitrarily correlated¹ Weibull distribution and by applying the Weibull power transformation $Z = R^{2/\beta}$ [8, eq. (2)] in the infinite series representation of the Rayleigh distribution [4, eq. (5)], the novel joint PDF of $\mathbf{Z}_3 = \{Z_1, Z_2, Z_3\}$ has been derived as

¹From now on and unless otherwise stated, it will be assumed that the Weibull distributions under consideration are arbitrary correlated.

follows [11], [12]

$$\begin{aligned}
f_{\mathbf{Z}_3}(z_1, z_2, z_3) &= \frac{\beta_1 \beta_2 \beta_3 \det(\Phi_3)}{z_1^{(2-\beta_1)/2} z_2^{(2-\beta_2)/2} z_3^{(2-\beta_3)/2}} \exp \left[- \left(z_1^{\beta_1} \phi_{11} + z_2^{\beta_2} \phi_{22} + z_3^{\beta_3} \phi_{33} \right) \right] \\
&\times \sum_{k=0}^{\infty} \epsilon_k (-1)^k \cos(k\chi) \sum_{\ell, m, n=0}^{\infty} \frac{|\phi_{12}|^{2\ell+k}}{\ell!(\ell+k)!} \frac{|\phi_{23}|^{2m+k}}{m!(m+k)!} \frac{|\phi_{31}|^{2n+k}}{n!(n+k)!} \\
&\times z_1^{\beta_1(\ell+n+k)+\beta_1/2} z_2^{\beta_2(\ell+m+k)+\beta_2/2} z_3^{\beta_3(m+n+k)+\beta_3/2}
\end{aligned} \tag{1}$$

where ϵ_k is the Neumann factor ($\epsilon_0 = 1, \epsilon_k = 2$ for $k = 1, 2, \dots$), $\chi = \chi_{12} + \chi_{23} + \chi_{31}$ and Φ_3 is the inverse covariance matrix given by

$$\Phi_3 = \Psi_3^{-1} = \begin{bmatrix} \phi_{11} & \phi_{12} & \phi_{13} \\ \phi_{12}^* & \phi_{22} & \phi_{23} \\ \phi_{13}^* & \phi_{23}^* & \phi_{33} \end{bmatrix} \tag{2}$$

where $\phi_{i\kappa} = |\phi_{i\kappa}| \exp(j\chi_{i\kappa})$ with $i, \kappa \in \{1, 2, 3\}$ and $|\cdot|$ denoting absolute.

By integrating (1), an infinite series representation for the CDF of \mathbf{Z}_3 is derived as [11], [12]

$$\begin{aligned}
F_{\mathbf{Z}_3}(z_1, z_2, z_3) &= \frac{\det(\Phi_3)}{\phi_{11}\phi_{22}\phi_{33}} \sum_{k=0}^{\infty} \epsilon_k (-1)^k \cos(k\chi) \sum_{\ell, m, n=0}^{\infty} C_3 \nu_{12}^{\ell+k/2} \nu_{23}^{m+k/2} \nu_{31}^{n+k/2} \\
&\times \gamma \left(\delta_1, z_1^{\beta_1} \phi_{11} \right) \gamma \left(\delta_2, z_2^{\beta_2} \phi_{22} \right) \gamma \left(\delta_3, z_3^{\beta_3} \phi_{33} \right)
\end{aligned} \tag{3}$$

where $C_3 = [\ell!(\ell+k)!m!(m+k)!n!(n+k)!]^{-1}$, $\nu_{i\kappa} = |\phi_{i\kappa}|^2 / \phi_{ii}\phi_{\kappa\kappa}$, $\delta_1 = \ell + n + k + 1$, $\delta_2 = m + \ell + k + 1$, and $\delta_3 = n + m + k + 1$ with $\gamma(\cdot, \cdot)$ denoting the incomplete lower Gamma function [13, eq. (3.381/1)].

The joint MGF of \mathbf{Z}_3 can expressed as $M_{\mathbf{Z}_3}(s_1, s_2, s_3) = \mathbb{E}(\exp(-s_1 Z_1 - s_2 Z_2 - s_3 Z_3))$. From (1) and following the integral solutions using the Meijer G-function presented in [8, pp. 3610], the following novel expression has been obtained [11], [14]

$$\begin{aligned}
M_{\mathbf{Z}_3}(s_1, s_2, s_3) &= \beta_1 \beta_2 \beta_3 \det(\Phi_3) \sum_{k=0}^{\infty} \epsilon_k (-1)^k \cos(k\chi) \\
&\times \sum_{\ell, m, n=0}^{\infty} C_3 \frac{|\phi_{12}|^{2\ell+k} |\phi_{23}|^{2m+k} |\phi_{31}|^{2n+k}}{s_1^{\beta_1(\ell+n+k+1)} s_2^{\beta_2(\ell+m+k+1)} s_3^{\beta_3(m+n+k+1)}} \\
&\times \Upsilon \left[\frac{\phi_{11}}{s_1^{\beta_1}}, \beta_1(\ell+n+k+1) \right] \Upsilon \left[\frac{\phi_{22}}{s_2^{\beta_2}}, \beta_2(\ell+m+k+1) \right] \\
&\times \Upsilon \left[\frac{\phi_{33}}{s_3^{\beta_3}}, \beta_3(m+n+k+1) \right]
\end{aligned} \tag{4}$$

where $\Upsilon(\cdot)$ is given in [8, eq. 8].

B. Quadrivariate Weibull Distribution

For the case of the quadrivariate (i.e. $L = 4$) Weibull distribution, we consider the inverse covariance matrix, Φ_4 , expressed as

$$\Phi_4 = \Psi_4^{-1} = \begin{bmatrix} \phi_{11}, \phi_{12}, \phi_{13}, & 0 \\ \phi_{12}^*, \phi_{22}, \phi_{23}, \phi_{24} \\ \phi_{13}^*, \phi_{23}^*, \phi_{33}, \phi_{34} \\ 0, & \phi_{24}^*, \phi_{34}^*, \phi_{44} \end{bmatrix} \tag{5}$$

where the $\phi_{i\kappa}$ $i, \kappa \in \{1, 2, 3, 4\}$ can take arbitrary values with the restriction of $\phi_{14} = \phi_{14}^* = 0$. Although this restriction is a mathematical assumption, necessary for the derivation of the equivalent statistics and does not correspond to a physical explanation, it is underlined that our approach is more general than of [15] for the multivariate Rayleigh distribution. More specifically, the statistical properties derived in [15] hold only under the assumption that Ψ is tridiagonal, i.e. when $\phi_{i\kappa} = 0$ for $|i - \kappa| > 1$. The same assumption was used in [9], where the correlation matrix was approached by the tridiagonal Green matrix.

In principle, an expression for the joint PDF of $\mathbf{Z}_4 = \{Z_1, Z_2, Z_3, Z_4\}$ can be derived using [4, eq. (16)] and by applying the power transformation described in [8, eq. (2)] as a product of the modified Bessel function of the first kind $I_n(u)$. However, this approach will not be adopted since expressions containing modified Bessel functions are difficult to be mathematically

manipulated, e.g. performing integrations. Instead, a more convenient approach is to use its infinite series expansion [13, eq. (8.447/1)]. Thus, the following PDF has been obtained [11]

$$\begin{aligned}
f_{\mathbf{Z}_4}(z_1, z_2, z_3, z_4) &= \beta_1 \beta_2 \beta_3 \beta_4 \det(\mathbf{\Phi}_4) \exp \left[- \left(z_1^{\beta_1} \phi_{11} + z_2^{\beta_2} \phi_{22} + z_3^{\beta_3} \phi_{33} + z_4^{\beta_4} \phi_{44} \right) \right] \\
&\times \sum_{j=0}^{\infty} \sum_{k=-\infty}^{\infty} \epsilon_j (-1)^{j+k} \cos(A) \sum_{\ell, m, n, p, q=0}^{\infty} C_4 |\phi_{12}|^{2\ell+j} |\phi_{13}|^{2m+j} |\phi_{24}|^{2n+|k|} |\phi_{34}|^{2p+|k|} |\phi_{23}|^{2q+|j+k|} \\
&\times z_1^{\beta_1(\ell+n+j+1)-1} z_2^{\beta_2[J_1]-1} z_3^{\beta_3[J_2]-1} z_4^{\beta_4(n+p+|k|/2+1)-1}
\end{aligned} \tag{6}$$

where $C_4 = [\ell!(\ell+j)!m!(m+j)!n!(n+|k|)!p!(p+|k|)!q!(q+|k+j|)!]^{-1}$, $A = j(\chi_{12} + \chi_{23} + \chi_{31}) + k(\chi_{23} + \chi_{34} + \chi_{42})$, $J_1 = \ell + n + q + (j + |k| + |j + k|)/2 + 1$ and $J_2 = m + p + q + (j + |k| + |j + k|)/2 + 1$.

Following a similar procedure as before and by using (6), the corresponding CDF and MGF have also been obtained [11], [16].

C. Performance Analysis

In this section important performance criteria for diversity receivers with three or four arbitrarily correlated diversity branches operating over Weibull fading and additive white Gaussian noise (AWGN) channels will be presented. In particular, by using the previously derived expressions for the statistical characteristics of the trivariate and quadrivariate Weibull distribution, the performance of MRC and SC diversity receivers have been studied and their OP and ABEP have been derived.

For the system model considered, the equivalent baseband signal received at the ℓ th branch can be mathematically expressed as $\zeta_\ell = wh_\ell + n_\ell$ where w is the complex transmitted symbol having average energy $E_s = \mathbb{E}\langle |w|^2 \rangle$, h_ℓ is the complex channel fading envelope with its magnitude $Z_\ell = |h_\ell|$ being a Weibull distributed RV and n_ℓ is the AWGN with single-sided power spectral density N_0 . The instantaneous, per symbol SNR, of the ℓ th diversity channel is $\gamma_\ell = Z_\ell^2 E_s / N_0$, while its average is $\bar{\gamma}_\ell = \mathbb{E}\langle Z_\ell^2 \rangle E_s / N_0 = \Gamma(d_{2,\ell}) \Omega_\ell^{2/\beta_\ell} E_s / N_0$ where $d_{\tau,\ell} = 1 + \tau/\beta_\ell$ with $\tau > 0$. Note that it is straightforward to obtain expressions for the statistics of γ_ℓ by replacing at the previously mentioned expressions for the fading envelope Z_ℓ , β_ℓ with $\beta_\ell/2$ and Ω_ℓ with $(\alpha_\ell \bar{\gamma}_\ell)^{\beta_\ell/2}$ [8]. Thus, denoting $\boldsymbol{\gamma}_L = \{\gamma_1, \gamma_2, \dots, \gamma_L\}$, and since the CDF $F_{\boldsymbol{\gamma}_L}(\gamma_1, \gamma_2, \dots, \gamma_L)$ and the MGF $M_{\boldsymbol{\gamma}_L}(s_1, s_2, \dots, s_L)$ of the SNR for the trivariate and quadrivariate

Weibull distribution can be easily obtained, but will not be presented here due to space limitation.

1) Performance of MRC Receivers

For MRC receivers the output, per symbol, SNR (SNR_o), is $\gamma_{mrc} = \sum_{\ell=1}^L \gamma_{\ell}$ [1]. To obtain the ABEP performance it is convenient to use the MGF-based approach. Hence, the MGF of the L -branch MRC output can be derived as $M_{\gamma_{mrc}}(s) = M_{\gamma_L}(s, s, \dots, s)$. By using the MGF-based approach, the ABEP of noncoherent binary frequency-shift keying (NBFSK) and binary differential phase-shift keying (BDPSK) modulation signaling can be directly calculated. For other types of modulation formats, numerical integration is needed in order to evaluate single integrals with finite limits.

2) Outage Probability of SC Receivers

The instantaneous SNR at the output of a L -branch SC receiver, will be the SNR with the highest instantaneous value between all branches, i.e. $\gamma_{sc} = \max\{\gamma_1, \gamma_2, \dots, \gamma_L\}$ [17]. Since the CDF of γ_{sc} , $F_{\gamma_{sc}}(\gamma_{sc}) = F_{\gamma}(\gamma_{sc}, \gamma_{sc}, \dots, \gamma_{sc})$, P_{out} can be easily obtained as $P_{out}(\gamma_{th}) = F_{\gamma_{sc}}(\gamma_{th})$ for both trivariate and quadrivariate cases.

D. Performance Evaluation Results

Using the previous mathematical analysis, in this section performance evaluation results for the SC and MRC receivers will be presented. Non-identical distributed Weibull channels, i.e., $\bar{\gamma}_{\ell} = \bar{\gamma}_1 \exp[-(\ell - 1)\delta]$ where δ is the power decay factor are considered and for the convenience of the presentation, but without any loss of generality, $\beta_{\ell} = \beta \forall \ell$ will be assumed. Considering a triple-branch diversity receiver with the linearly arbitrary normalized covariance matrix² given in [2, pp. 886] and SC diversity, the OP has been obtained as a function of the first branch normalized outage threshold $\gamma_{th}/\bar{\gamma}_1$ for different values of β and δ . The performance evaluation results, illustrated in Fig. 1, indicate that P_{out} degrades with increasing $\gamma_{th}/\bar{\gamma}_1$ and δ and/or decreasing β . Note that for $\beta = 2$ and $\delta = 0$ the obtained results are in agreement with previously known performance evaluation results presented in [9].

²Note that the covariance matrix specifies the fading correlation between two complex Gaussian RVs.

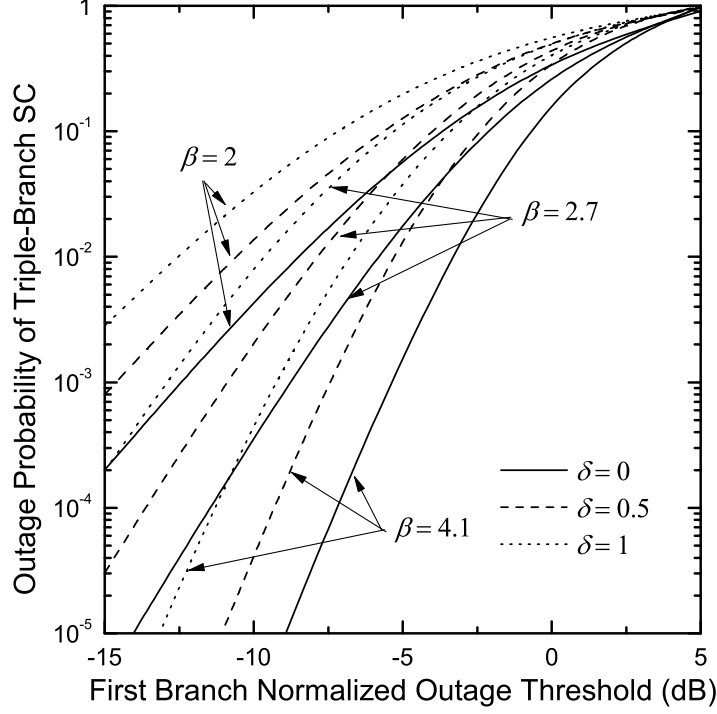


Figure 1: Outage probability of triple-branch SC receiver as a function of the first branch normalized outage threshold for different values of β and δ .

For MRC receiver and BDPSK signaling, the ABEP has been obtained and is illustrated in Fig. 2 for four receiving branches, assuming the covariance matrices presented in [4, eq. (34)]. As expected, the ABEP improves as the first branch average input SNR $\bar{\gamma}_1$ increases, while for a fixed value of $\bar{\gamma}_1$, similar to the SC diversity, a decrease of β and/or an increase of δ degrades the ABEP. Furthermore, performance evaluation results obtained by means of computer simulation also shown in Fig. 2 and have verified the accuracy of the analysis. It is finally noted that for the four-branch diversity reception and $\bar{\gamma}_1 > 5$ dB, only one term is required to achieve accuracy better than 10^{-5} .

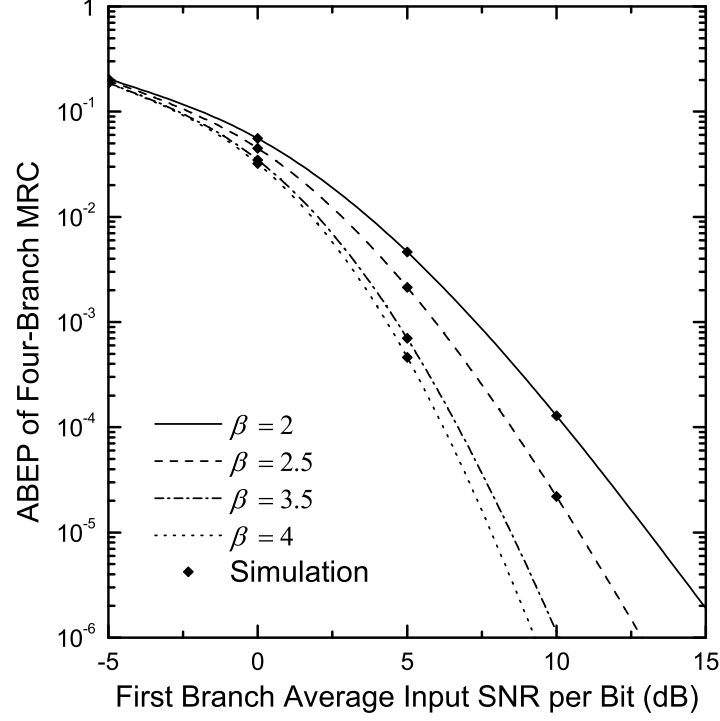


Figure 2: ABEP of four-branch MRC receiver as a function of the first branch average input SNR per bit for different values of β .

III. Conclusions

In this dissertation the novel statistical characteristics of the trivariate and quadrivariate Weibull fading distribution with arbitrary correlation, non-identical fading parameters and average powers have been derived using infinite series representations. The theoretical analysis has been also applied in order to evaluate the performance of SC and MRC diversity receivers.

References

- [1] M. K. Simon and M.-S. Alouini, *Digital Communication over Fading Channels*, 2nd ed. New York: Wiley, 2005.

- [2] G. K. Karagiannidis, D. A. Zogas, and S. A. Kotsopoulos, "An efficient approach to multivariate Nakagami- m distribution using Green's matrix approximation," *IEEE Trans. Wireless Commun.*, vol. 2, no. 5, pp. 883–889, Sep. 2003.
- [3] Y.-K. Ko, M.-S. Alouini, and M. K. Simon, "Outage probability of diversity systems over generalized fading channels," *IEEE Trans. Commun.*, vol. 48, no. 11, pp. 1783–1787, Nov. 2000.
- [4] Y. Chen and C. Tellambura, "Infinite series representation of the trivariate and quadrivariate Rayleigh distribution and their applications," *IEEE Trans. Commun.*, vol. 53, no. 12, pp. 2092–2101, Dec. 2005.
- [5] R. K. Mallik, "On the multivariate Rayleigh and exponential distributions," *IEEE Trans. Inf. Theory*, vol. 49, no. 6, pp. 1499–1515, Jun. 2003.
- [6] P. Dharmawanse, N. Rajatheva, and C. Tellambura, "Infinite series representation of the trivariate and quadrivariate Nakagami- m distributions," *IEEE Trans. Wireless Commun.*, vol. 6, no. 12, pp. 4320–4328, Dec. 2007.
- [7] P. S. Bithas, G. K. Karagiannidis, N. C. Sagias, P. T. Mathiopoulos, S. A. Kotsopoulos, and G. E. Corazza, "Performance analysis of a class of GSC receivers over nonidentical Weibull fading channels," *IEEE Trans. Veh. Technol.*, vol. 54, no. 6, pp. 1963–1970, Nov. 2005.
- [8] N. C. Sagias and G. K. Karagiannidis, "Gaussian class multivariate Weibull distributions: Theory and applications in fading channels," *IEEE Trans. Inf. Theory*, vol. 51, no. 10, pp. 3608–3619, Oct. 2005.
- [9] G. C. Alexandropoulos, N. C. Sagias, and K. Berberidis, "On the multivariate Weibull fading model with arbitrary correlation matrix," *IEEE Antennas Wireless Propag. Lett.*, vol. 6, pp. 93–95, 2007.
- [10] W. Weibull, "A statistical distribution function of wide applicability," *Appl. Mech. J.*, no. 27, 1951.
- [11] Z. G. Papadimitriou, P. T. Mathiopoulos, and N. C. Sagias, "The trivariate and quadrivariate weibull fading distributions with arbitrary correlation and their applications to diversity reception," *IEEE Trans. Commun.*, vol. 57, no. 11, pp. 3230–3234, Nov. 2009.

- [12] Z. G. Papadimitriou, N. C. Sagias, P. S. Bithas, P. T. Mathiopoulos, and L. Merakos, "The trivariate weibull distribution with arbitrary correlation," in *IEEE International Workshop on Satellite and Space Communications*, Leganes, Spain, Dec. 2006, pp. 249–253.
- [13] I. S. Gradshteyn and I. M. Ryzhik, *Table of Integrals, Series, and Products*, 6th ed. New York: Academic Press, 2000.
- [14] Z. G. Papadimitriou, P. S. Bithas, P. T. Mathiopoulos, N. C. Sagias, and L. Merakos, "Triple-branch mrc diversity in weibull fading channels," in *International Workshop on Signal Design and Its Applications in Communications*, Chengdu, China, Sep. 2007, pp. 247–251.
- [15] L. E. Blumenson and K. S. Miller, "Properties of generalized Rayleigh distributions," *Ann. Math. Statist.*, vol. 34, pp. 903–910, 1963.
- [16] Z. G. Papadimitriou, P. T. Mathiopoulos, N. C. Sagias, and L. Merakos, "On the weibull distribution with arbitrary correlation," in *3rd International Symposium on Communications, Control and Signal Processing*, St. Julians, Malta, Mar. 2009.
- [17] G. K. Karagiannidis, D. A. Zogas, and S. A. Kotsopoulos, "Performance analysis of triple selection diversity over exponentially correlated Nakagami- m fading channels," *IEEE Trans. Commun.*, vol. 51, no. 8, pp. 1245–1248, Aug. 2003.

Decision Management and Object-oriented Protocol and Services Reconfiguration in Future Internet Autonomic and Heterogeneous Telecommunication Environments

Eleni Patouni¹

¹ Department of Informatics & Telecommunications,
University of Athens, Athens, Greece,
elenip@di.uoa.gr

Abstract. In recent years, the vast evolution in telecommunication systems is remarkable, as regards the fast development and incorporation of new technologies in the heterogeneous networking environment. One major issue concerns the complexity management in user connectivity, in relation with two fundamental alternative solutions: the development of interworking functions for handovers between heterogeneous systems and the introduction of mechanisms for dynamic adaptation and reconfiguration. This phd thesis focuses on the second solution, which is called to overcome the design limitations of the first one. Specifically, the addressed issues concern the introduction and impact of reconfiguration in local (per device) and network levels, for the component-based dynamic adaptation of mobile devices and network elements. In the context of this thesis, special emphasis is paid on the specification and the detailed design of the reconfiguration deployment in the protocol stack and the service level, using object-oriented models. In addition, the mechanisms' evaluation and assessment was realized locally in the mobile devices, as regards the possibility of their deployment and the feasibility of the approach. Special focus was paid on the global evaluation and assessment of the introduced mechanisms for protocol reconfiguration in the heterogeneous network environment, taking into account different types of mobile devices with varying capabilities: reconfigurable and autonomous mobile devices.

Keywords: Reconfiguration, Protocol Component, Decision Making, Requests Management, Load

1 Dissertation Summary

Following the rapid proliferation in the development of Future Internet technologies and the increased complexity of telecommunication systems (mainly mobile and wireless), a major arising problem concerns the seamless mobility and the increased user QoS needs between these systems. In addition, a key challenge related

¹ *Dissertation Advisor: L. Merakos, Professor*

to the abovementioned developments in heterogeneous systems, is the introduction of flexible mechanisms for increased complexity management in the user connectivity.

In this direction, two major alternative solutions have been emerged: the first one concerns the introduction of interworking functionalities for handovers execution between heterogeneous systems and the second one concerns the introduction of dynamic adaptation and reconfiguration mechanisms [1]-[5]. The first solution is inline with the traditional approaches for the introduction of novel functionality in the underlying telecommunication infrastructure. As proven in the literature, this approach raises several limitations compared to the gradual incorporation of new radio access technologies and also imposes “domino” effect in the specification of telecommunication protocols for each access subsystem. The second solution is able to overcome the abovementioned limitations –at the same time the introduction and specification of mechanisms for reconfiguration realization is required (clean-slate approach).

This thesis deals with issues that are related to the introduction and impact of reconfiguration in local (per device) and network levels [6],[7]. This analysis deals with the specification, the detailed design and the evaluation of the object-oriented protocol stack and services reconfiguration approaches as regards their application/deployment. It should be noted that Future Internet autonomous and heterogeneous telecommunication environments are considered. Specifically, we propose a novel framework that allows the dynamic adaptation of mobile devices and network elements. It is worth noting that two alternatives have been developed, the first one is based on reconfigurable protocols while the second one proposes the introduction of intelligence through autonomous components. The mechanisms evaluation and assessment was realised locally in the mobile devices and has proven the applicability and the feasibility of the approach [7],[8]. Next, special focus was paid on the overall evaluation and assessment of the introduced mechanisms for protocol reconfiguration in the heterogeneous networking environment.

The overall evaluation and assessment has taken into account the existence of different types of mobile devices with different capabilities, considering two main categories of mobile devices: reconfigurable and autonomous. Such categories are distinguished mainly based on their capabilities for local decision making between the available alternatives (handover, protocol reconfiguration or joint solution for both handover and protocol reconfiguration). The local decision is then validated by the network side. Therefore the assessment and evaluation of these mechanisms in the overall heterogeneous networking environment has focused on the management of the produced decision-making requests originating from two categories of mobile devices [6].

As a whole, the introduced innovative mechanisms capture the two fundamental aspects of the reconfiguration procedure in mobile devices and network elements: a) the component-based, dynamic adaptation and reconfiguration of the protocol stack and the local evaluation and assessment and b) the overall evaluation for selecting the best alternative for the devices’ connectivity (e.g. handover, protocol stack reconfiguration). The two abovementioned technical challenges form the two fundamental technical aspects of this thesis.

In detail, after studying the evolution of mobile telecommunication systems and highlighting their limitations, we present the future heterogeneous mobile

telecommunication systems, their fundamental capabilities and their objectives. Such analysis also raises the technical challenges that form the framework of this thesis. Next, we introduce a methodology for the introduction and specification of the necessary functionalities for the protocol reconfiguration/self-configuration framework and the overall decision-making and management of the devices' connectivity in the heterogeneous networking environment. The initial analysis is realised using scenarios, deriving case studies and finally specifying the functionalities and the respective capabilities. Next, we specify and evaluate the framework and mechanisms for component-based protocol/services reconfiguration, considering two types of components, reconfigurable and autonomous and focusing on their dynamic binding and replacement during runtime. At this point, it should be noted that the introduction of dynamic reconfiguration capabilities in the protocol stack subsystem increases its flexibility but inevitably incurs performance penalties. In this direction, the qualitative and quantitative analysis of such mechanisms examines the applicability of the design approach.

As regards the second technical challenge, the decision-making approach for mobile devices' reconfiguration is specified, analysed and evaluated both on mobile devices and heterogeneous radio-network environments. The first case focuses on examining the alternative protocol configurations and identifying the best configuration. The cognitive decision-making approach for mobile devices is modelled using fuzzy-logic [9]. The produced results reveal that the introduced mechanisms do not affect the responsivity of the device or the user experience. Concerning the second case, we introduce and analyse the system model and the algorithmic framework for the network decision-making and management, as regards mobile devices' adaptation in heterogeneous radio-network environments [6]. We consider two main adaptation alternatives, handover and protocol reconfiguration. Two types of mobile devices are also assumed in our system: reconfigurable and autonomous. The goal of this analysis is to guide the mobile devices relocation for realising load balancing. In addition, the qualitative and quantitative evaluation of the introduced mechanisms and algorithmic framework is realised. The results also reveal that that transition to learning-capable dynamically self-managed mobile devices yields more efficient management of the decision-making requests. Moreover, the simulation results show the gain of using the proposed concepts in a system, in terms of applying load balancing techniques for requests management. Results show the number and percentage of dropped requests versus the amount of mobile devices and other key parameters. The outcome of this analysis reveals a quite unexpected conclusion: the introduction of autonomicity in the devices adversely impacts the requests management process in the network. The analysis quantifies how increasing the autonomicity level of the mobile devices affects the network load. At the same time, we propose a mechanism for maximizing the percentage of requests handled by the network, compared to the percentage of dropped requests. Moreover, our work reveals the degree of performance deterioration caused by increasing the autonomicity level in the management of requests [6].

1.1 Related Work

One path in achieving flexibility and intelligence in the systems and addressing the heterogeneity and complexity challenges, is realised through the emerging visions of reconfigurability, cognition and autonomic networking [1],[2],[3],[3] and [5]. The latter bring forward new adaptation capabilities in the different layers of the protocol stack and system resources. Such aspects have been addressed in the literature. Specifically, adaptable protocol stacks are seen as a technological enabler of next-generation networks which leverage the introduced adaptation and customization capacities to achieve two main goals: the dynamic adjustment of protocols' operation mode and the performance optimization of the operating protocols/protocol stacks. Such targets have been the main research objective of various approaches, which are classified into the following three main categories - a detailed survey analysis on dynamically adaptable protocol stack frameworks is available in our work in [10].

- **Adaptable protocols:** This design approach introduces an extension protocol layer besides the generic part, for the implementation of custom protocol functions. This category employs a coarse granularity design since the fundamental design unit is a protocol layer or a set of them. Adaptable protocol stack frameworks include Conduits, JChannels and POEM [10].
- **Composable protocols:** this concept employs flat, hierarchical and graph-based models for building a customizable protocol/protocol stack out of fundamental protocol functions. Composable protocol stack frameworks include DiPS/CuPS, x-kernel, Coyote/Cactus, Appia, Ensemble, Horus, RBA, Da CaPo, ADAPTIVE, DRoPS, DIPS+, ACCORD and DPS [10].
- **Reconfigurable protocols:** This design allows for extending the traditional protocol stacks' composition schemes to support the dynamic binding and replacement of protocol components or even entire protocol layers during runtime, enabling service continuity and no loss of protocol data. Reconfigurable protocol stack frameworks include THINK, FRACTAL, GRPSFMT, DRAPS and Alonistioti [10].

At this point it should be noted that our approach falls under the category of reconfigurable protocols. The main advantage of our work is the detailed specification of a framework enabling the dynamic, semantic-based binding and replacement of protocol components during runtime operation of the protocol stack. In addition, the necessary support and state management mechanisms were defined, targeting transparency, robustness and seamless operation.

The reconfiguration decision-making procedure imposes significant research challenges, which have been the objective of some research activities. [9] presents issues related to the management and control of reconfigurable radio systems, also addressing the decision-making procedure. In [6], we model a reconfigurable system as a distributed transactional system and examine the global bounds of the asymptotic network response time and throughput. Our work uses multiclass queuing networks for the system model and is based on the findings by Balbo and Serazzi [11], [12] and Litoiu [13], [14] for the derivation of the network bottlenecks and the bounds of the

response time under asymptotic and non-asymptotic conditions. Besides them, several approaches have been addressed for the development of approximation techniques to estimate performance measures such as queue lengths, sojourn times and throughput. The use of approximation techniques has greatly facilitated estimation and optimization of performance measures in finite queuing networks. Some of the defined optimization approaches in finite queuing networks are used in addressing some of these problems, as in [15], [16], [17] and [18]. In the aforementioned approaches the implications of reconfiguration decisions in network level are not addressed. Moreover, the system bounds for each framework are not discussed in consideration of the overall load and reconfiguration overhead in conjunction with the user and device classes and respective request patterns. In conclusion, the introduction and adaptation of such methods in order to discuss optimization issues in autonomic and reconfigurable telecommunication systems has not been considerably investigated in the literature and forms one of the key directions of this paper.

2 Results and Discussion

This section presents the key concepts and results of the second technical challenge addressed in this thesis. The first challenge is not elaborated herein due to space limitations - details are available in [7],[8] [19],[20] and [21].

2.1 Algorithmic Framework for Handling Decision-making Requests

The goal of this analysis is the definition and evaluation of the algorithmic framework for handling decision-making requests for protocol reconfiguration. The key phases of the algorithmic framework are presented in Figure 1.

The key metric of the proposed algorithmic framework is the user satisfaction metric, which forms a function of the network response time [6]. First of all, we define as network response time the response time experienced by a mobile device making a decision-making request to the network side. We differentiate the network response time per class of mobile devices. Therefore we define the response time of class c

R_c as the response time experienced by a class c mobile device making a decision-making request. We also define as user satisfaction SA_c , the normalized distance of the network response time R_c from the maximum value of the response time R_c^{\max} to the interval of the maximum response time minus the minimum response time R_c^{\min} . Therefore, user satisfaction is analysed as follows:

$$SA_c = \frac{R_c^{\max} - R_c}{R_c^{\max} - R_c^{\min}} \quad (1)$$

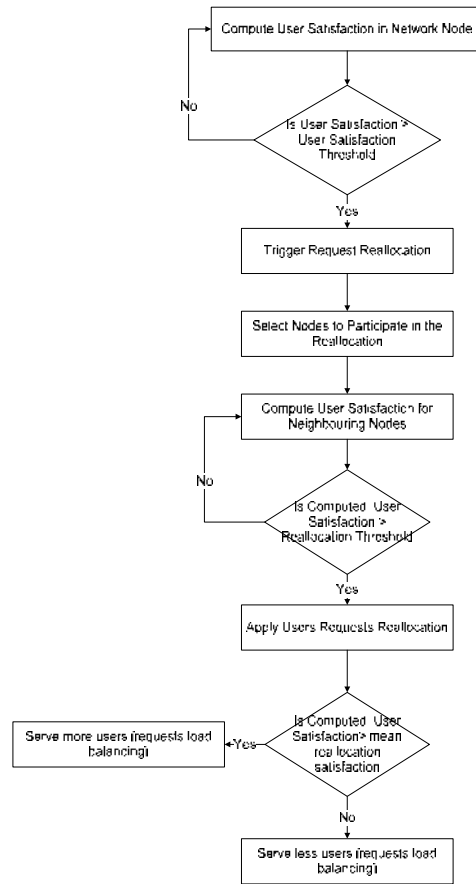


Figure 1: Flow diagram of the algorithmic framework for handling decision-making requests

At this point it should be noted that the maximum and minimum values of the response time are not static but dynamically varying based on the number of the decision-making requests and the average time between requests. For a given system, the response time and the respective maximum and minimum values of the response time can be computed using mean value analysis (MVA), an iterative technique for the analysis of closed queuing network models [17],[18]. This technique allows the computation of various performance metrics (e.g. response time) of any number of users iteratively (it introduces customers into the queuing network one by one, the cycle terminates when all customers have been entered).

However the computational complexity of MVA is very high and the storage requirements increase for networks with high numbers of mobile devices and classes. Therefore, instead of computing the user satisfaction, we compute an approximate value of the user satisfaction. This is realized by computing the bounds of the maximum and minimum response time. Therefore, the approximate value of user satisfaction is given below:

$$\overline{SA_c} = \frac{R_c^{Up} - R_c^{Ms}}{R_c^{Up} - R_c^L} \quad (2)$$

In this direction, the computation of user satisfaction requires to also compute the upper and lower bounds of the network response time and measure the network response time. To realize the bounds computation, we consider the analysis by Litoiu and Balbo, Serrazzi [12], [13], [14] and propose a methodology and respective analytical model for the computation of the approximate value of the user satisfaction metric. Such methodology concerns the bounds computation for the response time for distributed systems with multiple resources and workload mixes. The details on the methodology and analytical model for the bounds computation can be found at [6].

2.2 Results

The algorithmic framework for handling decision-making requests was evaluated through simulations. The simulations were realized using MATLAB Simulink tool. In this work, we developed four network nodes that manage the decision-making requests originating from mobile devices; such devices include both reconfigurable and autonomous mobile devices. In addition, we developed two separate load balancing systems that handle the decision-making requests per class.

At first, the use of the presented algorithmic framework was evaluated. Specifically, the user satisfaction degree is dynamically computed per class of mobile devices, using the outcomes of the previous subsections for the global bounds of the network response time for the case study system. Each of the network nodes is considered as a system with the same resource demands with the one analysed in [22].

In addition, in order to approach the behavior of a real system as regards the dynamic alteration of the network response time, we consider that it follows a gamma distribution.

Specifically, given a randomly generated value R_g , the next generated value follows a gamma distribution (as in [23]) in the interval $[A, B]$, where $A = R_g - e$ and

$$B = R_g + e, \quad e = \frac{R_c^{Up} + R_c^L}{\alpha}.$$

The shape parameter for the gamma distribution equals to $\left(\frac{A+B}{2}\right)/b$, where b is the scale parameter. R_c^{Up} and R_c^L are the upper and lower bounds of the network response time respectively computed for the specified number of mobile devices and think time. α is an integer – for this simulation work

we consider that $a \in [2,5]$. In addition, if $A < R_c^L$ then $A = R_c^L$; correspondingly if $B > R_c^{Up}$ then $B = R_c^{Up}$. This way A and B are always within the upper and lower bounds of the network response time. It should be also noted that each network node behavior can be dynamically altered during simulation, e.g. the number of mobile devices can be changed due to load balancing actions. Therefore, the global bounds of the network response time are dynamically computed for each class of mobile devices (using equations (15) and (17) in our analysis in [6]).

An important metric in this work is the think time metric, which represents the average time between requests and therefore affects the response time and respective bounds (details on the computations are available in [6]). In this model we define the think time z as the time interval between two handover or reconfiguration decision requests. Therefore

$$z = \frac{1}{E_{HO}} \quad (3)$$

where E_{HO} is the expected number of handovers in a system. Based on the analysis in [24], E_{HO} is a function of the call-to-mobility ratio denoted as ρ . E_{HO} is given below considering that the ratio of Access Routers (ARs) to Mobility Anchor Points (MAPs) equals 1.

$$E_{HO} = \frac{1}{\rho^2} + \frac{1}{\rho} \quad (4)$$

Therefore, the think time z for both classes of mobile devices is analysed as follows:

$$z = \frac{1}{\frac{1}{\rho^2} + \frac{1}{\rho}} \quad (5)$$

Next using equation (2) the user satisfaction is also dynamically computed for each class of mobile devices N . For the simulation, we consider that the user satisfaction threshold is 0.075. If the user satisfaction is found to be lower than this threshold, then the requests reallocation procedure is triggered. We also consider that the reallocation threshold of the nodes is 0.3 – this means that the nodes with user satisfaction lower than this value do not participate in the reallocation procedure. Thereafter, the reallocation procedure is applied following the concepts. The mean reallocation satisfaction is computed and the nodes with lower value than the mean reallocation satisfaction should allocate a percent of the serving mobile devices to the nodes with higher user satisfaction than the mean reallocation satisfaction. Simulation results include the dynamic variation of the total number of mobile devices per class, due to the load balancing procedure. It should be noted that the load balancing system may fail to reallocate some requests and drops them - this is expected when

the negotiation procedure fails, e.g. when the satisfaction threshold of all nodes is lower than 0.3.

In addition, we can dynamically compute if node requests are dropped and the actual percent of the times the load balancing system has to drop some requests over the total execution times of the load balancing. Secondly, the absence of the presented algorithmic framework in the same system was evaluated. More specifically, we consider that the mobile devices generate requests to the mobile nodes; we measure the network response time and using the global bounds of the response time derived from the analytical model we dynamically compute the user satisfaction. Since in this system we consider the absence of the proposed algorithmic framework for load balancing, we simply measure the user satisfaction over the simulation per class of mobile devices and we assume that when user satisfaction equals zero, then 5% of the node requests are dropped.

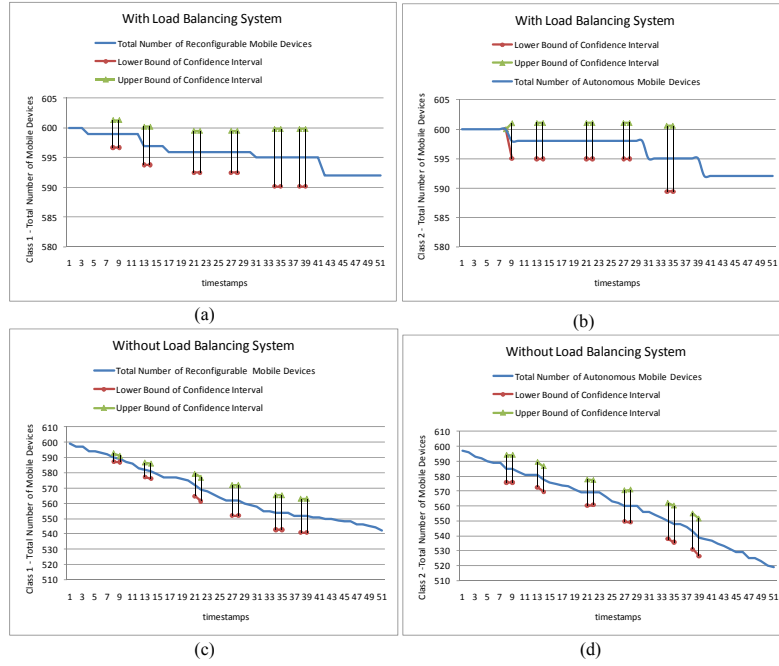


Figure 2: Total number of mobile devices: a) applying the load balancing system for reconfigurable mobile devices, b) applying the load balancing system for autonomous mobile devices, c) without applying the load balancing system for reconfigurable mobile devices, d) without applying the load balancing system for autonomous mobile devices

Figure 2 presents the variation of the total number of mobile devices for each class (also including the respective 95% confidence intervals), both with and without the application of the load balancing system. As seen in Figure 2 a and b, the application of the load balancing system results in dropping 1.33% of the total number of mobile devices (8 out of 600), whereas a legacy system results in dropping 9.67% of class 1 requests and 10.17% of class 2 mobile devices (58/600 and 61/600 respectively), as seen in Figure 2 c and d. At this point we should note how the type of mobile devices affects the load balancing behavior. The results show that the introduction of the load balancing system leads to the same behavior in terms of requests handling and dropping for both reconfigurable and autonomous mobile devices. Such outcome is quite unexpected and reveals that the optimization of the load balancing system is request- independent. In addition, if we consider the absence of load balancing system, we come to the conclusion that the system tends to drop more requests coming from autonomous mobile devices compared to the requests generated by reconfigurable mobile devices. This outcome of our analysis reveals one drawback of the introduction of autonomicity/intelligence in the mobile device.

In addition, simulation results include the load balancing failure percent when the load balancing system is applied– that is the percent of times the load balancing fails to reallocate the user requests over the total number of load balancing triggers. This percent was derived from simulation results versus the total number of mobile devices and the user satisfaction factor.

As analysed in [6], both classes of mobile devices have similar load balancing failure percentages, which do not increase in a linear manner as the number of mobile devices increases. The highest percentage is presented for mobile devices values between 200 and 300 – this also varies based on the SA threshold. Such low values of load balancing failures are expected since the load balancing system fails to reallocate the requests only when all system nodes are saturated (node SA lower than user satisfaction threshold) or close to being saturated (node SA lower than reallocation threshold). Again, we observe that the autonomous mobile devices tend to have greater load balancing failure percentages compared to reconfigurable mobile devices.

Conclusions

Reconfigurability is seen as one of the strong candidate concepts for the support of the convergence of heterogeneous systems, the evolution and migration of future communication systems, and the introduction of substantial flexibility in mobile systems. Furthermore, reconfigurability provides the ground for the development of yet more advanced concepts like cognitive and autonomic communications. In order to meet these expectations, a major issue is to establish a framework for enabling reconfiguration in all protocol layers, as well as plug-and-play solutions for protocol stack formation and activation. In this thesis, a generic architecture, respective interfaces, protocols and mechanisms for protocol stack and protocol component synthesis have been designed, implemented and illustrated. Finally, performance issues have been studied and the key performance metrics of protocol reconfiguration have been evaluated and discussed. Based on the quantitative and qualitative design

considerations set for the protocol reconfiguration attributes, as well as the discussion on limitations of other related frameworks, the proposed generic architecture satisfies the requirements related to flexibility, delay overhead, generic protocol component design, as well as plug and play capabilities.

In this thesis, we have also discussed the modeling and the impact of the network decision making process regarding handover and protocol reconfiguration in a heterogeneous networking environment, assuming two classes of mobile devices. The thesis has proposed an algorithmic framework for the management of the decision making requests for reconfiguration or handovers. Simulation results have also been presented for the alternative applications of the algorithmic framework for request relocation in a system, in terms of the percentage of reconfiguration or handover requests handled or dropped by the network. The outcome of this analysis shows how the increase in the autonomicity level of mobile devices affects the network load. Our work provides proof and tangible results of theoretical assumptions and statements relevant to the gains and applicability of autonomicity concepts for the first time in the literature, also addressing the pros and cons of introducing autonomicity in the mobile devices. Future work includes the extension of the presented concepts to advanced load balancing schemes (e.g. use of learning schemes) that will enable more proactive management of the decision-making requests handling.

References

1. M. Dillinger, K. Madami and N. Alonistioti, "Software Defined Radio – Architectures, Systems and Functions", Wiley series in Software Radio, 2003, ISBN 0-470-85164-3.
2. J. Mitola III, "Cognitive Radio for Flexible Mobile Multimedia Communications", *Mob. Netw. Appl Journal*, vol. 6, no.5, 2001, pp. 435–441.
3. S. Haykin, "Cognitive Radio: Brain-Empowered Wireless Communications", *IEEE J. Sel. Areas Commun.*, vol. 23, no. 2, Feb. 2005, pp. 201–220.
4. S. Dobson, S. Denazis, A. Fernández, D. Gaïti, E. Gelenbe, F. Massacci, P. Nixon, F. Saffre, N. Schmidt and F. Zambonelli, "A survey of autonomic communications", *ACM Trans. Auton. Adapt. Syst.*, vol.1, no. 2, 2006, pp. 223-259.
5. J. O. Kephart and D. M. Chess. "The Vision of Autonomic Computing," *IEEE Computer*, vol. 36, no.1, 2003, pp. 41-50.
6. E. Patouni, N. Alonistioti and L. Merakos, "Cognitive Decision Making for Reconfiguration in Heterogeneous Radio Network Environments", in *IEEE Transactions on Vehicular Technology (TVT)*, special issue on "Achievements and the Road Ahead: The First Decade of Cognitive Radio", vol. 59, issue 4, May 2010, pp. 1887-1900.
7. E. Patouni, S. Gault, M. Muck, N. Alonistioti, K. Kominaki, "Advanced Reconfiguration Framework based on Game Theoretical Techniques in Autonomic Communication Systems", *Annals of Telecommunication Journal*, vol. 62, no. 9-10, 2007, pp. 1099-1120.
8. N. Alonistioti, E. Patouni, V. Gazis, "Generic architecture and mechanisms for protocol reconfiguration", *Mobile Networks and Applications Journal*, Special Issue on Reconfigurable Radio Technologies in Support of Ubiquitous Seamless Computing, vol. 11, no. 6, December 2006, pp. 917-934.
9. A. Merentitis, E. Patouni, N. Alonistioti and M. Doubrava, "To Reconfigure or Not to Reconfigure: Cognitive Mechanisms for Mobile Devices Decision Making", in the

- Proceedings of the 68th IEEE Vehicular Technology conference, 21–24 September 2008, Calgary, Alberta,
10. V. Gazis, E. Patouni, N. Alonistioti and L. Merakos, "A Survey of Dynamically Adaptable Protocol Stacks", IEEE Communications Surveys and Tutorials, vol.12, no.1, January 2010, pp.3-23.
 11. G. Balbo, and G. Serazzi, "Asymptotic analysis of multiclass closed queueing networks: multiple bottlenecks", Performance Evaluation Journal, v.30 n.3, Sept. 1997, pp.115-152.
 12. G. Balbo, and G. Serazzi, "Asymptotic analysis of multiclass closed queueing networks: common bottleneck", Performance Evaluation Journal, vol.26, no.1, Jul. 1996, pp.51-72.
 13. E. Marin Litoiu, "A performance analysis method for autonomic computing systems", ACM Transactions on Autonomous and Adaptive Systems, vol.2 n.1, March 2007, p.3-es.
 14. M. Litoiu, J. Rolia and G. Serazzi, "Designing Process Replication and Activation: A Quantitative Approach", IEEE Transactions on Software Engineering, vol.26, no.12, Dec. 2000, pp.1168-1178.
 15. R. Onvural, "Survey of closed queueing networks with blocking", ACM Comput. Surv. vol. 22, no.2, 1990, pp. 83-121.
 16. H. G. Perros, "Queueing networks with blocking- a bibliography", ACM Sigmet. Perform. Eval. Rev, vol. 12, no. 2, 1984, pp. 8-12.
 17. J. M. Smith and S. Daskalaki, "Buffer space allocation in automated assembly lines", Op Res., vol.36, no. 2, 1988, pp.343-358.
 18. L. Kerbache and J. MacGregor Smith, "Asymptotic behavior of the expansion method for open finite queueing networks", Comput. Ops Res., vol. 15, no.2, 1988, pp.157-169.
 19. E. Patouni, O. Holland and N. Alonistioti, "Cognitive Functionalities for Mobile Terminal Self-Recovery and Protocol Auto-Configuration", in the Proceedings of the IEEE International Symposium on Personal, Indoor and Mobile Radio Communications 2008 (PIMRC 2008), 15-18 September 2008, Cannes, France,
 20. E. Patouni, S. Gault, M. Muck, N. Alonistioti, K. Kominaki, "Autonomic Communications: Exploiting Advanced and Game theoretical Approaches for RAT Selection and Protocol Reconfiguration", in the Proceedings of the Autonomic Networking Conference, Paris, France, September 2006.
 21. E. Patouni and N. Alonistioti, A Framework for the Deployment of Self-Managing and Self-Configuring Components in Autonomic Environments, in the Proceedings of the International IEEE WoWMoM Workshop on Autonomic Communications and Computing (ACC 06), Niagara-Falls / Buffalo-NY, 26-29 June, 2006
 22. P. Magdalinos, S. Polymeneas, P. Gliatis, X. Fafoutis, A. Merentitis and C. Polychronopoulos, "A Proof of Concept Architecture for Self-Configuring Autonomic Systems", presented at the ICT-MobileSummit 2008 Conf., Stockholm, Sweden, June 10 – 12, 2008.
 23. K. Y. Jo and S.R. Ali, "Design and analysis of large-scale wireless communications networks", in Proc. of the ACM Int. workshop on Performance monitoring, measurement, and evaluation of heterogeneous wireless and wired networks, New York, NY, USA, 2006, pp.18-24.
 24. T.C. Schmidt and M. Waehlich, "Analysis of Handover Frequencies for Predictive, Reactive and Proxy Schemes and Their Implications on IPv6 and Multicast Mobility", in Proc. ICN 2005. 4th International Conference on Networking, Part II, ser. Lecture Notes in Computer Science, P. Lorenz and P. Dini, Eds., Berlin Heidelberg: Springer-Verlag, vol. 3421, pp. 1039-1046, Apr. 2005.

A Web Usage Mining Framework for Web Directories Personalization

Dimitrios Pierrakos*

Department of Informatics & Telecommunications., University of Athens, Greece
Institute of Informatics & Telecommunications., NCSR “Demokritos”, Athens, Greece
`dpie@iit.demokritos.gr`

Abstract. In this thesis we propose a novel framework that combines Web personalization and Web directories, which results in the concept of Community Web Directories. Community Web directories is a novel form of personalization performed on Web directories, that correspond to “segments” of the directory hierarchy, representing the interests and preferences of user communities. The proposed approach is based on Web usage mining and the usage data that are analyzed here correspond to user navigation throughout the Web, rather than a particular Web site. For the construction of Community Web Directories, we introduce three novel techniques that combine the users’ browsing behavior with thematic information from the Web directories. Finally, we present OurDMOZ, a system that builds and maintains community Web directories.

1 Introduction

Information overload is one of the Web’s major shortcomings, placing obstacles in the way users access the required information. Web Personalization, i.e., the task of making Web-based information systems adaptive to the needs and interests of individual users or group of users like *user communities*, emerges as an important means to tackle information overload. The first step towards achieving personalization is the specification of user models. However, acquiring and creating accurate and operational user models is a difficult task. Web Usage Mining is one such approach which employs knowledge discovery from data to create user models, based on the analysis of usage data, as we presented in [11].

Another attempt to alleviate the problem of information overload is the organization of the Web content into thematic hierarchies. These hierarchies are known as *Web Directories*, and correspond to listings of topics which are organized and overseen by humans. However, the size and the complexity of the Web directory are canceling out the gains that were expected with respect to the information overload problem.

The contribution of this thesis is a novel approach to overcome the deficiencies of Web personalization and Web directories by combining their strengths,

* Dissertation Advisors: Yannis Ioannidis, Professor - Georgios Paliouras, Researcher, NCSR “Demokritos”

providing a new tool to fight information overload. In particular, we focus on the construction of usable Web directories that model the interests of user communities. The construction of user community models, with the aid of Web Usage Mining has primarily been studied in the context of specific Web sites [6]. In this thesis, we have extended this approach to a much larger portion of the Web, through the analysis of usage data collected by the proxy servers of an Internet Service Provider (ISP).

In the course of this thesis, we developed and evaluated three community modeling techniques, based on clustering and probabilistic learning. These techniques allowed us to take advantage of existing Web directories and specialize them to the interests of particular communities. In addition to handling the “global information overload” problem, the proposed methods also deal effectively with the “local overload” problem. This problem is a side-effect of the pruning of a number of leaf nodes of the initial Web directory, which pushes the information that they contained, i.e., the terminal links to Web pages, upwards in the hierarchy. This leads to increased information density in some leaf nodes of the personalized directory. In order to address this issue, the proposed methods combine usage data with thematic information from the original Web directories.

The proposed methodology is tested on two types of Web directory: an *artificial Web directory*, that was constructed using document clustering from the Web pages included in the log files themselves, and a “*real*” *Web directory*, namely the Open Directory Project (ODP). The main difficulty in the latter approach was the association of usage data, i.e. the Web pages, to categories of the directory, given the small proportion of Web pages that are explicitly assigned (manually) to categories of the directory. We approached this problem by an automated page classification method.

Finally, we present OurDMOZ, a system that integrates and implements the various components of the proposed framework. The thesis includes the results of a user evaluation study, which assessed the potential benefits of OurDMOZ and consequently of community Web directories to the end user.

2 Related Work

A number of studies exploit Web directories to achieve a form of personalization. Automatic profile construction is proposed in [5]. The user profiles, linked to categories of the directory are used typically for personalized Web search, while the directory itself is not personalized. The personalization of Web directories is mainly represented by services such as Yahoo! and Excite (www.excite.com), which support the manual selection of interesting categories by the user. An initial approach to automate this process, was the Montage system [1].

Our work differs from the above cited approaches in several aspects. First, instead of using the Web directory for personalization, it personalizes the directory itself. Compared to existing approaches to directory personalization, it focuses on aggregate or collaborative user models such as user communities, rather than

content selection for single user. Furthermore, unlike most existing approaches, it does not require a small set of predefined thematic categories, which could complicate the construction of rich hierarchical models. Finally, the work presented in [2], which is closest to ours is limited to the recommendation of short navigation paths in the ODP hierarchy, rather than the personalization of the whole Web directory structure. Moreover, that method makes the assumption that usage data are collected from the navigation of users within the Web directory. Thus, its applicability to independent services such as a Web portal is questionable.

3 Discovery of Community Web Directories from Web Usage Data

The construction of community Web directories is seen in this thesis as a fully automated process, powered by operational knowledge, in the form of user models that are generated by Web usage mining. User communities are formed using data collected from Web proxies as users browse the Web. The goal is to identify interesting behavioral patterns in the collected usage data and construct community Web directories based on those patterns. The stages of getting from the data to the community Web directories (Figure 1) are summarized below and described in the following sections:

- *Usage Data Preparation*, comprising the collection and cleaning of the usage data.
- *Web Directory Initialization*, providing the characterization of the Web pages included in the usage data, according to the categories of a Web directory. There are two approaches for the characterization of the Web pages. The first approach is to classify them on an existing Web directory, like ODP. The second approach is to map them onto an artificial Web directory constructed from the Web pages themselves using a hierarchical document clustering approach.
- *Community Web Directory Discovery*, being the main process of discovering the user models from data, using machine learning techniques and exploiting these models to build the community Web directories.

3.1 Usage Data Preparation

The usage data that form the basis for the construction of the communities are collected in the access log files of proxy servers, e.g. ISP cache proxy servers. These data record the navigation of the ISP subscribers through the Web. The first stage of data preparation involves data cleaning. The next stage is the identification of individual user sessions. The lack of user registration data or other means of user identification, such as cookies, led us to exploit a simpler definition of user sessions. A user session is defined as a sequence of log entries, i.e., accesses to Web pages by the same IP address, where the time interval between two subsequent entries does not exceed a certain time threshold.

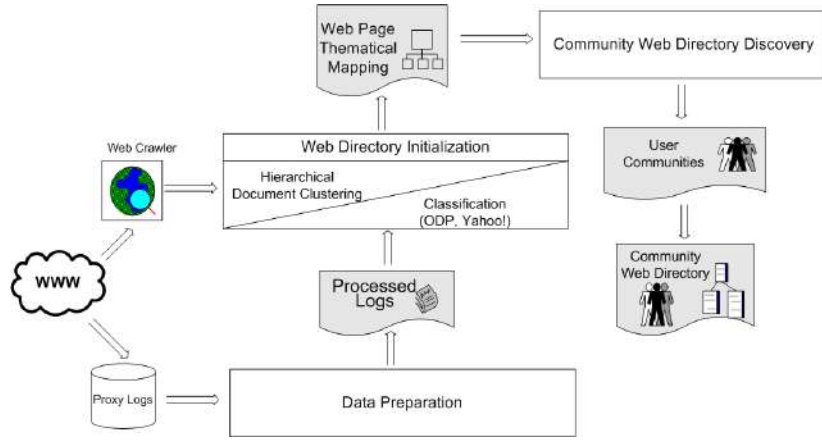


Fig. 1: The process of constructing Community Web directories.

3.2 Web Directory Initialization

The next stage towards the construction of community Web directories is the association of the users' browsing data with the Web directory. In order to personalize the Web directory, we need to "initialize" it with the users' data, i.e., "map" the Web pages onto the Web directory structure. This mapping requires the thematic categorization of the Web pages to the categories of the Web directory, since it is very unlikely to find many of the Web pages that appear in a log file in any directory.

As a first step in the categorization process, a crawler downloads the Web pages included in the usage data and encodes them using the vector space representation, by extracting the most important terms. There are two methodologies that we follow to realize the mapping of Web pages. The first one is based on document clustering and the second is based on document classification. Document clustering is performed following a hierarchical agglomerative approach, whilst document classification requires as a preliminary step, the extraction of keywords from the Web pages included initially in the categories of the Web directory.

The document clustering approach constructs a new Web directory from the usage data themselves. We call this directory an *Artificial Web directory*. The resulting hierarchy is a binary tree, representing clusters of Web pages that form thematic *categories*. This hierarchy corresponds to the initial, non-personalized Web directory, which provides directly a mapping between the Web pages and the categories that the pages are assigned to.

However and despite the similarity of the artificial directory to existing Web directories, there are also notable differences such as the artificial Web directory might not "cover" the semantics of new sessions, due to "overfitting" of the document clustering approach to the initial data. Thus we study the personalization

of a “real” Web directory and in particular the ODP. The main difficulty in this effort was the association of usage data, i.e., the Web pages, to categories of the Web directory, given the small proportion of Web pages that are explicitly assigned (manually) to categories of the directory. We approached this problem by an automated page classification method, which is described below.

Recall from the previous section that each Web page is represented as a vector of terms, we follow a similar vector space representation, for the node-categories of the ODP taxonomy. The Web pages are then classified onto the ODP hierarchy using cosine similarity.

3.3 Objective Category Informativeness

Each community directory includes only a subset of the categories of the initial Web directory, that represent the browsing preferences of the community. However, due to the structure of the Web directory, the community selection process leads to an undesirable side-effect: high-level categories that become leaves aggregate the Web pages of all of their sub-categories that are not in the community model, leading to a “cumulation” of information at the leaves of the reduced directory, which is bound to be overwhelming for its users. Therefore, although the “global” overload problem seems to be tackled well, a “local” overload arises.

To alleviate the “local” information overload problem, we introduce an additional criterion in the discovery of user communities. This criterion incorporates a measure of a-priori informativeness of the categories, which is taken into account when pruning leaf nodes from the Web directory. The inclusion of leaves that satisfy this criterion selectively reduces the generality of the directories, making them reflect more “fine-grained” user interests and resulting in a better distribution of the information that is indexed.

The new criterion is called *Objective Category Informativeness Association*, (*OCIA*), and is based on a measure of the *Mutual information*, (*MI*) of the leaf category l_n , to its parent category c_i . An improved version of MI is the *Symmetrical Uncertainty* (*SU*) measure, which normalizes MI by dividing by the sum of the entropies of \mathbf{C}_i and the leaf \mathbf{L}_n :

$$SU(\mathbf{C}_i, \mathbf{L}_n) = 2.0 \times \left[\frac{(H(\mathbf{L}_n) + H(\mathbf{C}_i) - H(\mathbf{C}_i, \mathbf{L}_n))}{H(\mathbf{C}_i) + H(\mathbf{L}_n)} \right]. \quad (1)$$

The value range of symmetrical uncertainty is [0..1]. Values closer to 0 indicate a weak association between the parent and the leaf category. Thus, leaf categories with a low association to their parents should be included in the community Web directories. OCIA is estimated by normalizing SU further, by the ratio of the number of pages of the leaf to the pages of the parent category, N_{l_n} , N_{c_i} respectively, in order to remove the bias towards leaf categories that contain a large number of Web pages. OCIA is given by the following equation:

$$OCIA(\mathbf{C}_i, \mathbf{L}_n) = \frac{N_{l_n}}{N_{c_i}} \times SU(\mathbf{C}_i, \mathbf{L}_n). \quad (2)$$

OCIA is the criterion that is used to decide whether a leaf node should be included in the community model. Only leaves for which *OCIA* is smaller than a designated *Parent-Children Association Threshold*, (*PCAT*), are selected. Thus, the subset $L'_i \subseteq L_i$ of these leaves is defined as:

$$L'_i = \{l_n \in L_i \mid OCIA(C_i, \mathbf{L}_n) \leq PCAT\}. \quad (3)$$

4 The *Objective Community Directory Miner (OCDM)* Algorithm

In this section we present the three algorithms that have been developed for the construction of community Web directories, as presented in [8], [9], and [10].

5 The *Objective Community Directory Miner (OCDM)* Algorithm

The first machine learning method that we employed for pattern discovery is the *Objective Community Directory Miner (OCDM)*, an enhanced version of the cluster mining algorithm [7]. Cluster mining discovers patterns of common behavior by looking for all maximal fully-connected subgraphs (cliques) of a graph that represents the users' characteristic features, i.e., thematic categories in our case.

OCDM enhances cluster mining so as to take into account the hierarchy of topic categories. This is achieved by updating the weights of the vertices and the nodes in the graph. Each category is mapped onto a set of categories, corresponding to its parent and grandparents in the thematic hierarchy. The frequency of each of these categories is increased by the frequency of the initial child category. The underlying assumption for the update of the weights is that if a certain category exists in the data, then its parent categories should also be examined for the construction of the community model. In this manner, even if a category (or a pair of categories) have a low occurrence (co-occurrence) frequency, their parents may have a sufficiently high frequency to be included in a community model. This enhancement allows the algorithm to start from a particular category and ascend the topic hierarchy accordingly. The result is the construction of a topic tree, even if only a few nodes of the tree exist in the usage data.

The connectivity of the resulting graph is usually high. For this reason we make use of a connectivity threshold that reduces the edges of the graph. This threshold is related to the frequency of co-occurrence of the thematic categories in the data. Once the connectivity of the graph has been reduced, the weighted graph is simplified to an unweighted one. Finally all maximal cliques of the unweighted graph are generated, each one corresponding to a community model. Θ_r . Then, for each community model, Θ_r i.e., clique, we examine the informativeness of the leaf categories of the initial Web directory that are not in the

clique. Using the OCIA criterion, we compare each such leaf against its closest ancestor that is included in the Θ_r .

6 The *Objective Probabilistic Directory Miner (OPDM)* Algorithm

In the *OCDM* algorithm discussed above, the constructed patterns are based solely on the “observable” behavior of the users, as this is recorded in the usage data. Generally, users’ interests and motives are less explicit. We are considering that the user’s choices are motivated by a number of latent factors that correspond to these subsets. These factors are responsible for the associations between users. The advantage of this approach is that it allows us to describe more effectively the multi-dimensional characteristics of user interests.

A powerful statistical methodology for identifying latent factors in data is Probabilistic Latent Semantic Analysis (PLSA) [4]. Applying PLSA to our scenario of Web directories we consider that there exists a set of user sessions $U = \{u_1, u_2, \dots, u_i\}$, a set of Web directory categories $C = \{c_1, c_2, \dots, c_j\}$, as well as their binary associations (u_i, c_j) which correspond to the access of a certain category c_j during the session u_i . The PLSA model is based on the assumption that each observation of a certain category inside a user session, is related to the existence of a latent factor, z_k that belongs to the set $Z = \{z_1, z_2, \dots, z_k\}$. We define the probabilities $P(u_i)$: the a priori probability of a user session u_i , $P(z_k|u_i)$: the conditional probability of the latent factor z_k being associated with the user session u_i and $P(c_j|z_k)$: the conditional probability of accessing category c_j , given the latent factor z_k . Using these definitions, we can describe a probabilistic model for generating session-category pairs by selecting a user session with probability $P(u_i)$, selecting a latent factor z_k with probability $P(z_k|u_i)$ and selecting a category c_j with probability $P(c_j|z_k)$, given the factor z_k . This process allows us to estimate the probability of observing a particular session-category pair (u_i, c_j) , using joint probabilities as follows:

$$P(u_i, c_j) = P(u_i)P(c_j|u_i) = P(u_i) \sum_k P(c_j|z_k)P(z_k|u_i). \quad (4)$$

Using Bayes’s theorem we obtain the equivalent equation:

$$P(u_i, c_j) = \sum_k P(z_k)P(u_i|z_k)P(c_j|z_k). \quad (5)$$

Equation 5 leads us to an intuitive conclusion about the probabilistic model that we exploit: each session-category pair is observed due to a latent generative factor that corresponds to the variable z_k and hence it provides a more generic association between the elements of the pairs. However, the theoretic description of the model does not make it directly useful, since all the probabilities that we introduced are not available a priori. These probabilities are the unknown parameters of the model, and they can be estimated using the *Expectation-Maximization* (EM) algorithm.

Using the above probabilities we can assign categories to clusters based on the k factors z_k that are considered responsible for the associations between the data. This is realized by introducing a threshold value, named *Latent Factor Assignment Probability*, (*LFAP*) for the probabilities $P(c_j|z_k)$ and selecting those categories that are above this threshold. More formally, with each of the latent factors z_k we associate the categories that satisfy:

$$P(c_j|z_k) \geq LFAP. \quad (6)$$

In this manner and for each latent factor, the selected categories are used to construct a new Web directory. This corresponds to a topic tree, representing the community model, i.e., usage patterns that occur due to the latent factors in the data. The *LFAP* criterion is combined with the *OCIA* criterion and the initial Web directory is traversed and each category-node which does not satisfy the *LFAP* and the *PCAT* thresholds is pruned. In this manner and for each latent factor, the selected categories are used to construct a new Web directory. This corresponds to a topic tree, representing the community model, i.e., usage patterns that occur due to the latent factors in the data.

7 Objective Clustering and Probabilistic Directory Miner (OCPDM)

In an attempt to combine the advantages of clustering and probabilistic modeling, we introduce here a new hybrid method for the discovery of community models. This method combines a clustering algorithm with PLSA. We apply the popular k-means clustering algorithm, for the creation of an initial set of communities. This approach differs from OCPDM clustering, as it produces non-overlapping clusters, i.e., each category belongs to a single cluster. However, as we have explained above for PLSA, the explicit modeling of latent factors is considered advantageous. Thus, we assume that in addition to the k-means clusters, further hidden associations exist in the data, i.e., sub-communities inside each cluster that are not directly observable. To discover this hidden knowledge, we map each cluster derived by k-means onto a new space of latent factors. In this manner, the community Web directories are constructed using a combination of observable and latent associations in the data, and potentially allow us to better model the interests of users. Thus, the new algorithm *Objective Clustering and Probabilistic Directory Miner (OCPDM)* invokes OCPDM for each of the K clusters. The categories of the cluster, on which each latent factor has the maximum impact, are selected using the *LFAP* threshold.

8 Community Web Directory Refinement

The result of the pattern discovery methods presented in Sections 5 to 7 is a set of hierarchies that correspond to the community Web directories, i.e., to a prototypical model for each community, which is representative of the participating

users. The construction of the directory is based on the selection of the categories by each algorithm and their mapping onto the original Web directory. However, the construction of useful community Web directories needs to go beyond the selection of categories by the pattern discovery algorithms. Further processing is required to improve the structure of the directory and this is achieved by the following operators: *Shortcut Operator*. If a category has a single descendant node, then a “shortcut” is created from the parent to the leaf node. *Absorb*. This operator applies to categories that are leaves in the community Web directory, but not in the initial Web directory. Since all of their descendant categories are excluded from the community Web directory, all the Web pages of the eliminated descendant leaves are absorbed by the shrinking category. This operator ensures that no information is lost, even when the “original” leaves are not included in the community Web directory. In the case though, where at least one descendant leaf is included in the community models, this operator is not applied, assuming that the users are not interested in the other leaf categories.

9 Experimental Setup

The evaluation process assessed the effectiveness of the algorithms on categories of the artificial Web directory and on the ODP categories. The evaluation employs mainly two measures: *Coverage* and *User Gain*. Coverage corresponds to the predictiveness of our model, i.e., the number of target Web pages that are actually *covered* by the session-specific community directories. On the other hand, user gain is an estimate of the actual gain that a user would have by following the community Web directory, instead of the initial non-personalized Web directory to get to the desired Web page.

An interesting measure of the effectiveness of our approach is the trade-off between coverage and user gain. The usual choice for such a trade-off measure is the use of Receiver Operating Characteristics (ROC) curves and we plot coverage against (1-User Gain). We name this plot a trade-off curve, since we are not measuring exactly sensitivity and specificity as commonly done in ROC analysis. In Figure 2, we present the trade-off curves for the algorithms for the artificial Web directory and the ODP. From this figure we conclude that the user seems to be benefiting from the personalization, both in terms of Coverage and User Gain. Regarding the comparison of the three directory methods, OPDM clearly outperforms the other two in both directories. The performance of the “hybrid” OCPDM method is lower in terms of coverage, due to the non-overlapping nature of the k-means algorithm.

Comparing the behavior of the methods on the two different directories, we have obtained a higher user gain at a smaller cost in coverage in the ODP directory. In particular, for the ODP directory and for the OPDM algorithm which exhibits the best performance overall, we obtain user gain around 0,50, maintaining coverage at the level of 0,75. For the artificial Web directory, the same algorithm results in a coverage of almost 0,90 but with a user gain value of 0,27. This level of user gain in the ODP directory, is attainable, due to the

size and the generic nature of the ODP. Thus, the use of a real directory has revealed the power of personalizing the directory.

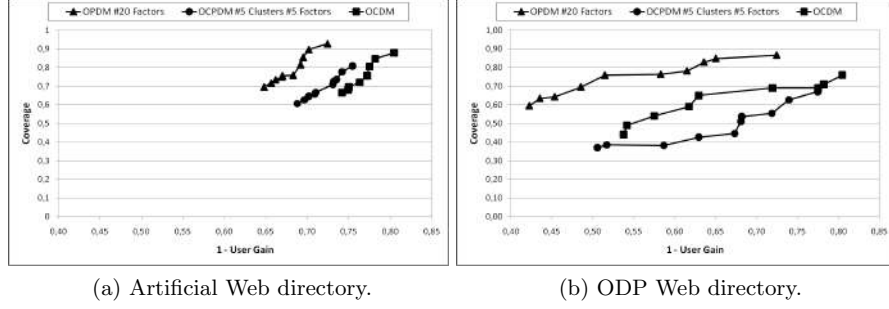


Fig. 2: Web Directory Coverage-User Gain Trade-Off with OCIA (Average PCAT Thresholds).

The results presented in this section provide a detailed picture of the benefits of our approach to personalizing Web directories. Regarding the various discovery methods that we tested, the “pure” PLSA technique (OPDM) outperforms the simple clustering algorithm (OCMD) and the combination of clustering and PLSA (OCPDM).

10 OurDMOZ

In this section, we present *OurDMOZ*, a system that implements and integrates the various components of the proposed methodology. In particular, *OurDMOZ* collects and processes usage data, maps the data onto the Web directory, uses machine learning techniques to extract the community models and finally builds the community Web directories. The main contribution of *OurDMOZ* is that it offers, through its Web application, a personalized view of ODP and consequently a personalized view of the Web. In *OurDMOZ*, a user can join a particular community either by specifying her preferences, or by using the system for some time and letting it decide on the most suitable community models. Thus, there is no requirement for personal information, or other private data, to be provided to the system.

We perform an actual user evaluation, where *OurDMOZ* is given to a set of users who interact with the system and use its personalization functionalities. One of the scenarios followed considered the fill-in of a small questionnaire. The questions included in the questionnaire were answered in a seven-level Likert scale from “Strongly disagree” to “Strongly agree”. The results of the users’ responses to this questionnaire are presented in Figure 3 and show that the users found *OurDMOZ* easy and helpful.

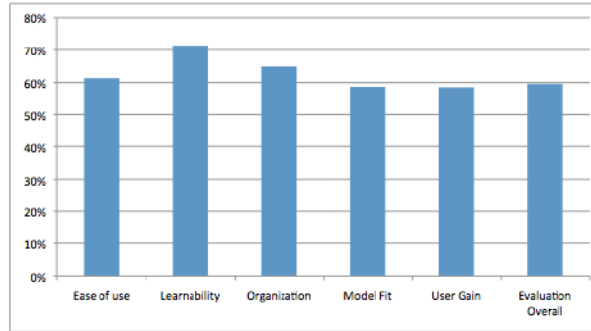


Fig. 3: Questionnaire Results.

11 Conclusions

Community Web directories exemplify a new type of Web personalization, beyond common Web personalization functions such as Web page recommendations and adaptive Web services. In this thesis, we present a complete methodology for the construction of such directories, with the aid of machine learning methods. User community models take the form of thematic hierarchies and are constructed by employing clustering and probabilistic learning approaches.

More specifically, the thesis has contributed in the following areas:

1. Presentation of a roadmap of the Web personalization.
2. Analysis of the main stages of the Web usage mining process and their relation to the Web personalization.
3. Extension of the Web usage mining approach to a much larger portion of the Web, through the analysis of usage data collected by proxy servers. These data correspond to traffic throughout the Web and they are not restricted within the context of a single Web site. The proposed methodology addresses the high dimensionality of the problem, through the classification of individual Web pages onto the categories of the directory.
4. Proposal of three novel pattern discovery algorithms based on clustering and probabilistic approaches (PLSA) for the extraction of community models from the usage data. These methods take into account, not only the browsing behavior of users, but also the structure and the distribution of information within a Web directory.
5. Proposal of a methodology for converting of community models to community Web directories.
6. Development of a complete system, named OurDMOZ, that constructs community Web directories and exploits them for offering personalization functionalities to Web users. These functionalities include not only a customized view of the Web directory, but also they offer recommendation services to Web users.

We hope that this thesis will contribute to the move from Web site personalization, to real Web personalization. In this direction, several issues remain open. These issues are related with all the stages of the community Web directory construction process, such as the exploitation of new techniques that might offer a more accurate view of the users' behavior, whilst respecting user's privacy.

References

1. C. R. Anderson and E. Horvitz. Web montage: A dynamic personalized start page. In *11th WWW Conference*, May 2002.
2. T. Dalamagas, P. Bouros, T. Galanis, M. Eirinaki, and T. Sellis. Mining user navigation patterns for personalizing topic directories. In *9th annual ACM international workshop on Web information and data management*, pages 81–88, 2007.
3. J. A. Hartigan. *Clustering Algorithms*. A Wiley-Interscience Publication, New York: Wiley, 1975, 1975.
4. T. Hofmann. Probabilistic latent semantic analysis. In *UAI*, 1999.
5. T. Oishi, K. Yoshiaki, M. Tsunenori, H. Ryuzo, F. Hiroshi, and M. Koshimura. Personalized search using odp-based user profiles created from user bookmark. In *10th Pacific Rim International Conference on Artificial Intelligence*, pages 839–848, 2008.
6. G. Paliouras, C. Papatheodorou, V. Karkaletsis, and C. D. Spyropoulos. Discovering user communities on the internet using unsupervised machine learning techniques. *Interacting with Computers Journal*, 14(6):761–791, 2002.
7. M. Perkowitz and O. Etzioni. Adaptive web sites: automatically synthesizing web pages. In *Proceedings of the fifteenth national/tenth conference on Artificial intelligence/Innovative applications of artificial intelligence*, pages 727–732. American Association for Artificial Intelligence, 1998.
8. D. Pierrakos and G. Paliouras. Exploiting probabilistic latent information for the construction of community web directories. In *User Modeling*, pages 89–98, 2005.
9. D. Pierrakos and G. Paliouras. Personalizing web directories with the aid of web usage data. *IEEE Transactions on Knowledge and Data Engineering*, 22:1331–1344, 2010.
10. D. Pierrakos, G. Paliouras, C. Papatheodorou, V. Karkaletsis, and M. Dikaiakos. Web community directories: A new approach to web personalization. In B. B. et al., editor, *Web Mining: From Web to Semantic Web, EMWF 2003*, volume 3209 of *LNCS*, pages 113–129. Springer, 2004.
11. D. Pierrakos, G. Paliouras, C. Papatheodorou, and C. D. Spyropoulos. Web usage mining as a tool for personalization: a survey. *User Modeling and User-Adapted Interaction*, 13(4):311–372, 2003.

Static and dynamic graph algorithms with applications to infrastructure and content management of modern networks

Gerasimos G. Pollatos*

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
gpoll@di.uoa.gr

Abstract. This dissertation focuses on networks primarily used for the dissemination of content. We model these networks using graph theory and study algorithms for the replication of content as well as for their infrastructure. In terms of infrastructure, we propose the first fully dynamic algorithm for maintaining a minimum spanning tree on a directed graph. In terms of content replication and focusing on networks with constant number of clients we propose optimal algorithms that solve the basic problem of replicating data over a network of clients with constrained local storage and extend our results to various important extensions. These algorithms constitute the first research work on networks with non-metric distances among clients. In addition and in order to study implications of the selfish behaviour of users participating in such networks, we define an appropriate non-cooperative strategic game and study existence of pure Nash equilibria as an indication of stabilization of the networks performance. We provide tight bounds for the prices of anarchy and stability, the standard measures of efficiency of such networks. We identify conditions under which equilibria might not be expensive and extend our results to more complex hierarchical networks.

1 Introduction

Typically speaking the term *content network* refers to a large-scale system of computers containing copies of data, placed at various points in a network so as to maximize bandwidth for access to the data from clients throughout the network. A client accesses a copy of the data near to the client, as opposed to all clients accessing the same central server, so as to avoid bottleneck near that server. This *maximization* of the bandwidth dedicated for data access can also be seen as a *minimization* of the access cost each client has to *pay* in order to gain access to desired content.

Content types include ordinary downloadable objects, such as media files, software and documents, as well as web objects and applications. Other components of internet delivery, such as DNS queries, routes and database queries, can

* Dissertation Advisor: Vassilis Zissimopoulos, Professor

also be considered as content. In order to refer to such content, we will use the abstract notion of a data object or simply *object*.

The most fundamental problem a designer of a large-scale content network has to face is the *quality of service*. A system is meaningful if it achieves to serve the purpose it was built for, and succesful if it manages to accomplish that as best as possible. The quality of service relies among others, in two major factors; reliability and performance. Reliability amounts to maintaining the following invariant: "problems arising in the system should have the smallest possible impact, if any, on the offered services". That is, the clients of a system should ideally be unaware of problems related to server availability issues, or sudden power failures in some of the system's database locations. On the other hand, performance is a more complicated issue which involves the basic aims of the network. If speed of access is the aim, then optimization of performance is equivalent to fast access rates, while if up-to-date content availability is the goal, then performance amounts to be able to efficiently update content in all cooperating locations.

In this thesis, we address from an algorithmic perspective both reliability and performance in content networks. We utilize the field of dynamic graph algorithms in order to surpass difficulties in terms of reliability. As for performance, we turn our attention to *cooperative* caching of objects among distributed computers and design algorithms for their coordination by external authorities as well as study implications in their operability due to absence of any external authority.

2 Infrastructure on content networks: the DMST problem

In terms of infrastructure, we study connectivity issues among nodes of the content network. We use graph theory and model the content network as a graph, with each vertex of the graph representing a single user and each edge representing a connection among two users. We assume that each edge is weighted and directed.

In the environment of a content network, the derived graph is obviously subject to discrete changes. Such changes happen if for example a user voluntarily exits the content network, or if a sudden power or hardware failure forces the user to abandon the network. Using graph theoretic terms this is a dynamic graph. We focus on the classic problem of computing a directed minimum spanning tree (DMST) on such a graph. Such a tree is a maximal subset of the edges of the graph, with minimum cost that is (a) acyclic, i.e. does not contain directed cycles, and (b) no vertex of the directed graph has more than one incoming edge in this subset. The need for maintaining such a tree in content networks is apparent due to the fact that the spanning tree essentially maintains a minimum set of links among network nodes (i.e. the minimum requirements for connectivity of all nodes).

The dynamic version of the problem amounts to exploiting a previously computed solution so that we can make use of it after a node failure, without having to recompute the new solution from scratch. The problem is studied for the first time here as opposed to the widely studied undirected version. A relevant result proved in this thesis for the hardness of the problem is the following:

Theorem 1. *[1] Dynamic maintenance of a DMST under edge deletions and/or insertions is as hard as recomputing a DMST from scratch if only the DMST information is retained and used between updates.*

The dynamic algorithm for the DMST is based on an appropriate data structure, which the algorithm utilizes for recording redundant edges and all vertices during the execution of the only known algorithm for the static version of the problem i.e. Edmonds' algorithm [2], [3], [4]. The augmented data structure, namely the *Augmented Tree* (ATree), appropriately encodes the redundant edge set H along with all vertices (contracted vertices or c-vertices and simple vertices) processed during execution of Edmonds' algorithm. Simple vertices are represented in the ATree by *simple nodes* while c-vertices are represented by *c-nodes*.

Since a digraph can be always transformed to be strongly connected, all vertices (simple and intermediate c-vertices) will be eventually contracted to a single c-vertex by the end of the algorithm's execution. This c-vertex is represented by the root node of the ATree, which has no parent. The parent of each other node N is the intermediate c-node N^c to which it was contracted. The parent of each node is unique, hence the described structure is indeed a tree.

The ATree has at most $O(n)$ nodes since the algorithm handles at most $O(n)$ contractions. Construction of an ATree can be easily embedded into the functionality of Edmonds' algorithm, without affecting its complexity. Furthermore, all intermediate c-vertices created by Edmonds' algorithm are made explicit in the ATree: For a given c-vertex v_c , we can engage a Breadth-First Search (BFS) starting from its representing c-node in the ATree and collect all encountered simple nodes, hence we gather all vertices of the original digraph that were eventually contracted to v_c . Since the ATree is of $O(n)$ size, BFS takes $O(n)$ time.

For any edge we want to delete from the original digraph, two cases must be considered: (a) the edge does not belong to H , in which case only a simple update on the recorded lists is needed, and (b) the edge belongs in H , in which case we proceed by *decomposing* the ATree, initialize Edmonds' algorithm w.r.t. the remainders of the ATree and execute it.

The purpose of engaging a decomposition of the ATree is to identify surviving c-nodes and hence surviving c-vertices. By this way we can re-execute Edmonds' algorithm on a partially contracted digraph with less vertices, considering less edges than re-evaluating all contractions from scratch.

The decomposition of the ATree begins from node N which is the head of the deleted edge and follows a path from N towards the ATree root, removing all c-nodes on this path except N . Each one of the children of a removed c-node forms its own subtree. By the end of this procedure, the initial structure

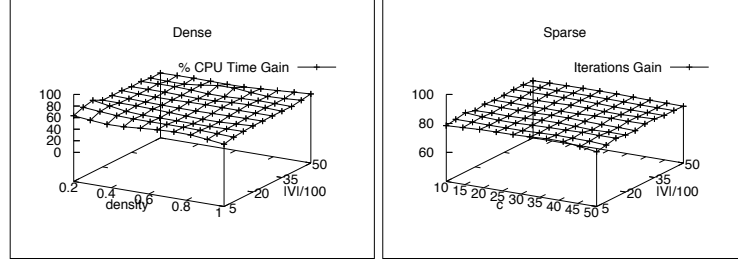


Fig. 1. Experimental evaluation of the dynamic algorithm for the DMST problem.

is decomposed into smaller ATrees, each corresponding to a contracted subset of the original digraph's vertices. All these ATrees remain intact after decomposition. Having performed the decomposition of the ATree, we proceed with a proper initialization and re-execution of Edmonds' algorithm, on a new partially contracted digraph.

We handle edge insertions by reducing them to edge deletions. We first check whether the newly inserted edge should replace some edge encoded in the ATree. This check involves *only* c -nodes of the ATree corresponding to c -vertices containing the head of the new edge. Given that we have found a candidate node N which should replace its current incoming edge with the new one, we proceed by engaging a virtual deletion of the old edge and re-execution of the algorithm on the remaining graph augmented with the newly inserted edge.

To analyze the theoretical performance of our algorithm we used the widely known output complexity model [5]. If we denote by ρ the set of changes in a previously computed solution (the number of affected constituents), made by an algorithm as a response to an update in the dynamic graph, then $\Omega(|\rho|)$ is a lower bound on the complexity of an update, where $|\rho|$ is the cardinality of the set ρ . In the output complexity model, the complexity of updates for a dynamic algorithm is measured as an asymptotic upper bound of a function of $|\rho|$, where the set of affected vertices is ρ , $|\rho|$ being its size and the cardinality of affected edges is denoted with $||\rho||$. Our basic result is:

Theorem 2. [1] *Fully dynamic maintenance of a directed minimum costs spanning tree can be achieved in $O(n + ||\rho|| + |\rho| \log |\rho|)$ output complexity per edge operation.*

In order to measure the practical performance, we employed our algorithm on sequences of edge operations on digraphs of varying order and density and recorded and compared the following two measures: (a) **average CPU time**, needed to complete computation of the new DMST and (b) **number of iterations**, performed by each one of the two algorithms. Sample experimental results are depicted in figure 1.

Our dynamic algorithm substantially achieves an update time reduced by a factor of more than 2 as opposed to solving the problem statically (from scratch)

on dense digraphs. However it is our belief that the case of sparse digraphs merits further theoretical investigation from an average case complexity perspective, since as the results indicate, there appears to be an asymptotic improvement on average.

3 Cooperative content replication

One of the most fundamental techniques for improving the performance of a large-scale content network is to cache popular objects close to the clients that potentially request them. This enables requests to be satisfied by a nearby copy and hence reduces not only the access latency but also the burden on the network as well as the remote servers, offering the objects. In the simplest caching scheme, nodes operating as caches never consult one another and when a cache miss occurs, the server is contacted directly. Improvement of the effectiveness of caching is usually accomplished through a powerful paradigm; cooperation. Nodes cooperate both in serving each other's requests as well as in making storage decisions.

Roughly speaking, the problem we study is the following: given a set of cooperating nodes, the amount of local storage at each node, the network distances between the nodes and the predictions of access rates from each node to a set of objects, determine where to place the copies of each object in order to minimize the total access cost over all nodes. We assume that an underlying mechanism exists that can route each node to the closest other node on the network that holds a replica of an object, when an object miss occurs.

From an optimization point of view, this problem is NP-hard since it is a direct generalization of the well known uncapacitated facility location problem where multiple different types of facilities are considered. In this dissertation we focused on the case where the number of cooperating nodes is constant and designed algorithms that efficiently address this problem.

The data placement problem ([6], [7]), is abstracted as follows: there is a network \mathcal{N} consisting of a set \mathcal{M} of $M = |\mathcal{M}|$ users (also referred to as clients) and a universe \mathcal{O} of $N = |\mathcal{O}|$ objects. Each object $o \in \mathcal{O}$ has length l_o and each user $j \in \mathcal{M}$ has a local capacity C_j for the storage of objects. The distance between the users can be represented by a distance matrix D , not necessarily symmetric, where d_{ij} denotes the distance from j to i . The matrix D models the underlying topology. In our work we do not assume any restrictions on the distances, e.g. metric, which is usually the case in the literature.

Each user i requests access to a set of objects $R_i \subseteq \mathcal{O}$, namely its *request set*. For each object o in its request set, client i has a *demand* of access $w_{io} > 0$. This demand can be interpreted as the frequency under which user i requests object o or even as the *preference* that i has for object o . The subset P_i of its request set, that i chooses to replicate locally is referred to as its *placement*. Obviously, $|P_i| \leq C_i$ for unit-sized objects. We assume an installation cost f_i^o for each object o and each cache (user) i .

The objective of the data placement problem is to choose placements of objects for every client such as the total induced access and installation costs for all objects and all clients is minimized. In the following, we will make the reasonable assumption that each object $o \in \mathcal{O}$ is requested by at least one user. If this is not the case, one can always formulate and solve an equivalent problem which has an object set containing only requested objects.

Up to know, the only way to tackle this problem, was the 10-approximation algorithm of [7] designed for the general case. When object lengths are uniform (or equivalently unit) our algorithm finds the optimum solution in polynomial time. When object lengths are non-uniform, our algorithm returns an optimum solution which violates the capacities of the clients' caches by a small, asymptotically tight additive factor.

We show how our results can be modified to handle various important extensions of the problem such as cases when bounds are imposed on the number of maximum replicas allowed for each object (a k -median variant for DP), or cases when upper and lower bounds are imposed on the number of users a single replica of an object can serve. Furthermore, we describe ways to cope with a well known generalization of DP, the page placement problem ([8]), in which bounds are imposed on the number of clients that can connect to a client's cache, as well as cases where object updates are frequent and consistency of all replicas of each object has to be guaranteed.

The latter problem is more commonly known as the connected data placement problem [7]. Our algorithms are applicable with uniform and non-uniform object lengths and can be employed independently of the underlying topology of the network, thus giving the first non-trivial results for non-metric DP problems. Most of the results described in the dissertation appear in [9].

Our basic result for objects of uniform (equivalently unit) length is the following:

Theorem 3. [9] *The non-metric data placement problem with uniform length objects and a fixed number of clients can be solved optimally in polynomial time $O(N^{M+1})$.*

The algorithm we designed is based on dynamic programming algorithm and is in fact pseudo-polynomial, since the complexity depends on the maximum cache size. In the case of unit-sized objects we are able to bound it by the total number of objects and thus obtain a polynomial time algorithm. In the case of objects of arbitrary length the same bound does not hold and the algorithm remains pseudo-polynomial. To tackle this issue, we let $\alpha = \varepsilon l_{max}/N$ where ε is an arbitrarily small positive constant and modify the object lengths and cache sizes appropriately as $l'_o = \lfloor \frac{l_o}{\alpha} \rfloor$ and $C'_j = \lfloor \frac{C_j}{\alpha} \rfloor$ respectively. To compute a solution we use the same algorithm with the modified sizes and obtain the following result:

Theorem 4. [9] *The non-metric data placement problem, with non-uniform object lengths and a fixed number of clients, can be solved optimally in polynomial*

	Known results	In this paper
	arbitrary M , metric	fixed M , no metric
uniform lengths DP	APX-hard, 10-approx [6, 7, 10]	optimal
uniform lengths with installation costs DP	APX-hard, 10-approx [6, 7, 10]	optimal
connected DP	14-approx [7]	optimal
k -median DP	10-approx [6, 7]	optimal
page placement	13-approx [10] *	optimal with cache augmentation ** εl_{max}
non-uniform lengths DP	10-approx with cache augmentation l_{max} [7]	optimal with cache augmentation εl_{max}

Table 1. The main known results on the DP problem (*non-uniform lengths with constant blow-up on both capacities, **on cache capacity only).

time using εl_{max} augmentation on the machines' capacity, where ε is an arbitrarily small positive constant and l_{max} is the length of the largest object.

We note that the augmentation in each user's cache stated in the above theorem is asymptotically tight. Furthermore, our results can be modified to handle well-known and important extensions of the DP problem with constant number of clients. The extensions we consider are the following: (1) existence of a distant server (user) whose cache can hold all the objects, (2) there is an upper bound k_o imposed on the number of copies of an object o that are replicated in the network, (3) each user j has a set RJ_j of *rejected* objects that cannot be placed on its cache, (4) the number of users u_{o_i} served by a single copy o_i of object o is lower bounded by $k_{o_i}^{min}$ and upper bounded by $k_{o_i}^{max}$, (5) there is an upper bound k_j on the number of users that can access a given user j 's cache, (6) objects updates are frequent and all copies of an object must be up-to-date. We summarize the main results on the DP problem and these extensions, in table 1.

4 Distributed selfish replication

While cooperation of nodes is an attractive and reasonable paradigm in environments where machines trust one another, such as within an Internet service provider, a cache service provider or even a corporate intranet, there are numerous content networks under which this assumption is not valid. For example, file sharing or peer-to-peer networks, that have become extremely popular nowadays, are formed by users that participate voluntarily and have no knowledge of the motives and goals of other participants.

In such networks, one cannot a-priori assume that a participant will voluntarily offer it's local storage for the replication of content that is of no interest to it. On the contrary the most logical assumption is that such users upon joining the network behave selfishly, i.e. aim to maximize their own benefit and thus

the total access cost depends on their selfish replication decisions. The two basic questions that arise under this setting are: (a) does the performance of the network ever stabilize? (b) what is the overall performance of a stable network, in absence of a central optimizing authority?

Using tools from the field of algorithmic game theory, we studied implications of this user behaviour in content networks. We formulated the basic problem of data replication as a strategic non-cooperative game and study this game in order to analyze the performance of these networks. We use the standard quantification measures, the prices of anarchy [11] and stability [12], [13], in order to measure the performance loss due to the lack of coordination.

The modelled strategic game is as follows [14], [15]:

Definition 1. *An instance of the Distributed Selfish Replication (DSR) game is specified by the tuple $\langle N, \{P_i\}, \{c_i\} \rangle$, where:*

- N is the set of the n players, which in our case are the nodes.
- $\{P_i\}$ is the set of strategies available to player i . Each strategy corresponds to a different placement which means that there is one strategy for each $size_i$ -cardinality subset.
- $\{c_i\}$ is the set of utilities for the players. The utility for each player is the access cost that the player wishes to minimize.

A placement X is then a *strategy profile* (or *configuration*) for the DSR game. DSR is a n -player, non-cooperative, non-zerosum game [16]. For the DSR game, we investigate configurations that are pure Nash equilibria.

Definition 2. *A pure Nash equilibrium (N.E.) for the DSR game is a placement X^* , such that for every node $i \in N$,*

$$c_i(X^*) \geq c_i(\{P_1^*, \dots, P_{i-1}^*, P_i, P_{i+1}^*, \dots, P_n^*\})$$

for all $P_i \in \{P_i\}$.

The definition essentially states that, under such a placement X^* , no node can modify its individual placement unilaterally and benefit from this modification. In what follows we use the terms node and player interchangeably. Let

$$X_{-i} = \{P_1, \dots, P_{i-1}, P_{i+1}, \dots, P_n\}$$

refer to the strategy profile of all players but i . For the DSR game, it is easy to see that given a strategy profile X_{-i} , player i can determine optimally in polynomial time its *best response* P_i to the other players' strategies. This computation amounts to solving a special 0/1 *Knapsack* problem, in which player i chooses to replicate the $size_i$ objects with the greatest value $w_{io}d_i(o)$.

The network topologies we consider and the results we obtained are the following. The minimum and maximum distances appearing in the network are also referred to with $d_{\min} = d_k$ and $d_{\max} = d_0$.

Ultrametric Replication Group: this is a network model that generalizes the one introduced by Leff, Wolf and Yu in [17] and studied in [18], involving k origin servers, instead of 1. The distance of every node i from server l is d_l for $l = 0 \dots k-1$, while two nodes i and j are at distance $d_{ij} = d_k$. In our study we assume that distances form an ultra-metric ¹ such that $d_k < d_{k-1} < \dots < d_0$. We designed a distributed protocol that upon finishing guarantees convergence to pure Nash equilibria.

Theorem 5. [14] *Pure strategy Nash equilibria exist for the DSR game on $LWY(k)$ networks, and can be found in polynomial time.*

Furthermore we studied the quality of such equilibria obtaining the following results:

Theorem 6. [14] *The price of anarchy for the DSR game is upper bounded by $\frac{d_{max}}{d_{min}}$. The Price of Stability for the DSR game has a lower bound arbitrarily close to d_{max}/d_{min} in the worst-case, even for 1 server and 0/1 demand weights.*

We should note that both results are valid even when the maximum experienced access cost by a user is measured instead of the sum of all access costs. Furthermore we studied networks with modestly demanding participants, that is participants offering storage capacity to the network asymptotically equal to their demand, and identified that in such networks pure equilibria can be of significantly less cost.

Balanced Hierarchies: this is a network model with distances that also form an ultrametric. The network's nodes are clustered hierarchically, so that at each clustering level the maximum distance between any two nodes in the same cluster is given and is smaller than the distance of any two nodes in different clusters. We designed an extension of our basic distributed protocol that also achieves convergence to pure Nash equilibria. Our main result on the quality of the pure equilibria is the following.

Theorem 7. [15] *The Price of Anarchy of the DSR game with 0/1 preference weights on balanced hierarchical networks is $O\left(\frac{\ln n}{\ln \ln n}\right)$. The Price of Stability of the DSR game on hierarchical networks with 0/1 preferences is $\Omega\left(\frac{\ln n}{\ln \ln n}\right)$.*

General Networks: in this case the distance matrix $[d_{ij}]$ can be arbitrary. Our study shows that pure equilibria are not guaranteed to always exist when such network topologies are considered.

Theorem 8. [15] *The DSR game on general networks is not a potential game.*

¹ An ultrametric is a metric which satisfies the strengthened version of the triangle inequality, $d(x, z) \leq \max\{d(x, y), d(y, z)\}$ for all x, y, z . This essentially states that at least two of the $d(x, z)$, $d(x, y)$ and $d(y, z)$ are the same.

5 Conclusions

This thesis investigated problems arising in content networks and described ways to effectively cope with them. Focusing on infrastructure problems arising in content networks and using tools from the field of dynamic graph algorithms we studied the problem of maintaining a directed minimum spanning tree on a graph that changes dynamically under edge insertions and deletions. After analyzing the hardness of maintaining such a tree, we described the first fully dynamic algorithm that maintains the DMST under edge updates and analyzed it in the output complexity model. The results of the extensive experimental evaluation, revealed the practical efficiency of the proposed algorithm. An important aspect of further research is to resolve the complexity of updates, a small step to which has been made in this thesis.

In terms of content replication we addressed the basic problem of replicating data among users of a content network, the problem most commonly referred to as the data placement problem. We focused on the case of constant number of users and described polynomial time algorithms that solve optimally the basic problem when the objects in consideration have uniform size. When object sizes are not uniform we described an algorithm that also finds the optimum solution to the problem, albeit a small and asymptotically tight augmentation in each user's local storage. The proposed technique was extended to handle various other common extensions of the basic problem such as the page-placement problem and the connected data placement problem among others. A significant characteristic of this technique is its ability to solve these problems independently of the underlying topology of the network thus giving the first non-trivial results for non-metric topologies. A challenging aspect of future work involves employing our results to models involving payments.

Finally, we studied implications on the process of replicating data over a content network when the users are autonomous and selfish and participate voluntarily. We formulated a proper strategic game that modeled the data replication problem and studied conditions under which the network stabilizes. We proved inability to converge to pure Nash equilibria for general underlying topologies. For simple hierarchical networks and balanced hierarchical networks with multiple servers we described an algorithm that reaches pure Nash equilibria and admits a distributed implementation. Furthermore, we analyzed the quality of achieved equilibria by computing the prices of anarchy and stability for the game and identified conditions under which these ratios can be different. The most significant aspect of future work is investigation of whether part or all of our results can be extended in the case of arbitrary-sized objects.

References

1. Gerasimos G. Pollatos, Orestis Telelis, and Vassilis Zissimopoulos. Updating directed minimum cost spanning trees. In *Workshop on Experimental Algorithms (WEA)*, pages 291–302, 2006.

2. F. Bock. An algorithm to construct a minimum spanning tree in a directed network. *Developments in Operations Research*, pages 29–44, 1971.
3. Y.J. Chu and T.H. Liu. On the shortest arborescence of a directed graph. *Scientia Sinica*, 14:1396–1400, 1965.
4. J. Edmonds. Optimum branchings. *Journal of Research of the National Bureau for Standards*, 69(B):125–130, 1967.
5. G. Ramalingam and Thomas W. Reps. On the computational complexity of dynamic graph problems. *Theoretical Computer Science*, 158(1&2):233–277, 1996.
6. I. Baev and R. Rajaraman. Approximation algorithms for data placement in arbitrary networks. In *Proceedings of the ACM-SIAM Annual Symposium on Discrete Algorithms (SODA)*, pages 661–670, 2001.
7. Ivan D. Baev, Rajmohan Rajaraman, and Chaitanya Swamy. Approximation algorithms for data placement problems. *SIAM Journal on Computing*, 38(4):1411–1429, 2008.
8. Adam Meyerson, Kamesh Munagala, and Serge Plotkin. Web caching using access statistics. In *Proceedings of the ACM-SIAM Annual Symposium on Discrete Algorithms (SODA)*, pages 354–363, 2001.
9. Eric Angel, Evripidis Bampis, Gerasimos G. Pollatos, and Vassilis Zissimopoulos. Optimal data placement on networks with constant number of clients. *Submitted to COCOON 2010*, 2010.
10. Sudipto Guha and Kamesh Munagala. Improved algorithms for the data placement problem. In *Proceedings of the ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 106–107, 2002.
11. Elias Koutsoupias and Christos H. Papadimitriou. Worst-case Equilibria. In *Proceedings of the Symposium on Theoretical Aspects of Computer Science (STACS)*, pages 404–413, 1999.
12. E. Anshelevich, A. Dasgupta, E. Tardos, and T. Wexler. Near-optimal network design with selfish agents. In *Proceedings of the ACM Annual Symposium on Theory of Computing (STOC)*, pages 511–520, 2003.
13. E. Anshelevich, A. Dasgupta, J. M. Kleinberg, E. Tardos, T. Wexler, and T. Roughgarden. The Price of Stability for Network Design with Fair Cost Allocation. In *Proceedings of IEEE Annual Symposium on Foundations of Computer Science (FOCS)*, pages 295–304, 2004.
14. Gerasimos G. Pollatos, Orestis Telelis, and Vassilis Zissimopoulos. On the social cost of distributed selfish content replication. In *Proceedings of IFIP-TC6 Networking*, pages 195–206, 2008.
15. Gerasimos G. Pollatos, Orestis Telelis, and Vassilis Zissimopoulos. The social cost of distributed selfish content replication. *Submitted to Computer Communications*, 2010.
16. M. J. Osborne and A. Rubinstein. *A course in game theory*. MIT Press, 1994.
17. A. Leff, J. Wolf, and P. S. Yu. Replication Algorithms in a Remote Caching Architecture. *IEEE Transactions on Parallel and Distributed Systems*, 4(11):1185–1204, 1993.
18. N. Laoutaris, O. A. Telelis, V. Zissimopoulos, and I. Stavrakakis. Distributed Selfish Replication. *IEEE Transactions on Parallel and Distributed Systems*, 17(12):1401–1413, 2006.

Algorithms for Space-Time Equalization of Wireless Channels

Constantinos Rizogiannis *

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
`krizog@di.uoa.gr`

Abstract. In this thesis we investigate receiver techniques for maximum likelihood (ML) joint channel/data estimation in flat fading multiple-input multiple-output (MIMO) channels. The performance of iterative least squares (LS) for channel estimation combined with sphere decoding (SD) for data detection is examined for block fading channels, demonstrating the data efficiency provided by the semi-blind approach. The case of continuous fading channels is addressed with the aid of recursive least squares (RLS). The observed relative robustness of the ML solution to channel variations is exploited in deriving a block QR-based RLS-SD scheme. For the multi-user MIMO scenario, the gains from exploiting temporal/spatial interference color are assessed. We also derive the optimal training sequence for ML channel estimation in the presence of co-channel interference (CCI). In the second part of the thesis we propose two new adaptive equalizers for direct sequence code division multiple access (DS-CDMA) systems operating over time-varying and frequency selective channels. The equalizers consist of a number of serially connected stages and detect users in an ordered manner, applying a decision feedback equalizer (DFE) at each stage. Both the equalizer filters and the order in which the users are extracted are updated in a RLS manner, efficiently realized through time- and order-update recursions.

1 Introduction

During the last two decades, there has been an explosion in the services offered by wireless telecommunication networks, which is boosted from the relevant growth of the technologies of Informatics and Telecommunications. At the same time, there are new challenges for the development of the next generation telecommunication systems. Two of the most basic technologies for the evolution of the new services in the wireless networks are the multiple-input multiple-output (MIMO) and the code division multiple access (CDMA) systems. In this PhD thesis, we worked on the design and analysis of space-time signal processing algorithms for this kind of systems.

Specifically, we investigate iterative and recursive least squares (LS) algorithms for maximum-likelihood (ML) joint channel/data estimation, that are

* Dissertation Advisor: Sergios Theodoridis, Professor

both data efficient and computationally attractive. The proposed schemes use the sphere decoding (SD) algorithm for data detection, and short training sequences for an initial channel estimation. We studied three new algorithms [10] for block- and continuous-fading frequency flat MIMO channels. Moreover, in the case of DS-CDMA systems, we propose two new adaptive equalizers of the successive interference cancellation (SIC) type operating over time-varying and frequency selective channels. Their development relies on the formulation of a DS-CDMA system as one with multiple inputs and multiple outputs and the adoption of existing adaptive solutions of the BLAST-type for MIMO channel equalization [1, 12, 7].

2 Semi-blind maximum-likelihood joint channel/data estimation for correlated channels in multiuser MIMO networks

2.1 Signal and System Model

Consider a MIMO communications system, with M_T transmit and M_R receive antennas, where $M_R \geq M_T$, and frequency flat fading channels. The received signal vector at time n is given by

$$\mathbf{x}(n) = \mathbf{H}_0(n)\mathbf{s}_0(n) + \mathbf{v}(n) \quad (1)$$

where $\mathbf{H}_0(n) \in \mathbb{C}^{M_R \times M_T}$ is the channel matrix, assumed of full column rank, $\mathbf{s}_0(n) \in \Omega^{M_T \times 1}$ denotes the input signal vector taking values from a finite alphabet (FA) Ω with cardinality $Q = |\Omega|$, and $\mathbf{v}(n) \in \mathbb{C}^{M_R \times 1}$ is composed of colored interference (CCI) and additive, temporally and spatially white, zero mean Gaussian noise.

2.2 Single-User Case

Maximum Likelihood Estimation. In the absence of multiuser interference, $\mathbf{v}(n)$ in (1) is only composed of white Gaussian noise. Thus, the problem of ML estimation can be formulated as

$$\min_{\mathbf{s}_0(n) \in \Omega^{M_T \times 1}, \mathbf{H}_0(n) \in \mathbb{C}^{M_R \times M_T}} \|\mathbf{x}(n) - \mathbf{H}_0(n)\mathbf{s}_0(n)\|^2 \quad (2)$$

It is clear that, given the input data $\mathbf{s}_0(n)$, the solution for the channel $\mathbf{H}_0(n)$ is given by its least squares (LS) estimate. For a known channel, the ML-optimal input vector is to be searched among all Q^{M_T} candidate M_T -tuples from $\Omega^{M_T \times 1}$. Sphere decoding (SD) [2] is known to be a computationally efficient alternative to exhaustive enumeration [4]. The basic idea is to reduce the number of candidates by searching only within a hypersphere centered at $\mathbf{x}(n)$ using a QR decomposition (QRD) of the channel matrix.

Block Fading. Assuming block fading and dropping time indices, (1) can be re-written as

$$\mathbf{X} = \mathbf{H}_0 \mathbf{S}_0 + \mathbf{V} \quad (3)$$

where \mathbf{X} denotes the $M_R \times N$ output matrix, \mathbf{S}_0 is the $M_T \times N$ input matrix, and $\mathbf{V} \in \mathbb{C}^{M_R \times N}$ is the noise matrix. N denotes the length of the data block, over which the channel matrix is assumed constant. Let $\mathbf{H}_0^{(0)}$ denote the estimate of \mathbf{H}_0 that may have resulted from a (short) training period. This can be improved, and as a consequence the data estimates as well, via an iterative procedure consisting of alternately optimizing the data estimate based on the current channel estimate and vice versa. Table 1 summarizes the general ALS ML scheme, where $\mathbf{H}_0^{(i)}$ and $\mathbf{S}_0^{(i)}$ are the channel and data estimates in the i -th

Table 1. ALS for joint ML channel estimation/data detection.

Given:	$\mathbf{X}, \mathbf{H}_0^{(0)}$
Step 0	$i = 0$
	Repeat until convergence
Step 1	$i = i + 1$
Step 2	$\mathbf{S}_0^{(i)} = \arg \min_{\mathbf{S}_0 \in \Omega^{M_T \times N}} \ \mathbf{X} - \mathbf{H}_0^{(i-1)} \mathbf{S}_0\ ^2$
Step 3	$\mathbf{H}_0^{(i)} = \mathbf{X} \mathbf{S}_0^{(i)\dagger}$

iteration. Two well-known examples are *iterative least squares with projection (ILSP)* and *iterative least squares with enumeration (ILSE)* [16]. ILSP is a simple approach, where the ML solution is only approximated, by projecting onto the FA each of the entries of the soft LS data estimate, $\mathbf{H}_0^{(i-1)\dagger} \mathbf{X}$. In ILSE, full enumeration is performed, thus obtaining the exact ML solution but at a very high computational cost. An exact ML ALS scheme of lower expected complexity would result if SD were utilized instead to detect each of the columns of $\mathbf{S}_0^{(i)}$ in Step 2 above. This algorithm will henceforth be referred to as *iterative least squares with SD (ILS-SD)*.

Continuous Fading. For continuous fading channels an online version of the ILS-SD algorithm is employed using the recursive least squares (RLS) algorithm. Considerable computational savings would result if the Q, R factors were tracked instead of the channel matrix itself [8]. An additional reduction in the computational burden of the receiver can be achieved by performing the Q, R update once in every T samples instead of on a sample by sample basis. Between two consecutive updates, SD is based on the available QRD as if the channel remained constant in the meantime. This ‘sub-sampled’ channel tracking scheme is suggested by the observed robustness of the ML solution to mild changes in the channel [10] and leads to significant computational savings with little or no performance loss. The proposed algorithm, will be referred to hereafter as RLS-SD.

2.3 Multi-User Case

Interference Color. Here, we consider the case where the interference component, $\{\mathbf{v}(n)\}$, may be correlated both in time and space. Therefore the received signal process $\{\mathbf{x}(n)\}$ is also temporally correlated. To exploit this fact, we employ more than one consecutive received samples to *jointly* detect the corresponding input vectors [14]. Stacking N consecutive received samples together we can write:

$$\underbrace{\begin{bmatrix} \mathbf{x}(n) \\ \mathbf{x}(n-1) \\ \vdots \\ \mathbf{x}(n-N+1) \end{bmatrix}}_{\bar{\mathbf{x}}} = \underbrace{\begin{bmatrix} \mathbf{H}_0 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{H}_0 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{H}_0 \end{bmatrix}}_{\bar{\mathbf{H}}_0} \underbrace{\begin{bmatrix} s_0(n) \\ s_0(n-1) \\ \vdots \\ s_0(n-N+1) \end{bmatrix}}_{\bar{\mathbf{s}}_0} + \underbrace{\begin{bmatrix} \mathbf{v}(n) \\ \mathbf{v}(n-1) \\ \vdots \\ \mathbf{v}(n-N+1) \end{bmatrix}}_{\bar{\mathbf{v}}}$$

or

$$\bar{\mathbf{x}} = \bar{\mathbf{H}}_0 \bar{\mathbf{s}}_0 + \bar{\mathbf{v}} \quad (4)$$

All the interferers' channels are assumed to obey the well-known Kronecker model [13, 18], with the same receive fading correlation matrix. Assume, moreover, an *interference-limited environment*, where CCI overwhelms background noise [14]. One can then show that the interference correlation matrix is given by [17]

$$\mathcal{R}_{\bar{\mathbf{v}}} = E(\bar{\mathbf{v}}\bar{\mathbf{v}}^H) = \mathcal{R}_{\mathbf{t}}^* \otimes \mathcal{R}_{\mathbf{s}} \quad (5)$$

where $\mathcal{R}_{\mathbf{t}}^*$ and $\mathcal{R}_{\mathbf{s}}$ are the *temporal* and the *spatial* interference colors, respectively.

Maximum-Likelihood Estimation. Under the assumption of Gaussianity for $\{\mathbf{v}(n)\}$ [15], the ML joint channel estimation / data detection problem for (4) can be formulated as

$$\min_{\bar{\mathbf{s}}_0, \bar{\mathbf{H}}_0} \left[\mathcal{R}_{\bar{\mathbf{v}}}^{-1/2} (\bar{\mathbf{x}} - \bar{\mathbf{H}}_0 \bar{\mathbf{s}}_0) \right]^H \left[\mathcal{R}_{\bar{\mathbf{v}}}^{-1/2} (\bar{\mathbf{x}} - \bar{\mathbf{H}}_0 \bar{\mathbf{s}}_0) \right]$$

where $\mathcal{R}_{\bar{\mathbf{v}}}^{-1/2}$ is a Hermitian square root of $\mathcal{R}_{\bar{\mathbf{v}}}^{-1}$. Utilizing the relation $\mathcal{R}_{\bar{\mathbf{v}}}^{-1/2} = \mathcal{R}_{\mathbf{t}}^{-*/2} \otimes \mathcal{R}_{\mathbf{s}}^{-1/2}$, resulting from (5), the ML problem for \mathbf{H}_0 and \mathbf{S}_0 is now formulated as

$$\min_{\mathbf{S}_0, \mathbf{H}_0} \left\| \mathcal{R}_{\mathbf{s}}^{-1/2} \mathbf{X} \mathcal{R}_{\mathbf{t}}^{-1/2} - \mathcal{R}_{\mathbf{s}}^{-1/2} \mathbf{H}_0 \mathbf{S}_0 \mathcal{R}_{\mathbf{t}}^{-1/2} \right\|^2 \quad (6)$$

where $\bar{\mathbf{x}} = \text{vec}(\mathbf{X})$ and $\bar{\mathbf{s}}_0 = \text{vec}(\mathbf{S}_0)$. Hence, the solution is given by the Gauss-Markov estimator (GME) [5]:

$$\hat{\mathbf{H}}_0 = \mathbf{X} \mathcal{R}_{\mathbf{t}}^{-1} \mathbf{S}_0^H \left(\mathbf{S}_0 \mathcal{R}_{\mathbf{t}}^{-1} \mathbf{S}_0^H \right)^{-1} \quad (7)$$

from which we can observe that the channel estimate involves only the temporal correlation of the interference. Note also, that this is an *unbiased* estimate of \mathbf{H}_0 , that is, $E(\hat{\mathbf{H}}_0) = \mathbf{H}_0$ and the corresponding covariance matrix is given by [5]:

$$\mathbf{C}_{\hat{\mathbf{H}}_0} = E \left[\left(\hat{\mathbf{H}}_0 - \mathbf{H}_0 \right) \left(\hat{\mathbf{H}}_0 - \mathbf{H}_0 \right)^H \right] = \left(\tilde{\mathbf{S}}_0 \tilde{\mathbf{S}}_0^H \right)^{-1} = \left(\mathbf{S}_0 \mathbf{R}_t^{-1} \mathbf{S}_0^H \right)^{-1}. \quad (8)$$

Optimal Training for Channel Estimation. To save bandwidth, one would like to devote as few as possible symbols to training the channel estimator. Thus, given a fixed training sequence length, N_t , we want to compute the $M_T \times N_t$ training matrix \mathbf{S}_0 that minimizes $\text{tr}(\mathbf{C}_{\hat{\mathbf{H}}_0})$, subject to a constraint on the total energy consumed for training. Formally:

$$\min_{\mathbf{S}_0} \text{tr} \left[\left(\mathbf{S}_0 \mathbf{R}_t^{-1} \mathbf{S}_0^H \right)^{-1} \right] \quad (9)$$

$$\text{s.t. } \text{tr}(\mathbf{S}_0 \mathbf{S}_0^H) \leq E_T \quad (10)$$

The solution to this problem is provided in the following:

Theorem 1 *The class of training matrices optimizing the criterion (9), (10) is given by*

$$\mathbf{S}_0^{\text{opt}} = \mathbf{U} \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_{M_T} \end{bmatrix} \mathbf{G}_{M_T}^H \quad (11)$$

where \mathbf{U} can be any unitary $M_T \times M_T$ matrix and

$$\sigma_i = \sqrt{\frac{\sqrt{\lambda_i}}{\sum_{j=1}^{M_T} \sqrt{\lambda_j}}} E_T, \quad i = 1, 2, \dots, M_T, \quad (12)$$

with $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_{M_T}$ being the M_T smallest eigenvalues of \mathbf{R}_t and \mathbf{G}_{M_T} the $N_t \times M_T$ matrix of corresponding (orthonormal) eigenvectors, in that order.

Iterative Joint Channel / Data Estimation. In practice, due to the highly increased complexity that the detection of a longer sequence entails, detection is performed in pairs of vectors ($N = 2$). The ML data detection problem will then be written as in (6):

$$\begin{aligned} \min_{\mathbf{s}_0(n-1), \mathbf{s}_0(n)} & \left\| \mathbf{R}_s^{-1/2} [\mathbf{x}(n) \ \mathbf{x}(n-1)] \mathbf{R}_t^{-1/2} \right. \\ & \left. - \left(\mathbf{R}_s^{-1/2} \mathbf{H}_0 \right) [\mathbf{s}_0(n) \ \mathbf{s}_0(n-1)] \mathbf{R}_t^{-1/2} \right\|^2 \end{aligned} \quad (13)$$

or

$$\min_{\bar{\mathbf{s}}_0 \in \Omega^{2M_T \times 1}} \left\| \left(\mathcal{R}_t^{-*/2} \otimes \mathcal{R}_s^{-1/2} \right) \bar{\mathbf{x}} - \left(\mathcal{R}_t^{-*/2} \otimes \mathcal{R}_s^{-1/2} \mathbf{H}_0 \right) \bar{\mathbf{s}}_0 \right\|^2 \quad (14)$$

In the training period, \mathcal{R}_t is of size $N_t \times N_t$, with N_t being the training sequence length as above. However, in the detection phase, described by (14), the temporal correlation matrix is 2×2 . One can simply compute these two matrices separately, with the aid of sample averaging.

Then, the ML channel estimation problem can be written as

$$\min_{\mathbf{H}_0} \left\| \mathcal{R}_s^{-1/2} \mathbf{X} \left(\mathbf{I}_{N/2} \otimes \mathcal{R}_t^{-1/2} \right) - \mathcal{R}_s^{-1/2} \mathbf{H}_0 \mathbf{S}_0 \left(\mathbf{I}_{N/2} \otimes \mathcal{R}_t^{-1/2} \right) \right\|^2 \quad (15)$$

with a 2×2 matrix $\mathcal{R}_t^{-1/2}$. Solving for \mathbf{H}_0 , we obtain an estimate as in (7) where \mathcal{R}_t^{-1} should be replaced by $\mathbf{I}_{N/2} \otimes \mathcal{R}_t^{-1}$. The proposed iterative procedure will be henceforth referred to as *ILS-SD-R*.

2.4 Simulation Results

The effectiveness of ILS-SD, as compared to SD detection based on the trained channel estimate can be seen in Fig. 1(a) for uncorrelated Rayleigh block fading channels. The performance of ILSP is also shown. As expected, both trained SD

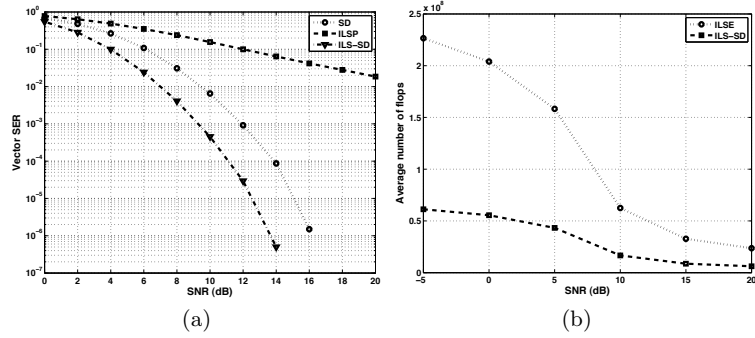


Fig. 1. (a) Comparing ILS-SD with SD based on training only and ILSP. Uncorrelated Rayleigh flat fading channels with $M_T = M_R = 4$ and QPSK input. (b) Computational requirements for convergence of the ILS-SD algorithm as compared to ILSE.

and ILS-SD perform much better than ILSP. ILS-SD is seen to greatly outperform trained-only SD. Moreover, we have seen [10] that ILS-SD converges faster than ILSP on the average, with the difference being more noticeable for low and medium values of the (per antenna) SNR. The computational savings in ILS-SD as compared to exhaustive enumeration (ILSE) are significant, as can be seen in the example of Fig. 1(b).

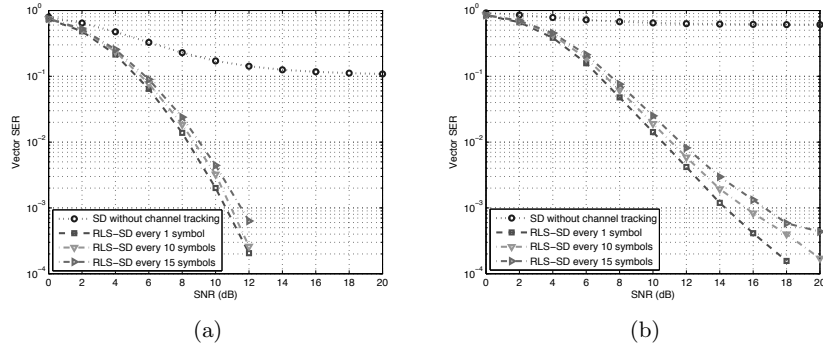


Fig. 2. Performance of RLS-SD. Results of simulation over uncorrelated Rayleigh flat fading channels using ten training symbols and for two mobile speeds: (a) 10 m/s and (b) 30 m/s. Probability of outage: 5%.

Some representative results from simulating RLS-SD for uncorrelated Rayleigh channels are shown in Fig. 2. For the sake of comparison, the results of employing SD with no channel tracking are also included. The loss in performance for RLS-SD when the update is done every $T > 1$ symbol periods is seen to be insignificant for sufficiently small values of T . It is worthwhile to notice that, similar results with that of Figs. 1(a), 2 have been also obtained [10] for correlated Rayleigh and Ricean channels.

The results for the multiuser case are demonstrated in Figs. 3, 4. The training-based SD scheme is also evaluated in Fig. 3. The results of ignoring the interference colors (temporal [9] and spatial/temporal) are also included. A considerable reduction in SER is seen to be achievable by employing optimal training in estimating the channel, especially for moderate to high SINR values compared with orthogonal (DFT) training. One can conclude that exploiting CCI color can result in significant performance gains. Moreover, as expected [14], it appears that the interference *spatial* correlation accounts for most of this benefit.

Using ILS-SD-R iterations results in Fig. 4. One can see that the gain from employing optimal training in ILS-SD-R initialization is now canceled by the iterative improvement procedure, especially for sufficiently long training sequences. Note, however, that, as seen in Fig. 4(b), optimal training can still result in faster convergence, at least in the moderate to low SINR regime.

3 Adaptive BLAST-type Decision-Feedback Equalization Schemes for Wideband DS-CDMA Systems

In the second part of the thesis we study adaptive equalization algorithms for DS-CDMA systems. We propose two new adaptive equalizers [11] for time-varying and frequency selective channels.

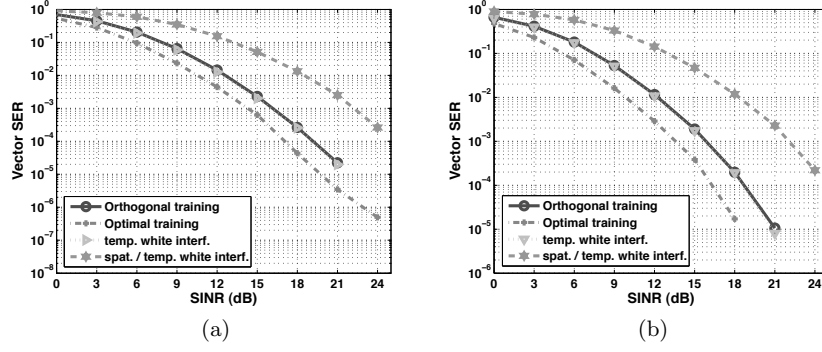


Fig. 3. Performance of SD scheme based on training only, with orthogonal and optimal training. Interference-limited environment (INR=20 dB) with (weakly) correlated Rayleigh channels. The effects of not taking the temporal and the spatial/temporal interference colors into account are also shown. (a) 8 training symbols (b) 12 training symbols.

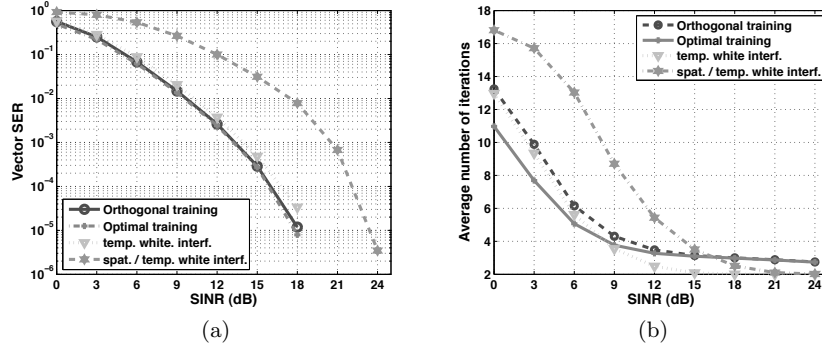


Fig. 4. Performance of the ILS-SD-R in the setup of Fig. 3. (a, b) 8 training symbols.

3.1 System Model

We consider the uplink of a symbol-synchronous DS-CDMA system with a spreading factor of P chips per symbol, K single-antenna users, and a single-antenna receiver. The users transmit independently symbol sequences which are spread through a short P -periodic spreading code $\mathbf{c}_i = [c_i(0) \ c_i(1) \ \cdots \ c_i(P-1)]^T$. The spreading codes are assumed to be known at the receiver. The transmission is through time-varying frequency selective channels, of length L with $L \leq P$.

Sampling at chip rate and collecting P successive measurements of the received signal x in a $P \times 1$ vector, a multiple-input multiple-output (MIMO) formulation with K inputs and P outputs [19] results for the DS-CDMA system. Similarly, collecting $P + L - 1$ successive samples of x instead of P , a new MIMO formulation with K inputs and $P + L - 1$ outputs results.

3.2 An Adaptive BLAST-type Equalization Scheme

An adaptive MIMO DFE detection scheme with variable detection order was proposed in [1] for flat time-varying channels. It was shown that the proposed technique performs similarly to V-BLAST algorithm with RLS channel tracking but at a reduced computational complexity. At each time instant, the receiver carries out the equalization in K serially connected stages. The users are detected in an ordered manner, applying a DFE at each stage. The stronger users are detected first, allowing easier detection for the weaker users [3]. The equalizer filters and the order of detection are updated at each stage by minimizing a set of LS cost functions for all candidate users. The user which attains the minimum cost is selected to be the next detected user.

An algorithm which exhibits the same BER performance as the above method but with reduced computationally complexity and favourable numerical behaviour was proposed in [12], based on the updating of the inverse Cholesky factor of the input autocorrelation matrix. An extension of this method to include frequency selective channels was developed in [7], where expanded input and weight vectors are used in order to eliminate both MAI and ISI.

Viewing a frequency selective DS-CDMA system as a MIMO system, as referenced in Section 3.1, the efficient square root LS algorithm of [7] can be straightforwardly applied for multiuser data detection. The resulting scheme will henceforth be referred to as the square root multiuser detection (SR-MUD) algorithm.

3.3 A RAKE-based Adaptive SIC Scheme

It is important to notice that, in the course of the SR-MUD algorithm, knowledge of users' code sequences is not required. An improved version of the SR-MUD algorithm can be developed, through incorporating knowledge of the code sequences by exploiting the RAKE receiver concept.

The structure of the new adaptive scheme is similar to SR-MUD. In this scheme the second MIMO formulation of the DS-CDMA system, presented in

Section 3.1, is used. However, a modified input signal is applied to the feedforward filter, utilizing the RAKE receiver idea. Specifically, the received signal is multiplied by a convolution matrix containing one-chip shifts of \mathbf{c}_i and the output of this product consists the input of the feedforward filter. Hence, we take advantage of the known code sequences to lessen the effect of the other users. However, due to the non-orthogonality of the distorted code sequences, the feed-back filter is necessary so as to eliminate the effect of the residual ISI and MAI. Moreover, based on the fact that the equalizer input vector can be expressed in an order-recursive manner an efficient order-update relation for the equalizer weights and the LS error energies can be obtained. Finally, through time- and order-update equations, we efficiently calculate the weights of the equalizers and determine the detection ordering. The proposed algorithm will henceforth be referred to as RAKE-RLS.

3.4 Simulation Results

The performance of the proposed adaptive schemes is compared with the RAKE receiver, the ASIC algorithm, and the linear receiver adapted via the exponentially weighted conventional RLS. The single user bound (SUB), is also shown as a benchmark. In our experiments, we simulate a near-far scenario, where the received amplitude of each user is determined such that $10 \log_{10}(A_i/A_{i+1})^2 = N$ dB, and the amplitude of the first user is set to 1.

The BER performance versus E_b/N_0 (dB) is depicted in Fig. 5 for $K = 7$, $L = 6$, $N = 2$ dB, and for different values of spreading factor P . The superiority of the proposed schemes is evident in the higher E_b/N_0 regime. Specifically, for small values of P ($P = 8$) and at high E_b/N_0 , SR-MUD outperforms RAKE-RLS, while for large values of P ($P = 128$) RAKE-RLS attains the best performance. The superiority of the proposed schemes have also been demonstrated [11] for different values of the channel length and the number of users.

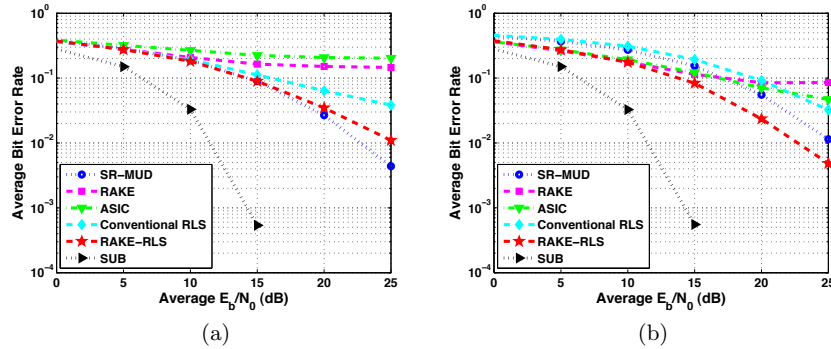


Fig. 5. BER vs. E_b/N_0 (dB) for $K = 7$, $L = 6$, $N = 2$ dB and spreading factor (a) $P = 8$ (b) $P = 128$.

4 Conclusions

Semi-blind schemes for ML joint channel estimation and data detection in MIMO flat fading channels were examined in this thesis. Both block-iterative and recursive algorithms were considered, to address block and continuous fading scenarios, respectively. The multiuser MIMO scenario, resulting in temporally/spatially colored CCI, was also addressed and the gains from exploiting CCI were assessed. The presented simulation results demonstrated the practical applicability of the investigated schemes in realistic environments. Moreover, two new adaptive equalization algorithms for time-varying and frequency selective channels in a DS-CDMA system were derived, based on the BLAST idea. The first algorithm results from a straightforward application of the idea of [1] to a MIMO-formulated DS-CDMA system, while the second one arises by incorporating the RAKE receiver concept to the first scheme. Both the equalizer filters and the optimum detection ordering are efficiently updated through time- and order-update equations. Improved BER performance is offered compared to existing adaptive DS-CDMA equalizers, in a near-far mobile environment and over a wide range of spreading factors, channel lengths and numbers of users.

References

1. J. Choi, H. Yu, and Y. H. Lee, "Adaptive MIMO decision feedback equalization for receivers with time-varying channels," *IEEE Trans. Signal Processing*, vol. 53, no. 11, pp. 4295–4303, Nov. 2005.
2. M. O. Damen, H. E. Gamal, and G. Caire, "On maximum-likelihood detection and the search for the closest lattice point," *IEEE Trans. Information Theory*, vol. 49, no. 10, pp. 2389–2402, Oct. 2003.
3. A. Duel-Hallen, J. Holtzman, and Z. Zvonar, "Multiuser detection for CDMA systems," *IEEE Personal Communications*, vol. 2, no. 2, pp. 46–58, April 1995.
4. J. Jaldén and B. Ottersten, "On the complexity of sphere decoding in digital communications," *IEEE Trans. Signal Processing*, vol. 53, no. 4, pp. 1474–1484, Apr. 2005.
5. T. Kailath, A. H. Sayed, and B. Hassibi, *Linear Estimation*, Prentice-Hall, 2000.
6. T. Kaiser *et al.* (eds.), *Smart Antennas – State of the Art*, Hindawi, 2005.
7. V. Kekatos, A. A. Rontogiannis, and K. Berberidis, "Cholesky factorization-based adaptive BLAST DFE for wideband MIMO channels," *EURASIP Journal on Advances in Signal Processing*, vol. 2007, Article ID 45789.
8. T. Koike, H. Murata, and S. Yoshida, "Adaptive MLSE equalizer with per-survivor QR decomposition for trellis-coded MIMO transmission," *Wireless Personal Communications*, vol. 35, pp. 213–225, Oct. 2005.
9. L. Li, H. Li, and Y.-D. Yao, "Channel estimation and interference suppression for space-time coded systems in frequency-selective fading channels," *Wireless Communications and Mobile Computing*, vol. 2, pp. 751–761, 2002.
10. C. Rizogiannis, E. Kofidis, C. B. Papadias, and S. Theodoridis, "Semi-blind maximum-likelihood joint channel/data estimation for correlated channels in multiuser MIMO networks," *Signal Processing*, vol. 90, no. 4, pp. 1209–1224, April 2010.

11. C. Rizogiannis, E. Kofidis, A. A. Rontogiannis, and Sergios Theodoridis, "Adaptive BLAST-type Decision-Feedback Equalizers for DS-CDMA Systems," *Proc. ICCS-2010*, Singapore, Nov. 2010.
12. A. A. Rontogiannis, V. Kekatos, and K. Berberidis, "A square-root adaptive V-BLAST algorithm for fast time-varying MIMO channels," *IEEE Signal Processing Letters*, vol. 13, no. 5, pp. 265–268, May 2006.
13. D. Shiu, G. J. Foschini, M. J. Gans, and J. M. Kahn, "Fading correlation and its effect on the capacity of multi-element antenna systems," *IEEE Trans. Communications*, vol. 48, no. 3, pp. 502–513, March 2000.
14. Y. Song and S. D. Blostein, "Channel estimation and data detection for MIMO systems under spatially and temporally colored interference," *EURASIP Journal on Applied Signal Processing*, 2004:5, pp. 685–695.
15. J. V. Steele, E. Lindskog, T. Chauvin, and S. K. Wilson, "Multiple access interference as spatially colored noise," *Proc. VTC-2002 (Fall)*, Vancouver, Canada, Sept. 2002.
16. S. Talwar, M. Viberg, and A. Paulraj, "Blind separation of synchronous co-channel digital signals using an antenna array – Part I: Algorithms," *IEEE Trans. Signal Processing*, vol. 44, no. 5, pp. 1184–1197, May 1996.
17. F. Wong and B. Park, "Training sequence optimization in MIMO systems with colored interference," *IEEE Trans. Communications*, vol. 52, no. 11, pp. 1939–1947, Nov. 2004.
18. K. Yu, M. Bengtsson, and B. Ottersten, "MIMO channel models," chap. 14 in [6].
19. W. Zha and S. D. Blostein, "Modified decorrelating decision-feedback detection of BLAST space-time system," *Proc. ICC-2002*, vol. 1, pp. 335–339, May 2002.

Technoeconomic analysis of Next Generation Networks

Theodoros Rokkas*

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
trokkas@di.uoa.gr

Abstract. In this thesis a methodological framework for the technoeconomic analysis of telecommunication networks is presented and then used in order to address problems concerning the deployment of Next Generation access Networks (NGN). Issues concerning the technoeconomic evaluation are presented such as economic indexes for the evaluation of investments, discounted cash flow analysis, real option analysis and risk analysis. The methodology is used in order to address several case studies for NGN deployments using both wired (FTTx) and wireless (FSO) architectures. In order to clarify the uncertainty real option analysis is used along with risk and sensitivity analysis.

1 Introduction

Telecom operators are skeptic in introducing fiber to the home (FTTH), due to the high investment costs associated with civil works, especially in urban and rural areas. Therefore, their current strategy is to exploit at the highest possible level their existing copper-based networks as long as possible. This strategy leads to fiber to the cabinet (FTTC) and fiber to the node (FTTN) deployments with VDSL access at last mile. A number of research and policy questions have arisen as different architectures and technologies are discussed, such as the upgrade possibilities from FTTC to FTTH. During recent years, an increasing number of research papers besides the consultancy reports, have been developed within national and international collaborative projects appeared aiming to contribute to this broadband debate. Most of these works deal with the installation first cost (IFC).

However, a complete analysis related to the FTTH deployment scenarios aiming to offer quantitative results and to analyze the associated attitude from incumbent and greenfield operators is still absent. This dissertation aims to offer these quantitative results by incorporating both “traditional” Discounted Cash Flows (DCF) analysis and Real Options Analysis (ROA).

* Dissertation Advisor: Thomas Sphicopoulos, Professor

1.1 The TONIC-ECOSYS methodology and tool

The technoeconomic methodology adopted for the evaluation of NGA deployments is based on the ECOSYS tool developed within the IST-TONIC [1] and the CELTIC-ECOSYS [2] European projects. The ECOSYS tool and its antecedent, the TONIC tool have been used for several studies among European telecom operators and universities for many years.

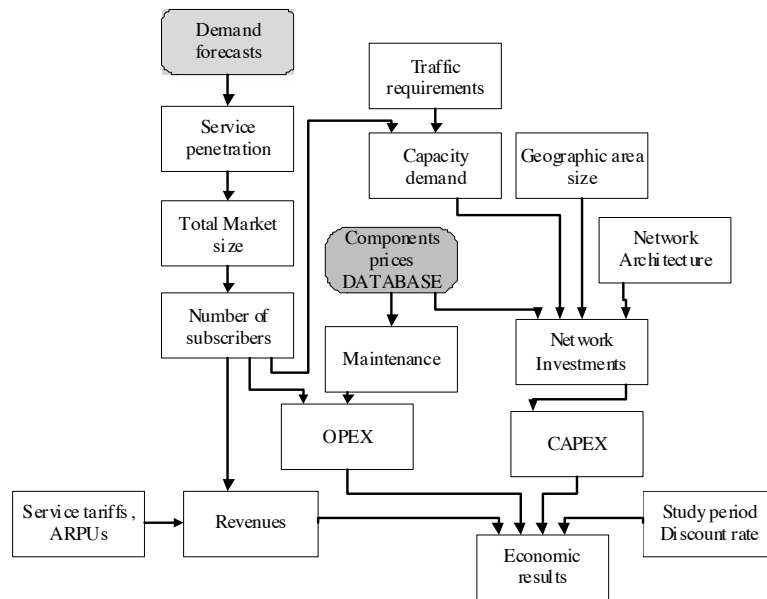


Fig. 1: Techno-economic methodology and tool

The model's operation (fig. 1) is based on its database, where the cost figures of the various network components are kept and updated from data gathered from the biggest European telecommunication companies and vendors as well as from telecom market.

The study period is best adapted to the case at hand. For a fixed network deployment case an eight to ten-year period is reasonable, considering the time a fixed network or service takes to reach market maturity and to payback the investments. The services offered are specified as well as their market penetration over the study period. In addition, the service tariffs are defined taking into account both econometric and price forecasts models and the part of the tariff that is attributed to each network under study can also be modeled. From the combination of yearly market penetration and tariff evolution, with the ECOSYS tool the revenue side of network deployments for each year given the selected service set can be calculated.

For the expenditures side, the architecture scenarios that provide the selected service set are selected from the candidate ones. This kind of modeling needs network planning expertise and is mostly outside of the framework of the ECOSYS

methodology. However many network architectures can be accounted for such as tree, mesh or ring architectures, incorporated within the tool, which includes a set of geometric models that assist in the network planning by automatically calculating lengths for cables and ducting. These geometric models are optional parts of the methodology and the ECOSYS tool can be used without them, as in the technoeconomic case of radio access technologies. Network data from other planning tools can also be used. The output of the architecture scenario definition is the so-called shopping list, which is calculated for each year of the study period and shows the volumes of all network cost elements (equipment, cables, cabinets, ducting, installation etc.) and the distribution of these network components over different network levels and layers.

In order to estimate the number of network components required throughout the study period, the necessary forecasts (both demand and price forecast) are carried out according to existing methodologies or market studies and incorporated in the technoeconomic model. The Operation Administration and Maintenance (OA&M) cost for each network element is estimated from the price of each of its constitutive parts. For example, in the case of an Ethernet switch, the model includes the switch basic equipment (switching fabric, power supply, rack and line interface cards) taking into account list price information of several vendors. The price evolution of the network components is estimated using the extended learning curve model. As far as the cost of repair parts is concerned, it is calculated by the model as a fixed percentage of the total investments in network elements while the cost of repair work is calculated based on Mean-Time-Between-Failure (MTBF) and the Mean-Time-To-Repair (MTTR).

By combining the revenues and expenditures sides, namely service revenues, investments, operating costs and general economic inputs (e.g. discount rate, tax rate), the tool calculates the results necessary for DCF analysis such as cash flows, Net Present Values (NPV), Internal Rate of Return (IRR), payback period and other economic figure of merits.

1.2 Area characteristics

The FTTH and FTTC architectures with a combination of Gigabit Ethernet and Ethernet over VDSL for the last mile are investigated, under the incumbent operators' point of view. Two area types in an average European country, namely Dense Urban (DU) and Urban (U) are under study. These areas share the same network topology but differ in several characteristics such as area dimension, population density, average cable and duct lengths, these characteristics are presented at Table 1. One common assumption is that one Central Exchange (CEX) is connected to four Local Exchanges (LEX) serves each area. Furthermore each LEX is located in the center of the service area and has a number of Cabinets connected to it. Finally all the customers are connected through their nearest cabinet (Fig 2). The fiber lengths have been calculated with the use of a geometric model. In order to model an entire European-type city, the appropriate number and pattern of dense urban and urban areas matching the city's characteristics should be assembled and added accordingly.

Table 1 Area characteristics

Area type	Dense Urban	Urban
Number of Central Exchanges (CEEx)	1	1
Number of Local Exchanges (LEEx)	4	4
Cabinets	256	256
Number of buildings	1024	2048
Subscribers per building	64	32
Total population per area	65536	65536
Total Service area (km ²)	12	32
Density (Houses/Km ²)	5641	2048

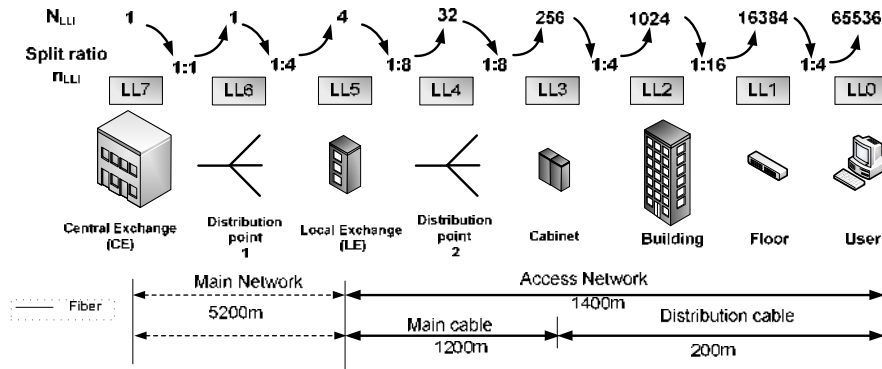


Fig. 2. Dense Urban area architecture and calculated lengths

1.3 Demand model

In the analysis, a logistic model is used to perform the demand forecasts for the selected services. This model is recommended for long-term forecasts and for new services both for fixed and mobile networks. To achieve a good fit, a four- parameter model including the saturation level is used. The model is defined by the following expression:

$$Y_t = \frac{M}{(1 + e^{a+bt})^c} \quad (1)$$

where the variables are as follows:

Y_t : is the demand actual or forecasted at time t as a population percentage

M : is the demand saturation level as a population percentage

t : is time in years and

a, b, c : are the diffusion parameters which can be estimated based on existing market data, related to broadband penetration across Europe.

2 Results and discussion

In the second part of the thesis the developed methodology was applied in selected case studies for both wireline (FTTx) and wireless (FSO) access networks.

2.1 Technoeconomic evaluation of FTTC and FTTH deployment scenarios using DCF and Real options valuation

For the case of FTTC, the incumbent operator makes a strategic decision at the first year of the project (in this model year 2009 has been used as the first year) to invest on a VDSL upgrade on the network in the dense urban areas. Part of these results was published in [3]. On the other hand, if the FTTH scenario is chosen all the copper lines are replaced with fiber ones reaching the customer premises. The decision that should be taken by an incumbent operator is whether or not should invest also to the urban areas and if yes when in the following years is the optimum time to do it. In order to answer these questions, initially the case of building these new networks simultaneously at both dense urban and urban areas was examined and then the impact of the delay of expanding the network to the urban areas as a function of time (e.g. if operator delays the expansion at urban areas for $T=1, \dots, 6$ years after the initial deployment at dense urban areas at $T=0$) was studied. The analysis was performed both with the traditional DCF analysis but also with the application of Real Options analysis.

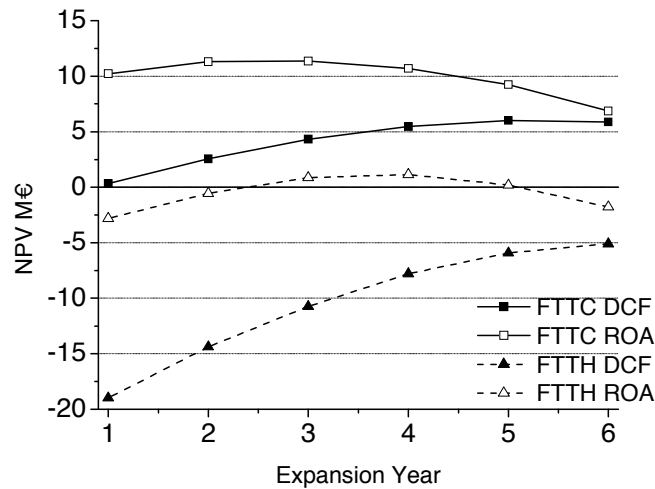


Fig. 3. ROA and DCF analysis results

At Fig.3, the NPVs for both scenarios for all the possible year of expansion with both DCF and ROA methods are presented. It can be observed that the difference in the calculations are significant for the first years and then both methods seem to converge as the years pass, which can be expected as any option value decreases as

the time reach the expiration date. However, an important finding of the ROA is that in the case of investment subsidization, the option value to expand in a later phase can significantly improve the financials of the business cases and this additional value should be taken into account. The results are financially improved and there are cases that the NPV turns positive. According to ROA the optimal strategy is to wait for four years, but the difference between the NPV values between the 2nd and the 4th year are marginal and thus, if the decision is made one or two years sooner it can be justified. This is mainly the explanation of the incumbent's current attitude, especially in Europe. Most of the incumbent operators did not invest in FTTH in urban areas and wait for either state-aid subsidizations via national funds or significant economic developments. For both cases the incumbent before making the decision must also take account the presence of the competitors in these areas, and can further benefit from an earlier investment by taking over the potential market share and have accessional economical advantages which have not been captured in the analysis made.

2.2 On the economics of Time and Wavelength Domain Multiplexed Passive Optical Networks

In the analysis made in [4], we discuss how the technoeconomic framework developed in the European project ECOSYS, can be applied to study and compare the business prospects of TDM/PON and WDM/PON architectures and provide support to the decision making process regarding access network installations.

Using the model described in the previous sections, the two PON scenarios were evaluated in the case of a greenfield telecom operator, where the PON main and access network must be build from scratch (i.e. no fiber ducts are available from any previous deployments), in the dense urban area case. The Net Present Value (NPV) was calculated for a study period of eight years starting at 2010 and is shown in Table 2. According to these results, the TDM/PON scenario appears better compared to the WDM/PON one.

Table 2 Greenfield operator results

Scenario	NPV (M€) (base case)
TDM/PON	-8,5
WDM/PON	-14,5

In order to gain a further understanding on these results and consider the uncertainties involved, a sensitivity analysis is performed for both PON architectures. Sensitivity analysis consists of the study of the impact of changes in a single parameter while all other parameters are kept constant. The parameters chosen for the sensitivity analysis were customer tariffs, service penetration, optoelectronic component prices, calculated duct length, the household density and the duct

availability. All the parameters (except the duct availability that was studied separately in a later section) were varied within an interval of $\pm 50\%$ of their initially assumed values (figures 4 and 5). In both scenarios, the most crucial parameter affecting the NPV is the customer's tariff price. In the case of TDM/PON, a 50% increase in the monthly rates (e.g. bringing it to 51€/month in the case of a silver residential subscriber) results in a positive NPV, while in the case of WDM/PON, it can improve its NPV by as much as 10M€. On the other hand, if the value is reduced by 50% (e.g. 17€/month in the case of a silver residential subscriber), the investment projects attain an NPV about two times less than the base case. It should be pointed out that, although the dependency on tariff pricing is strong, no major variations in these prices are expected due to the competition from other operators.

The next crucial parameter under consideration is the duct length that the operator has to dig in order to deploy the network and connect the customers. For the TDM/PON, if the duct length drops by half, then the project has an NPV of -0,3M€ while if it increases by 50%, then the NPV becomes twice as low (-16,6M€). The total duct length used in these calculations was estimated using an accurate geometric model and so no major variation in the initial value is expected.

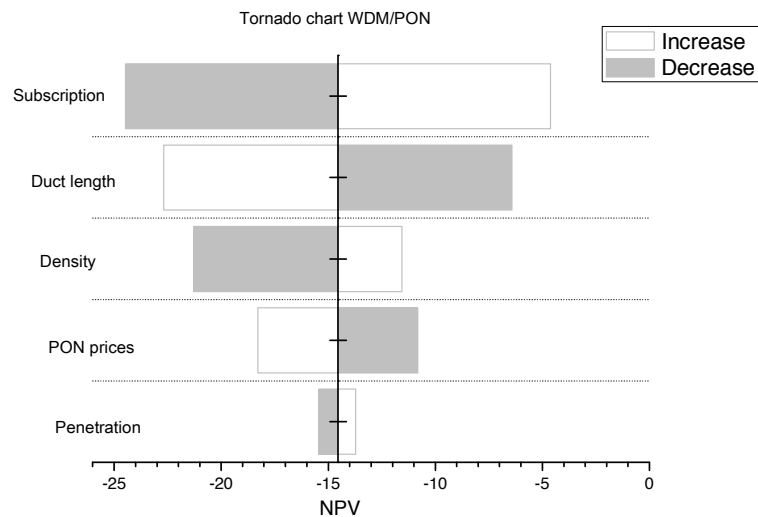


Fig. 4. Tornado chart for WDM/PON

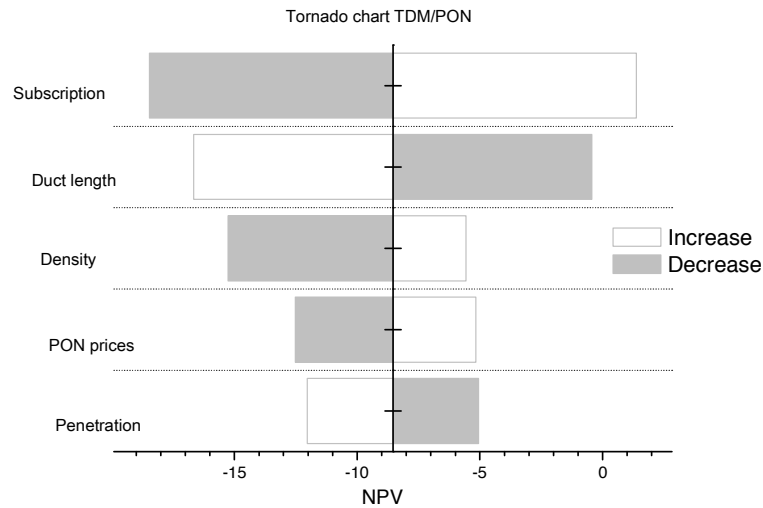


Fig. 5. Tornado chart for TDM/PON

In order to identify the most critical cost components, figures 6 and 7 illustrate the distribution of the capital investments (sum all over the study period) for WDM/PON and TDM/PON architectures respectively, for the case of a Greenfield operator. The largest part of the total costs for both examined architectures is the infrastructure expenses (digging trenches, installing cables and fibers). The cost of the optoelectronic components includes the necessary network equipment need to be installed in the OLT premises (transceivers, switches, etc) and the RN (AWGs, splitters, etc) and is more expensive for the WDM/PON than the TDM/PON. The ONT includes all equipment installed in the user premises (ONT, indoor fibers, transceivers, etc.). In the case of WDM/PON, the ONT is more expensive because of the more advanced optical components that must be installed and is almost 12% of the total investments compared to 5% in the TDM/PON case. This difference can be explained due to the higher prices of optoelectronic components used at the ONTs.

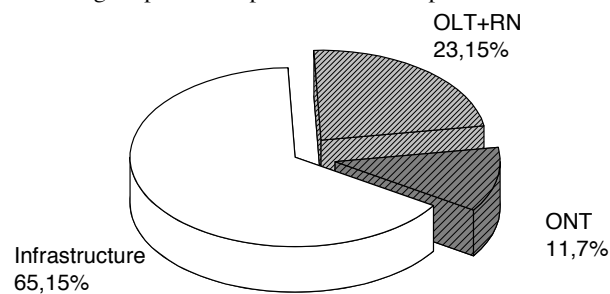


Fig. 6. Capex for TDM/PON

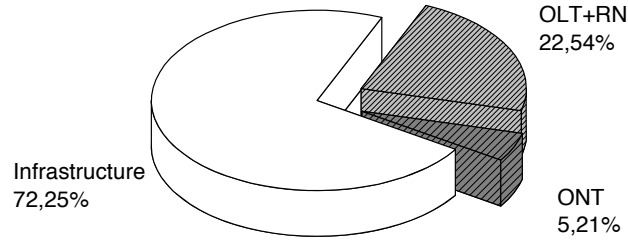


Fig. 7. Capex for WDM/PON

2.3 Business Prospects of Free Space Optics (FSO) in Dense Urban Areas

To investigate the business prospects of FSO technology for wide spread deployment in the access network, three alternative scenarios are considered [5]. In the FSO Lex-Cabinet scenario, a number of relatively Long Range (LR) Gigabit FSO links (GigFSOLR) with Automatic Gain Control (AGC) for increased reliability are deployed to connect the LEX and the Cabinet. Similar to the FTTC scenario, the users are connected through VDSL modems using the existing copper infrastructure. In the FSO1 Lex-Building scenario, GigFSOLR are used to provide wireless Gigabit connections between the buildings and the cabinet. Alternatively, in the FSO2 Lex-Building scenario, the Cabinet/customer connections have been established using less expensive, Shorter Range (SR) Gigabit FSO links (GigFSOSR) without AGC.

Table 3 compares the NPV for eight year study period for the two fiber-based and the three FSO-based scenarios assuming a dense urban area and $d_a=0$ (corresponding to limited past investments) or $d_a=70\%$ (corresponding to increased past investments). The FSO Lex-Cabinet scenario is better than FTTC and both FSO Lex-Building scenarios resulting in greater NPVs than FTTH/O in the case of no duct availability ($d_a=0$). Note also that the FSO1 Lex-Building scenario, results in a marginally positive NPV value. Based on the above remarks, it can be deduced that FSO technology can lead to more favourable business opportunities compared to its fiber counterpart, especially if no trenches for fiber ducts have been dugged up. However, the fiber-based scenarios are better than their FSO alternatives if $d_a \geq 70\%$.

Table 3. NPV Comparison of scenarios

Scenario Type	NPV (M€)
FTTC ($d_a=0$)	+15,08
FTTC ($d_a=70\%$)	+25,0
FSO Lex-Cabinet	+19,95

FTTH/O ($d_a=0\%$)	-8,86
FTTH/O ($d_a=70\%$)	+4,97
FSO1 Lex-Building	-3,67
FSO2 Lex-Buiding	+1,14

3 Conclusions

In the thesis the alternatives of FTTC/VDSL and FTTH roll-outs in dense urban and urban areas from an incumbent's point of view have been investigated. The analyzed business cases are reflecting the current stance of incumbent telecom operators regarding their decision to upgrade their infrastructure towards FTTH architecture.

Both classical DCF analysis and ROA have been used in order to evaluate the options that the incumbent has. ROA seems more suitable for capturing these effects comparing to DCF analysis. These results reveal that for FTTC the expansion can be made even for one year after the deployment at Dense Urban areas while in the FTTH case after two years.

The technoeconomic analysis carried out in the thesis revealed that infrastructure installation remains the higher cost component. If these costs can be reduced, say by using existing duct availability, or reducing the cost of civil works, then the prospects of both FTTH deployment scenarios are improved significantly. The analysis also suggests that the optimal strategy would be first to commence the installation of TDM and implement the costly infrastructure civil works and later upgrade to WDM solution when the WDM components prices will probably fall. This upgrade will use the infrastructure already deployed and will upgrade the optoelectronic equipment (ONTs, OLTs etc).

Also a technoeconomic evaluation of the business prospects of wide-scale deployment scenarios of FSO technology in the access network was carried out. Using key economic figure of merits, it was shown that in areas with limited fiber duct availability, FSO can provide an interesting, economically viable broadband alternative to FTTH and FTTC, since FSO installation does not require any civil works, which are mainly time invariable costs. As FSO equipment is mainly installed in the beginning of the deployment project, larger productions volumes and therefore better performance and price reductions due to vendor-operators agreements can be expected, which will result in further improvement of their business prospects.

References

1. IST-TONIC, 2000 - 2004, "Techno-economics of IP Optimised Networks and Services", European Union, EU IST, IST-2000-27172, Project information available at: <http://www-nrc-nokia.com/tonic/>
2. CELTIC-ECOSYS, 2002 - 2006, "techno-ECONomics of integrated communication SYStems and services", CELTIC, Project information available at: <http://www.celtic-ecosys.org/>

3. Th. Rokkas, D. Katsianis and D. Varoutas, "Techno-economic evaluation of FTTC/VDSL and FTTH roll-out scenarios: discounted cash flows and real option valuation" in *IEEE/ Journal of Optical Communications and Networking*, Vol. 2 Issue 9, pp.760-772 (2010)
4. Th. Rokkas, D. Katsianis, Th. Kamalakis and D. Varoutas, "On the economics of Time and Wavelength Domain Multiplexed Passive Optical Networks" in *IEEE/OSA Journal of Optical Communications and Networking*, Vol. 2 Issue 12, pp.1042-1051 (2010).
5. Th. Rokkas, Th. Kamalakis, D. Katsianis, D. Varoutas and Th. Sphicopoulos, "Business Prospects of Wide-Scale Deployment of Free Space Optical Technology as a Last-Mile Solution: A Techno-Economic Evaluation", *OSA Journal of Optical Networking* Vol. 6, Iss. 7, pp. 860–870 (2007)

Context Information Management for Pervasive Computing

Odysseas Sekkas*

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
`sekkas@di.uoa.gr`

Abstract. Pervasive Computing systems have to deal with the contextual information (context), which characterizes the current situation of the involved entities (e.g., users, mobile devices, environment, etc.). This dissertation studies context management issues related to the capability of a pervasive system on adapting its behavior to the involved entities context. Specifically, the interaction between the user and such system has to be less intruding as long as the latter recognizes the current user situation and adapts its functions accordingly. Such issues comprise the concept of Context Awareness. The dissertation focuses on context knowledge representation and reasoning as well as on approximate reasoning (Fuzzy Sets Theory). The management of lower level environmental information that emanates from sensors is also of great importance and is achieved with novel data fusion and decision fusion techniques that are proposed. Moreover, have been studied issues regarding collaborative context awareness and reasoning as well as bio-mimetic contextual dissemination. Adaptive algorithms are proposed so that is rendered possible the efficient dissemination of information in distributed environments. The objective is the minimization of communicating costs and the enhancement of context quality. Consequently, have been designated issues such as context discovery, context representation and inference, context fusion and collaborative context awareness

Keywords: pervasive computing, context management, context-awareness, context reasoning, epidemic algorithms

1 Introduction

In recent years we have witnessed a rapid progress in the pervasive (ubiquitous) computing paradigm. Specifically, pervasive computing is emerging as the future computing paradigm in which infrastructure and services are seamlessly available anywhere and anytime, improving human quality of life transparently to the underlying technologies. A system that is unobtrusively embedded in the environment, intuitive, constantly available and realizes the so-called ambient intelligence is defined as a pervasive system. The most profound technologies

* Dissertation Advisor: Stathes Hadjiefthymiades, Assistant Professor

are those that disappear [1]. This exciting paradigm steps from an amalgamation of information and communication technology. It is not only the consequence of convergence of advanced electronics but is also the result of contemporary research and technological advances in wireless and sensor networks, distributed systems, mobile and agent computing, autonomic and context-aware computing.

Context-awareness is one of the basic factors of the pervasive computing. It is defined as the ability of a system to use any piece of information (context) by sensing the physical environment and adapt accordingly its behavior. In order to engineer context-aware systems, it is highly important to understand and define the ingredients of context from an engineering perspective. Context defines ambient conditions and describes the situation of an entity [2]. Contextual information might change over time, describing human behaviors, application and environmental states. Context fusion is the method of deducing new and relevant information from a variety of sources in order to be used by applications and users.

In this dissertation we studied context-related issues regarding the current situation of a user (e.g. location, actions). In such case is proposed a system which exploits data streams derived from sensors, in order to accurately estimate the location of a user. The term sensors includes Wi-Fi adapters, IR receivers, RFID tag readers, etc. The core of the system is the fusion engine which is based on Dynamic Bayesian Networks (DBNs), a powerful mathematical tool for integrating heterogeneous sensor observations [3]. An extension to this system is novel context fusion engine that models, determines and reasons about the user situation. This engine which is based on Dynamic Bayesian Networks and Fuzzy Logic, deals with the reliability of sources and approximate contextual reasoning [4],[5].

For the ambient context awareness is proposed a two-level fusion scheme. To cope with heterogeneous sensors (e.g. temperature, humidity) and deliver alarms with increased accuracy and confidence, a layered fusion scheme has been adopted [6]. Different sensor feeds are processed in the two layers of the fusion scheme thus improving the reliability of the system in detecting various events. On the lower layer, the statistical behavior of sensor data is constantly assessed. On the higher layer, Dempster-Shafer (DS) theory of evidence is adopted in order to mix the indications coming from the lower layer. The proposed system has been tested for fire detection [7],[8].

Moreover, have been studied issues regarding collaborative context awareness and reasoning as well as bio-mimetic contextual dissemination. Adaptive algorithms are proposed so that is rendered possible the efficient dissemination of information in distributed environments. The objective is the minimization of communicating costs and the enhancement of context quality. Consequently, have been designated issues such as collaborative context awareness [9],[10].

The following sections describe analytically an event detection schema for context awareness. Specifically, we adopt the CUSUM test for change detection in sensor data. An improvement of this technique is also proposed. Fusing the retrieved data data from neighboring sensors we are able to mitigate problems

that lead to missed events and false alarms. Simulation results reveal the appropriateness of the mechanism in order to detect an event (particularly fire) as soon as possible and with low false alarm rate.

2 Data Fusion for Context Awareness

2.1 Event Detection

Let $\{X_i\}$ denote a sequence of random variables, i.e., a sequence of independent measurements of a sensor. We assume that X_i have density $f(x_i; \mu_0, \sigma)$ for $i = 1, \dots, \tau - 1$ and density $f(x_i; \mu_F, \sigma)$ for $i \geq \tau$, where parameter μ_0 is known and μ_F and σ are generally unknown. The time index τ signals the event (e.g. fire) in which a change in the distribution of the measurements occurs.

The parameter, μ_0 may denote the mean data value which is estimated every T_0 sec (i.e., $T_0 = 30$ min) based on sensor measurements. This time window length is in general variable and it is advisable to decrease it during daily periods that are characterized by large variations (i.e. temperature from 5:00am to 12:00am). The parameter μ_F denotes the mean value in case of event and it is considered unknown. Similarly, σ^2 denotes the unknown variance of the measurements. For example, the Sensirion SHT11 temperature sensor has an accuracy of $\pm 2.5^\circ\text{C}$ in the range from -40°C to 120°C . Adding a margin of 3°C to accommodate variations due to clouds etc., we may assume that $\sigma = 5.5^\circ$. Nevertheless, σ is a nuisance parameter and it is generally unknown. If an event occurs then the parameter τ is the time index indicating a change of densities. Sequential tests can deal with this detection of change as discussed below.

One of the most promising algorithms to sequentially detect the change is the CUSUM test [11]. For instance, if the parameter of interest is the mean value, we can monitor the partial sums

$$S_n - \min_{1 \leq k \leq n} S_k, \quad n = 1, 2, \dots$$

where $S_n = \sum_{i=1}^n X_i$ and conclude that a change from the initial μ_0 mean value to μ_F occurs at time n (as long as the previous statistic is large enough).

Gombay [12] adapted Page's CUSUM test ([11]) for change detection in the presence of nuisance parameters. Gombay proposed statistics based on the efficient score (Rao's statistics), on the maximum likelihood estimator (Wald's statistics), or on the log likelihood ratio. The efficient score vector is defined as

$$V_k(\mu, \sigma) = \sum_{i=1}^k \nabla_{\nu} \log f(X_i; \mu, \sigma), \quad \nu = (\mu, \sigma) \quad (1)$$

As it can be proved, if the density $f(\cdot)$ belongs to the exponential family, i.e., Gaussian, then once some regularity conditions hold under the null hypothesis, there exists a Wiener process $W(t)$ that approximates

$$W_k = \Gamma^{-1/2}(\mu_0, \sigma) V_k(\mu_0, \hat{\sigma}_k) \quad (2)$$

where $\hat{\sigma}_k$ is the maximum likelihood estimation of σ and $\Gamma(\mu_0, \sigma)$ is the Fisher information matrix. The test statistic W_k in (2) can be used to check if a change in densities has occurred at some time instant $\tau \leq k$. Under the alternative hypothesis, i.e., event at time τ , this statistic drifts for $k \geq \tau$ with the size of the drift proportional to the rate at which the test statistic moves in the direction of the alternative density. Moreover, in order to make decisions after n observations have been obtained, we use the following result (Darling, Erdos [13])

$$\lim_{n \rightarrow \infty} P\{a(\log(n)) \max_{1 \leq k \leq n} k^{-1/2} W_k \leq t + b(\log(n))\} = \exp(-e^{-t}) \quad (3)$$

where $a(x) = (2 \log(x))^{1/2}$ and $b(x) = 2 \log(x) + 0.5 \log(\log(x)) - 0.5 \log(\pi)$. To make use of this result we set a false alarm rate f , i.e. $f = 0.001$, where $1 - f = \exp(-e^{-t})$ and we compute the threshold

$$T(f) = (2 \log(\log(n)))^{-1/2} [-\log(-\log(1 - f)) + 2 \log(\log(n)) + 0.5 \log(\log(\log(n))) - 0.5 \log(\pi)] \quad (4)$$

Then, we conclude that the alternative hypothesis is supported by the data at the first k , if

$$k^{-1/2} W_k \geq T(f) \quad (5)$$

If no such k exists for $k \leq n$ we do not reject the null hypothesis. For $n = 900$ and the two indicative values of $f = 0.01$ and $f = 0.001$ we obtain $T(f) = 4.1$ and $T(f) = 5.3$ respectively. In what follows we assume that all measurements $X_i, i \geq 1$ are independent normal random variables. In this case the test statistic in (2) is considerably simplified. Let

$$f(x_i; \mu_0, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x_i - \mu_0)^2 / 2\sigma^2}$$

and under the alternative hypothesis (event occurrence)

$$f(x_i; \mu_F, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x_i - \mu_F)^2 / 2\sigma^2}$$

where $\mu_F > \mu_0$. The only known parameter is μ_0 that is the average value in the absence of the event.

Let $Y_i = X_i - \mu_0$ and $\mu_d = \mu_F - \mu_0$. It is clear that in the absence of event $Y_i \sim N(0, \sigma^2)$, whereas under the alternative hypothesis $Y_i \sim N(\mu_d, \sigma^2)$. In this case the test statistic is

$$k^{-1/2} W_k = k^{-1/2} \frac{\sum_{i=1}^k Y_i}{\left(\sum_{i=1}^k Y_i^2 / k\right)^{1/2}} \quad (6)$$

Under the alternative, the drift of $k^{-1/2} W_k$ after a change at time τ is

$$\text{Drift of Statistic} = k^{-1/2} \frac{(k - (\tau - 1))\mu_d}{(\sigma^2 + \frac{k - (\tau - 1)}{k} \mu_d^2)^{1/2}} \quad (7)$$

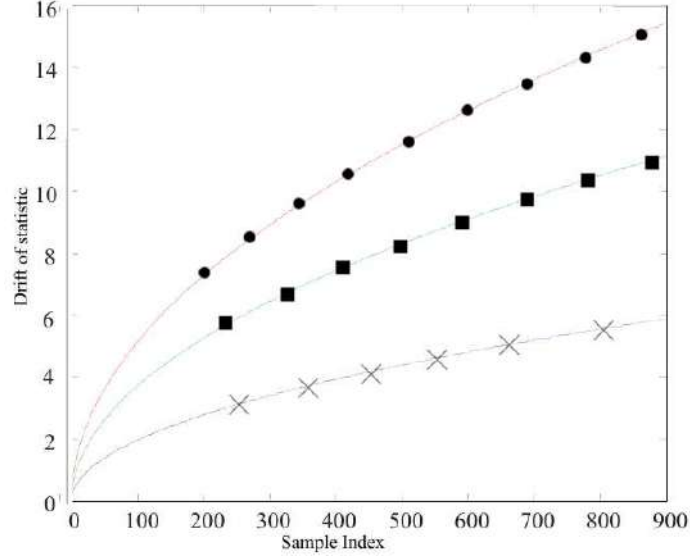


Fig. 1. Drift of statistic for various values of μ_d

Figure 1 shows the drift for $\tau = 1$, $n = 900$, $\sigma = 5$, $\mu_d = 1$ (blue-cross), $\mu_d = 2$ (green-square), and $\mu_d = 3$ (red-circle). As it is observed the greater the excess value ($\mu_d = \mu_F - \mu_0$) the largest the slope of the drift.

Table 1 shows the time instants that the test statistic crosses the thresholds of $T(f) = 4.1$ and $T(f) = 5.3$ provided that the change occurred at $\tau = 1$. As it

Table 1. Threshold cross time instants

$T(f)$	$\mu_d = 1$	$\mu_d = 2$	$\mu_d = 3$
4.1	435	120	65
5.3	740	205	100

is observed from Table 1, if the mean value excess of the alternative density is $\mu_d = 3$, it will be detected after 100 samples with a false alarm rate of 0.001.

The parameter μ_0 is estimated every T_0 based on all sensors in certain area. The period T_0 should be large enough to apply the sequential detection with as many samples as possible but small enough in order to capture the frequent changes of data.

2.2 Enhancement of the Detection Mechanism

Several issues arise when the previously described detection process is adopted. First of all, the method assumes that all sensors are calibrated. It will be a problem if one of the sensors (S_i) presents a relatively large positive offset compared to the rest of the sensors. What happens in this case is that μ_0 measured at the start of the time window T_0 is constantly smaller than the measurements of sensor S_i and, therefore, this sensor will falsely indicate an event after some time, depending on the size of the offset. A remedy to this problem is the periodic calibration of the sensors. During certain periods when exogenous parameters have no effect in the sensor measurements, offsets may be calculated and taken into account in the detection process. Thus, if a sensor presents an offset of μ_{off} compared to the average value, then the detection process of this sensor will use the value $\mu_0 + \mu_{\text{off}}$ instead of μ_0 .

A second issue is the correlation of the measurements. Criterion (6) was developed under the assumption that measurements are independent Gaussian distributed random variables. However, in real life measurements are correlated and this may cause a problem as shown in the following scenario. At the start of the interval T_0 , when the average value μ_0 is calculated various exogenous parameters may result in an underestimation of μ_0 . When such parameters do not exist anymore the average value of data will naturally increase and it will remain higher than its initial value for several samples. Depending on the relative increase and the correlation window the detection thresholds may be falsely crossed. In order to quantify and simulate the aforementioned situation we consider the following model:

We assume that sensor measurements X_i , are written as

$$X_i = \mu_0 + z_i + r_i \quad (8)$$

where z_i represents the noise due to the sensor's electronics and can be modeled as a Gaussian process of zero mean and variance σ_z^2 . The random variable r_i is the sample at time i of a process $r(t)$ which models the data readings variations due to exogenous parameters. We assume that this process is Gaussian having an autocorrelation function of the form

$$R_r(\tau) = E[r(t)r(t+\tau)] = \sigma_m^2 e^{-\alpha|\tau|} \quad (9)$$

The smaller the constant α , the greater the correlation between successive samples. The process $r(t)$ can be generated by passing white Gaussian noise $w(t)$ through a system with one pole at α , that is

$$\frac{dr(t)}{dt} = -\alpha r(t) + w(t) \quad (10)$$

where the autocorrelation function of $w(t)$ is $R_w(\tau) = 2\alpha\sigma_m^2\delta(\tau)$. From equation (10) we have

$$\frac{d}{dt} (e^{\alpha t} r(t)) = e^{\alpha t} w(t) \implies \int_t^{t+T_s} \frac{d}{dt} (e^{\alpha t} r(t)) = \int_t^{t+T_s} e^{\alpha \tau} w(\tau) d\tau$$

or else

$$e^{\alpha(t+T_s)}r(t+T_s) - e^{\alpha t}r(t) = \int_t^{t+T_s} e^{\alpha\tau}w(\tau)d\tau$$

Evaluating the previous at $t = iT_s$ we obtain

$$r_{i+1} = e^{-\alpha T_s}r_i + w_i \quad (11)$$

where

$$w_i = \int_{iT_s}^{(i+1)T_s} e^{-\alpha((i+1)T_s-\tau)}w(\tau)d\tau$$

The random variable w_i is Gaussian with zero mean and variance

$$\sigma_{w_i}^2 = E[w_i^2] = \sigma_m^2(1 - e^{-2\alpha T_s})$$

Figure 2 shows a sample function of r_i which was obtained for $\alpha = 1/120\text{sec}$, $T_s = 2\text{sec}$ and $\sigma_m = 1$. As can be seen from Figure 2, even small values of σ_m^2

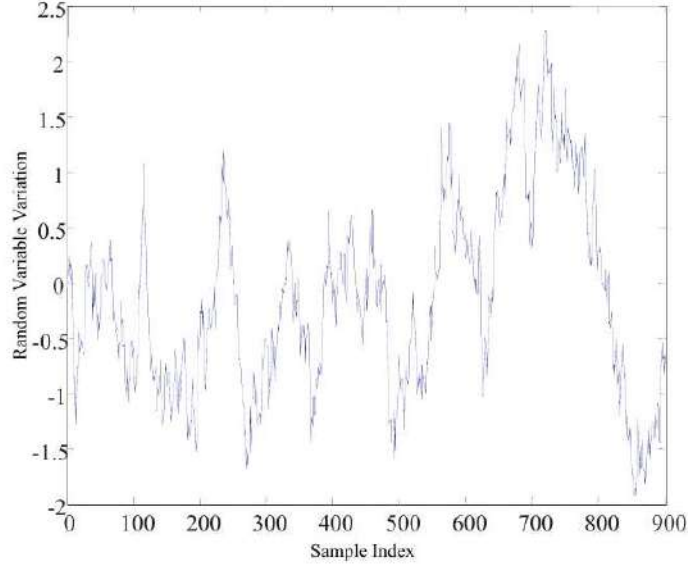


Fig. 2. Sample function of the variation random variable r_i modeling the temperature variations due to clouds, etc.

can cause large deviations from the zero mean value. The choice of the constant α indicates an average correlation window of 120 sec, that is deviations are persistent for 60 and more samples.

The problem introduced by the correlated measurements may be circumvented in one of the following ways:

1. One method is to increase the thresholds so that temporal crosses due to correlation will be avoided. A good practice is to rely on real data to set up the thresholds. However, this increase of thresholds may postpone the event detection or even cause a miss once the cross is outside the time window T_0 . Increasing the thresholds will work properly only if the assumed excess mean value μ_d is quite large.
2. A more promising solution is to rely on the cooperation of neighboring sensors to minimize correlation. As the measurements of nearby sensors undergo the same variations the term r_i can be estimated from neighboring nodes and subtracted from X_i .

Based on the second approach in order to deal with the correlated measurements, consider a sensor S_i and its neighbors S_j , $j = 1, \dots, |S_i|$, where $|S_i|$ denotes the cardinality of the neighbor set of sensor S_i . We also assume that sensor S_i senses an increase in the average value of data that is

$$X_i = \mu_0 + \mu_d + z_i + r_i \quad (12)$$

For the neighboring sensors we assume that their measurements are of the form

$$X_j = \mu_0 + z_j + r_{i-D_j} \quad (13)$$

where the noise term z_j is independent of z_i , and the term r_{i-D_j} expresses the same variation r_i that the measurements of sensor S_i undergo, delayed or advanced by D_j . Then, for sensor S_i we apply the proposed test statistic on the data.

$$Y_i = X_i - \frac{1}{|S_i|} \sum_{j=1}^{|S_i|} X_j = \mu_d + z_i + r_i - \frac{1}{|S_i|} \sum_{j=1}^{|S_i|} z_j - \frac{1}{|S_i|} \sum_{j=1}^{|S_i|} r_{i-D_j} \quad (14)$$

The term $\frac{1}{|S_i|} \sum_{j=1}^{|S_i|} z_j$ will be close to zero whereas the term $\frac{1}{|S_i|} \sum_{j=1}^{|S_i|} r_{i-D_j}$ acts as a predictor to r_i and, therefore, it almost cancels this term. Note that in applying the test statistic on data Y_i we do not have to subtract μ_0 since this term has already been cancelled.

3 Simulation results

In what follows, we will present some simulation results based on hypothetical scenarios, which emphasize the potential of the CUSUM test for the early detection of hazardous phenomena and particular fire detection. We assume two states: NOTIFY and ALERT. In the NOTIFY state the system is notifying about a possible change in the temperature mean value signaling a probable threat of a fire event. In the ALERT state, the system has to be notified on a sufficient belief for a fire event. Hence, we may use the results of Table 1 and select a false alarm rate of $f = 0.01$ to enter the NOTIFY state ($T(f) = 4.1$) and $f = 0.001$ to enter the ALERT state ($T(f) = 5.3$). When different types

of sensors (e.g. temperature and humidity) exist, we can aggregate the decisions on a fire event made from the sensed contextual data (derived from temperature and humidity sensors) in order to conclude the occurrence of a fire event.

We assume that the sampling rate is $F_s = 0.5$ Hz, that is one sample every 2 sec. We renew the estimation of the average temperature every 30 min (T_0) and therefore the time window to make a decision is $n = 30 \times 60 \times 0.5 = 900$ samples.

Scenario 1. In this scenario, a fire takes place 10 min ($\tau = 300$) after the estimation of the ambient temperature μ_0 . This fire causes an average temperature increase from $\mu_0 = 30$ to $\mu_F = 32$ Celsius degrees ($\mu_d = 2$) at the measurements of one sensor and the standard deviation is taken $\sigma = 5$. Note that μ_F is an unknown parameter that affects the slope of the drift statistic change. Figure 3a shows a sample function of the measurements, whereas Figure 3b shows the evolution of the test statistic.

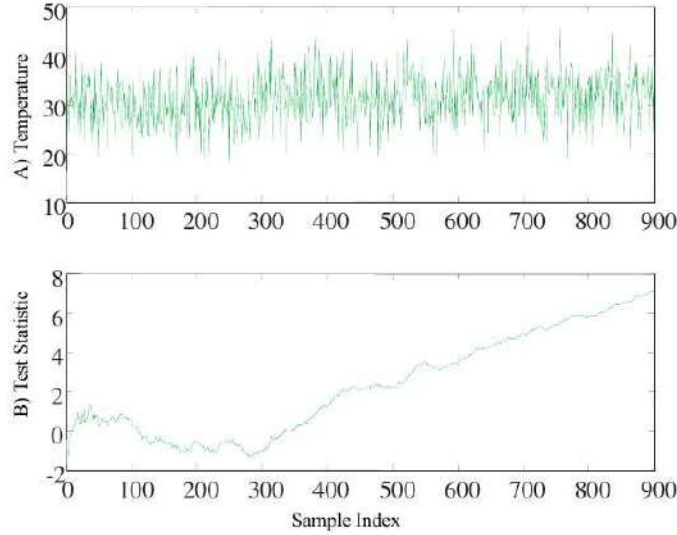


Fig. 3. (a) Sensed temperature in Celsius degree, (b) the evolution of the proposed test statistic 10 minutes after the occurrence of a fire (i.e., $\tau = 300$) with an excess value $\mu_d = 2$.

Note that by using the aforementioned false alarm rates of $f = 0.01$ (threshold $T(f) = 4.1$) and $f = 0.001$ (threshold $T(f) = 5.3$), the system will enter the NOTIFY state after approximately 300 samples (10 min) and in the emergency state after approximately 400 samples (13.5 min). Note, also, that although temperature varies greatly ($20^\circ - 40^\circ$), due to the large standard deviation, the test statistic used is insensitive to instantaneous changes.

Scenario 2. This scenario simulates the case of correlated measurements. Figure 4a shows a sample function of the process X_i in Equation (12). The mean value μ_0 was set to 30°C , $\sigma_z = 2$, $\sigma_m = 1$ and $\alpha = (120\text{sec})^{-1}$. A change of densities occurs at $\tau = 300$ with the excess mean value being $\mu_d = 2^\circ\text{C}$. This value is used only indicatively for simulation purposes. Higher values, make the test statistic, to drift faster and cross preset thresholds using fewer samples. Figure 4b shows the evolution of the test statistic. As it is observed from the figure the test statistic starts to drift after $\tau = 300$ but it might be that the thresholds $T(f)$ are crossed earlier, thus producing false alarms.

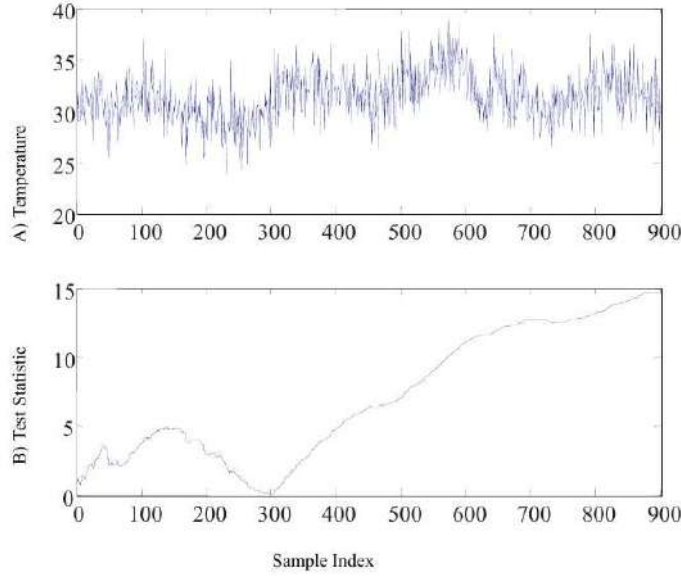


Fig. 4. (a) Sample function of the measurement process X_i , (b) the evolution of the proposed test statistic.

Figure 5 shows the simulation results for the technique proposed to mitigate correlated measurements. We assume that sensor S_i has three neighbors with corresponding delays, measured in samples, $D_1 = -3$ (time advance), $D_2 = 2$, and $D_3 = 5$. The noise z_j for each sensor is independent Gaussian with zero mean and standard deviation $\sigma_z = 2$ and $\mu_d = 2^\circ\text{C}$. A change of densities occurs at $\tau = 300$ and as illustrated in Figure 5, no false crossings of the thresholds occur prior to τ .

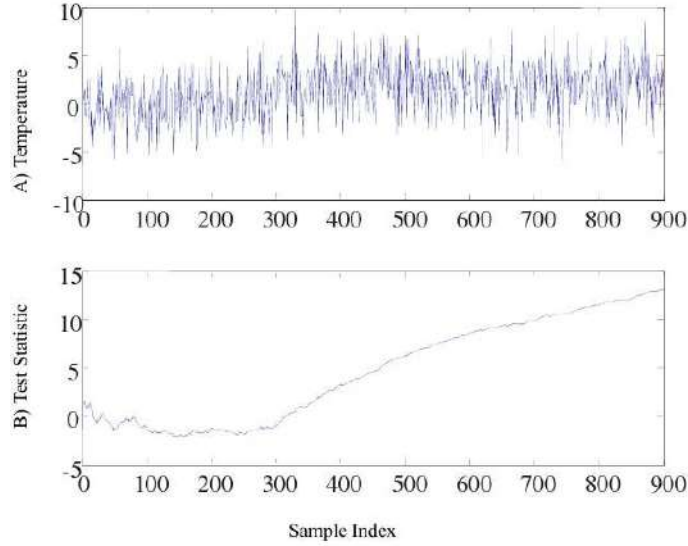


Fig. 5. (a) Sample function of the measurement process Y_i , (b) the evolution of the proposed test statistic processing the sensed contextual data from neighboring nodes.

4 Conclusions

In this dissertation we studied context management related issues for pervasive computing. A part of the dissertation deals with an event detection mechanism which is based on sensor data fusion. A cumulative sum sequential test is adopted that combines data of neighboring sensor nodes and detects changes of the underlying data distribution. The detected changes are then, compared against suitably chosen thresholds, according to a desired false alarm rate, which when crossed, the system sets it's internal notification state machine in an ALERT or a NOTIFY state. Simulation results for fire detection are also presented that verify the analysis of the proposed techniques. The synthetic trace used in our simulations contained Gaussian distributed data. CUSUM criterion is stimulated by an appropriate change of mean value of the Gaussian distribution. As a future work, we propose the enhancement of the implemented algorithms with alternative combination rules, e.g., [14], and the adoption of the Fuzzy Set theory to deal with uncertainty, imprecision and incompleteness of the underlying data.

References

1. Weiser, M.: The computer for the 21st century, Scientific American, September 1991,66-75.

2. Abowd, D.: Towards a Better Understanding of Context and Context-Awareness, Proc. International Conference of Human Factors in Computing Systems, 2000.
3. Sekkas, O., Hadjiefthymiades, S., Zervas, E.: Fusing sensor information for location estimation, in Proceedings of the 10th East-European Conference on Advances in Databases and Information Systems (ADBIS 2006), Thessaloniki, Greece, September 2006.
4. Anagnostopoulos, C., Sekkas, O., Hadjiefthymiades, S.: Context Fusion: Dealing with Sensor Reliability, in Proceedings of the 2nd International Workshop on Information Fusion and Dissemination in Wireless Sensor Networks (SensorFusion07 - MASS 2007), Piza, Italy, October 2007.
5. Sekkas, O., Anagnostopoulos, C., Hadjiefthymiades, S.: Context Fusion through Imprecise Reasoning, in Proceedings of the IEEE International Conference on Pervasive Services (ICPS 2007), pp.88-91, Istanbul, Turkey, July 2007.
6. Zervas, E., Mpimpoudis, A., Anagnostopoulos, C., Sekkas, O., Hadjiefthymiades, S.: Multisensor Data Fusion for Fire Detection, accepted for publication in Elsevier's Information Fusion, 2009.
7. Zervas, E., Sekkas, O., Hadjiefthymiades, S., Anagnostopoulos, C.: Fire Detection in the Urban Rural Interface through Fusion techniques, in Proceedings of the 1st International Workshop on Mobile Ad hoc and Sensor Systems for Global and Homeland Security (MASS-GHS 2007), Pisa, Italy, October 2007.
8. Sekkas, O., Manatakis, D., Manolakis, E., Hadjiefthymiades, S.: Sensor and Computing Infrastructure for Environmental Risks - The SCIER system, in Advanced ICTs for Disaster Management and Threat Detection: Collaborative and Distributed Frameworks, (Eds. Dr. E. Asimakopoulou and Dr. N. Bessis), IGI Global, November 2009.
9. Marias, G.F., Papapanagiotou, K., Tsetsos, V., Sekkas, O., Georgiadis, P.: Integrating a Trust Framework with a Distributed Certificate Validation Scheme for MANETs, in EURASIP Journal on Wireless Communications and Networking (EURASIP JWCN), Article ID 78259, 18 pages, 2006, Hindawi Pub. Corp.
10. Sekkas, O., Piguet, D., Anagnostopoulos, C., Kotsakos, D., Alyfantis, G., Kassapoglou-Faist, C., Hadjiefthymiades, S., Probabilistic Information Dissemination for MANETs: the IPAC Approach, in Proceedings of the 20th Tyrrhenian International Workshop on Digital Communications, Pula, Sardinia, Italy, September, 2009
11. Page, E.S.: Continuous Inspection Schemes, *Biometrika* vol. 41, 1954, pp. 100-115.
12. Gombay, E., Serban, D.: An adaptation of Pages CUSUM test for change detection, *Periodica Mathematica Hungarica*, vol. 50, 2005, pp. 135-147.
13. Darling, D.A., Erdos, P.: A limit theorem for the maximum of normalized sums of independent random variables, *Duke Math. J*, vol. 23, 1956, pp. 143-155.
14. Yager, R.R.: On the Dempster-Shafer Framework and New Combination Rules, *Information Sciences* 41:93-137.

Quality of Service Provision for IP Traffic over Wireless Local Area Networks

Dimitris Skyrianoglou¹

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
dimiski@di.uoa.gr

Abstract. This PhD dissertation deals with the provision of guaranteed quality of service (QoS) to the users of a Wireless LAN (WLAN) and the interworking between WLANs and the 3rd generation and IP networks. The work is divided into three main parts:

- i) Study and development of Wireless Adaptation Layer (WAL), a shim, transparent –both for the IP layer and the underlying WLAN- layer that supports the interworking of WLANs with IP networks and provides guaranteed QoS over WLANs by utilizing the QoS mechanisms of WAL..
- ii) Study of traffic scheduling algorithms for WLANs that are based on IEEE 802.11e protocol, the extension of legacy IEEE 802.11 protocol for supporting guaranteed QoS over 802.11 WLANs. More specifically a new traffic scheduling algorithm for 802.11e called ARROW (Adaptive Resource Reservation Over Wireless) and an extension of ARROW called ARROW (P-ARROW) was developed and evaluated.
- iii) Study of interworking between WLANs and UMTS for the provision of guaranteed QoS for the mobile users that perform a handover from one network to the other. The focus was on how the QoS mechanisms of UMTS and WLANs can interwork and combine so as to offer guaranteed QoS service to the users that perform a handover.

Keywords: Quality of Service (QoS), DiffServ, Wireless Adaptation Layer (WAL), IEEE 802.11e, Traffic Scheduling Algorithms, ARROW Scheduler, UMTS/WLAN Interworking, Seamless Handover

1 Introduction

The rapid development and the high transmission rates attained by the Wireless Local Area Networks (Wireless LANs - WLANs) have established them as one of the most attractive choices for supporting alternative access to large 3rd generation networks (3G) like UMTS or metropolitan IP networks. The installation of WLANs in places with a dense mobile user population (i.e. hot-spots like malls, airports, hospitals etc.) relieves the traffic load towards the metropolitan networks while, at the same time, achieves an improved level of quality of service for the mobile users.

¹ Dissertation Advisor: Lazaros Merakos, Professor

The work in hand deals with the provision of guaranteed quality of service (QoS) to the users of a WLAN and the interworking between WLANs and the 3rd generation and IP networks. The provision of quality of service to WLAN users at a level at least equal to that offered by the metropolitan network is deemed as especially important since the objective is to offer the mobile users a uniform level of quality of service regardless of their current location.

In this respect this work proposes the introduction of a new shim-layer called *Wireless Adaptation Layer (WAL)* that lies between the IP and the underlying wireless LAN DLC layer and aims at providing or complementing the QoS support for the underlying WLAN platform [1]-[3]. Further to this, the work delved into the QoS support mechanism of IEEE 802.11e WLAN protocol and proposed a novel traffic scheduling algorithm named *ARROW (Adaptive Resource reservation Over Wireless)* together with an extension of ARROW called *P-ARROW (Prioritized-ARROW)* [4]-[8]. Finally the work examined the interworking of WLANs with 3G networks like UMTS focusing again on the provision of QoS and proposing an architecture for supporting seamless handover for voice and video streams from one platform to the other [9]-[15].

1.1 Wireless Adaptation layer (WAL)

Several solutions are available in the literature, coping with limitations of the wireless links. Most of these solutions propose enhancements at the Transport or Application layers, while others focus on the Link Layer trying to transparently improve higher layers performance and thus avoid modifications. A number of these solutions fall into the category of Performance Enhancing Proxies (PEPs) that are defined as elements used to improve the performance of Internet Protocols on network paths where native performance suffers due to characteristics of a link or subnetwork path.

The approach proposed in this work is in line with the idea of PEPs but also tries to expand and generalize it. More specifically, it is based on the introduction of an intermediate layer called Wireless Adaptation Layer (WAL) between the IP and the Link Layer. WAL incorporates a set of functional modules, viewed as generalized PEPs, that can be dynamically combined and adapted to the special characteristics of the wireless link and the transport protocol.

WAL architecture is shown in Fig 1. A novel and key feature of the WAL is that it is an abstraction used for service provisioning at the link layer [1]-[3]. Each IP packet is classified by WAL into classes and associations. A WAL class defines the service offered to a particular set of IP packets and corresponds to a particular sequence of WAL modules that provide such a service. A WAL association identifies a stream of IP packets classified for the same WAL class and destined to or originated from a specific mobile terminal (MT). In other words, a WAL association corresponds to a particular type of service offered to a particular MT. In this way, we can differentiate the operation of WAL on a per-user basis. In addition, services for particular users can be customized to meet their specific QoS requirements and to implement a differentiated-charging policy.

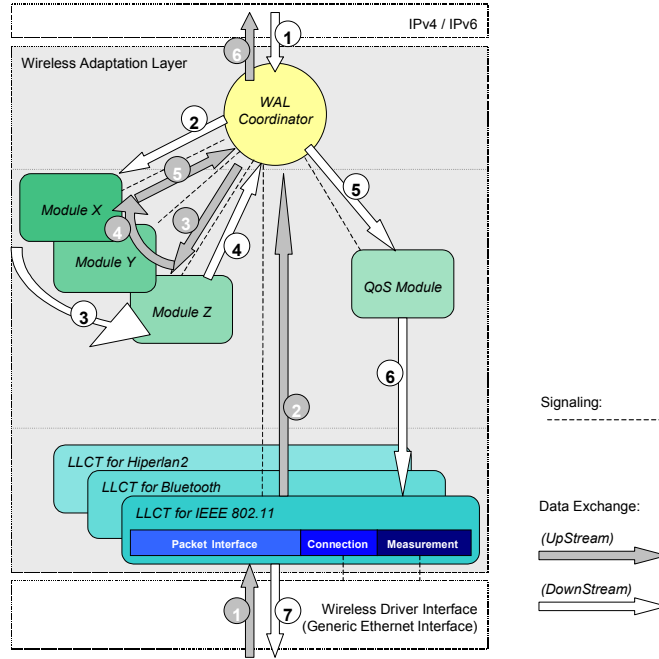


Fig. 1. WAL Architecture

The WAL Coordinator shown in Fig. 1 can be viewed as the central “intelligence” of the WAL. Both downstream (from IP layer) and upstream (to IP layer) traffic passes through the WAL Coordinator before being processed by other modules.

The QoS module (shown in Fig. 1) provides flow isolation and fairness guarantees through traffic shaping and scheduling. On the other hand, modules X/Y/Z comprise a pool of functional modules, aiming to improve performance in a number of ways. The modules that have been identified so far are: ARQ module, FEC module, Fragmentation module, IP Header Compression module, and SNOOP module.

Finally, in order to interface with a number of wireless drivers of different wireless technologies (such as IEEE 802.11, Bluetooth, HiperLAN/2, etc.), one Logical Link Control Translator (LLCT) module for each different wireless technology has been introduced. The main functions of this module manage the connection status with the wireless driver, and ensure the stream conversions toward the wireless driver.

For the classification of the IP packets to WAL classes a service differentiation is needed. Service differentiation in WAL is based on the DiffServ architecture. In this respect, the wireless access system can be viewed as a DiffServ domain with the Access Point acting as the DiffServ boundary node, interconnecting the wireless access system with the core network or other DiffServ domains.

1.2 Traffic Scheduling in IEEE 802.11e

The IEEE 802.11 standard is considered today the dominant technology for wireless local area networks (WLANs). Besides great research interest, 802.11 has enjoyed widespread market adoption in the last few years, mainly due to low-price equipment combined with high bandwidth availability. Recent improvements in the physical (PHY) layer provide transmission speeds up to hundreds of Mbps per cell, facilitating the use of broadband applications. However, one of the main weaknesses of the original 802.11, towards efficient support of multimedia traffic, is the lack of enhanced Quality of Service (QoS) provision in the Medium Access Control (MAC) layer.

In order to eliminate these weaknesses and respond to business requirements for multimedia over WLANs, IEEE is currently working on a set of QoS-oriented specification amendments, referred to as IEEE 802.11e, that enhance the existing MAC protocol and facilitate the multimedia QoS provision. In IEEE 802.11e, the QoS mechanism is controlled by the Hybrid Coordinator (HC), an entity that implements the so-called Hybrid Coordination Function (HCF). The HC is typically located in an Access Point (AP) and utilizes a combination of a contention-based scheme, referred to as Enhanced Distributed Coordination Access (EDCA), and a polling-based scheme, referred to as HCF Controlled Channel Access (HCCA), to provide QoS-enhanced access to the wireless medium. EDCA provides differentiated QoS services by introducing classification and prioritization among the different kinds of traffic, while HCCA provides parameterized QoS services to Stations (QSTAs) based on their traffic specifications and QoS requirements. To perform this operation, the HC has to incorporate a scheduling algorithm that decides on how the available radio resources are allocated to the polled QSTAs. This algorithm, usually referred to as the Traffic Scheduler, is one of the main research areas in 802.11e, as its operation can significantly affect the overall system performance [4]. Traffic Schedulers allocates resources to the QSTAs in the form of Transmission Opportunities (TXOPs). A TXOP is an interval of time when a QSTA obtains permission to transmit onto the shared wireless channel.

In the open technical literature, only a limited number of 802.11e traffic schedulers have been proposed so far and this work partially aims at filling this gap. The draft amendment of IEEE 802.11e includes an example scheduling algorithm, referred to as the Simple Scheduler, to provide a reference for future, more sophisticated solutions. The idea of this algorithm is to schedule fixed batches of TXOPs at constant time intervals. Each batch contains one fixed length TXOP per QSTA, based on the mean data rates as declared in the respective Traffic Specifications (TSPECs). With this discipline the Simple Scheduler respects the mean data rates of all TSs and performs well when the incoming traffic load does not deviate from its mean declared value (e.g., constant bit rate traffic). On the other hand, its performance deteriorates significantly when it comes to bursty traffic, as it has no means to adjust TXOP assignments to traffic changes. Identifying the weaknesses of the Simple Scheduler, SETT-EDD (Scheduling based on Estimated Transmission Times - Earliest Due Date) scheduler provides improved flexibility by allowing the HC to poll each QSTA at variable intervals, assigning variable length TXOPs. With SETT-EDD TXOP assignments are based on earliest deadlines, to reduce transmission delay and packet

losses due to expiration. SETT-EDD is a flexible and dynamic scheduler, but it lacks an efficient mechanism for calculating the exact required TXOP duration for each QSTA transmission. Each TXOP duration is estimated based on the mean data rate of each TS and the time interval between two successive transmissions.

In order to overcome the disadvantages of Simple and SETT-EDD schedulers this work proposes a new scheduling algorithm, referred to as *Adaptive Resource Reservation Over WLANs (ARROW)* [4], [6], [7], that adapts TXOP durations based on the backlogged traffic reports issued by QSTAs. The novel characteristic of ARROW is that it exploits the Queue Size (QS) field, introduced by 802.11e as part of the new QoS Data frames, not supported by legacy 802.11 systems. The QS field can be used by the QSTAs to indicate the amount of buffered traffic for their TSs, i.e., their transmission requirements. Furthermore, in order to take advantage of the periodic nature of CBR streams, a CBR-enhancement of ARROW was also developed.

Simulation results show that ARROW achieves much more efficient use of the available bandwidth, compared to Simple and SETT-EDD, leading to better channel utilization and higher throughput. The increased transmission overhead percentage of the proposed scheduler turned to be not a significant performance issue. Finally, it is important to note that ARROW does not mandate any standards changes. It could be readily deployed and implemented in practice, provided that STAs populate the QS field as defined in the 802.11e standard.

An important extension of ARROW is *P-ARROW (Prioritized ARROW)* [5]. The main enhancement of P-ARROW compared to ARROW is its ability to efficiently handle different traffic classes. The novel characteristic of P-ARROW is the introduction of Priority Factor a , and the use of traffic priorities based on delay constraints. Performance results extracted from simulation models, show that P-ARROW is very efficient in supporting the desired level of service differentiation and prioritization among different traffic classes.

1.3 UMTS/WLAN Interworking

As the Internet technologies evolve, more sophisticated and Quality of Service (QoS) demanding multimedia services are being requested by the users. The Internet Protocol (IP), together with its QoS enhancement frameworks (namely the Integrated Services - IntServ - and the Differentiated Services - DiffServ), is currently the main transport technology for supporting all these new services and in this respect the motto "Everything over IP and IP over everything" has become the trend of the day. On the other hand, both UMTS and Wireless LANs (WLANs) are already commercially available and become increasingly popular. The number of mobile users is growing rapidly and so does the demand for wireless access to the Internet services, imposing the need for a unified QoS support framework in both UMTS and WLANs.

Despite the initial impression, expressed by several network technology vendors, that UMTS and WLANs will be competing technologies it appears that they can be combined and complement each other in an effective way. The approach followed in this work is that both UMTS and WLANs can act as access systems to one common

IP core network, efficiently covering both wide areas and hot-spots. One of the main requirements of this system is a unified QoS support for IP traffic. As RSVP is considered the dominant signaling protocol of IP traffic, the discussion focuses on the adoption of RSVP messages and parameters by UMTS or WLAN QoS mechanisms [12]-[15].

Further to this, as the next-generation networks (NGN) are expected to support a wide variety of service types, especially broadband multimedia services, including video conference, streaming, and advanced telephony services, a major objective is how these services should operate seamlessly across all diverse access systems (e.g. WLAN, UMTS, fixed broadband, WiMAX, cable, etc). This seamless operation presents several challenges especially when the different access systems are loosely coupled and therefore lack the tight integration we experience in GSM/UMTS radio environments for instance. To address this issue for the case of UMTS/WLAN interworking this work proposes a specific architecture for the support of seamless voice and video handover between the two platforms [9]-[11]. The basic idea of the proposed architecture is that a new internal entity of UMTS called Seamless Handover Control Function, located at the IMS (Internet-Multimedia System) will act as an anchor-point hiding user mobility from the external IP network. Both UMTS and WLAN are also equipped with appropriate entities that take care of interworking procedures such as re-routing of traffic and authentication of the roaming users. The proposed architecture for the case of seamless voice handover is depicted in Fig. 2.

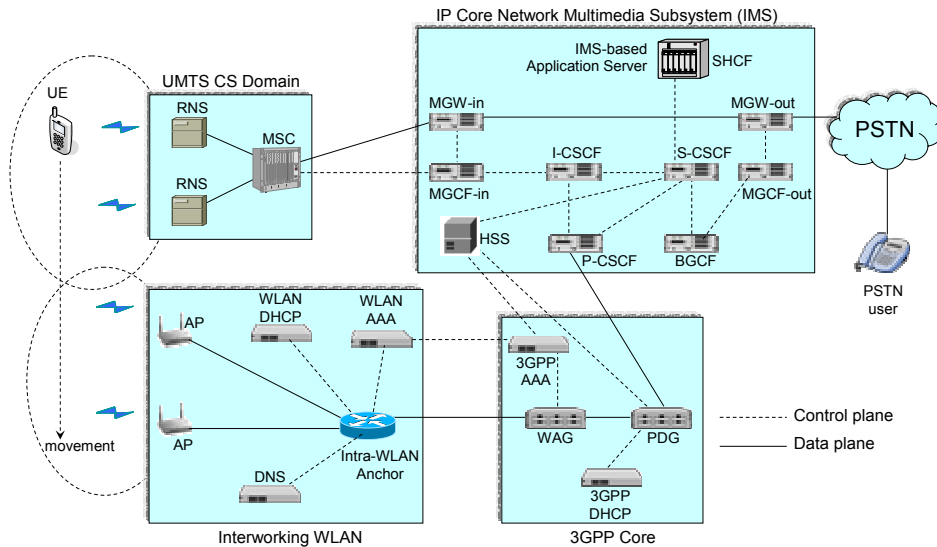


Fig. 2. IMS-based architecture for enabling seamless voice handover across UMTS and WLAN.

Simulation results indicate that WLAN can accommodate a limited number of UMTS roamers (i.e. users that perform a handover from UMTS to WLAN). This number depends on the bandwidth allocated for these users, their QoS requirements and also on the QoS support mechanism of WLAN [9]-[11].

2 The ARROW Scheduler

In IEEE 802.11e the traffic scheduler has to decide on the next TXOP assignment taking into account traffic characteristics and QoS requirements expressed through TSPEC parameters. As already mentioned, TXOP assignments are performed per QSTA, while TSPECs are defined per TS. Therefore, for each $QSTA_i$ having n_i active TSs (where i is the index of the QSTA), the traffic scheduler has to utilize some aggregate parameters, derived from the individual TSPECs, which are calculated as follows:

Minimum TXOP duration (mTD): This is the minimum TXOP duration that can be assigned to a QSTA and equals the maximum time required to transmit a packet of maximum size for any of the QSTA's TSs. Thus, mTD_i of $QSTA_i$ is calculated as:

$$mTD_i = \max \left(\frac{M_{ij}}{R_{ij}} \right), j \in [1, n_i] \quad (1)$$

where M_{ij} and R_{ij} are the maximum MSDU size and the minimum physical rate for TS_{ij} , respectively.

Maximum TXOP duration (MTD): This is the maximum TXOP duration that can be assigned to a QSTA. It should be less than or equal to the transmission time of the **Aggregate Maximum Burst Size ($AMBS$)** of a QSTA. The $AMBS$ is the sum of the maximum burst sizes (MBSs) of all TSs of a QSTA. Thus for $QSTA_i$ it holds:

$$AMBS_i = \sum_{j=1}^{n_i} MBS_{ij} \quad (2)$$

and,

$$MTD_i \leq \frac{AMBS_i}{R_i} \quad (3)$$

where R_i is the minimum physical bit rate assumed for $QSTA_i$ ($R_i = \min(R_{ij}), j \in [1, n_i]$).

Minimum Service Interval (mSI): It is the minimum time gap required between the start of two successive TXOPs assigned to a specific QSTA. It is calculated as the minimum of the $mSIs$ of all the QSTA's TSs:

$$mSI_i = \min(mSI_{ij}), j \in [1, n_i] \quad (4)$$

If not specified in the TSPEC, mSI_{ij} of TS_{ij} is set equal to the average interval between the generation of two successive MSDUs, i.e., $mSI_{ij} = L_{ij}/\rho_{ij}$.

Maximum Service Interval (MSI): It is the maximum time interval allowed between the start of two successive TXOPs assigned to a QSTA. Although no specific guidelines for calculating MSI are provided, an upper limit exists to allow an MSDU generated right after a TXOP assignment to be transmitted at the next TXOP. Accordingly:

$$MSI_i \leq D_i - MTD_i \quad (5)$$

where D_i is defined as the minimum delay bound of all TSs of $QSTA_i$ ($D_i = \min(D_{ij}), j \in [1, n_i]$). This is an upper limit that ensures that successive TXOPs will be assigned close enough to preserve delay constraints.

2.2 Operation of ARROW Scheduler

Both Simple and SETT-EDD, as briefly described above, decide on TXOP durations using some kind of estimation of the amount of data waiting to be transmitted by every QSTA. ARROW tries to overcome this drawback by adapting TXOP durations based on traffic feedback reports issued by QSTAs. The novel characteristic of ARROW is that it exploits the *Queue Size (QS)* field, introduced by 802.11e as part of the new *QoS Data* frames, not supported by legacy 802.11 systems [4], [6], [7]. The QS field can be used by the QSTAs to indicate the amount of buffered traffic for their TSs, i.e., their transmission requirements.

An example of the use of the QS field in ARROW is depicted in Fig. 3. The allocation procedure will be described in detail later in this section. For simplicity reasons, one TS per QSTA is assumed. At time $t_i(x)$, $QSTA_i$ is assigned $TXOP_i(x)$, according to requirements expressed through the QS field of the previous TXOP as well as traffic characteristics and QoS requirements declared in the respective TSPEC. Using a *QoS Data* frame, $QSTA_i$ transmits its data together with the current size of its queue in the QS field ($QS_i(x)$). At time $t_i(x+1)$ the scheduler assigns $TXOP_i(x+1)$ to $QSTA_i$, in order to accommodate the requirements of $QS_i(x)$. During the interval $[t_i(x), t_i(x+1)]$ new data is generated in $QSTA_i$, therefore $QSTA_i$ uses the *QoS Data* frame transmitted at $TXOP_i(x+1)$ to indicate the new queue size ($QS_i(x+1)$). In the same manner, at $t_i(x+2)$ the scheduler assigns $TXOP_i(x+2)$ to $QSTA_i$, accommodating the requirements of $QS_i(x+1)$ and gets the new queue size from $QSTA_i$ ($QS_i(x+2)$). As clearly shown, by utilizing the QS field, ARROW has very accurate information about the time varying properties of each TS, and is able to adapt the TXOP duration accordingly. This is considered essential, especially in the case of bursty and VBR traffic, where transmission requirements feature large time variations.

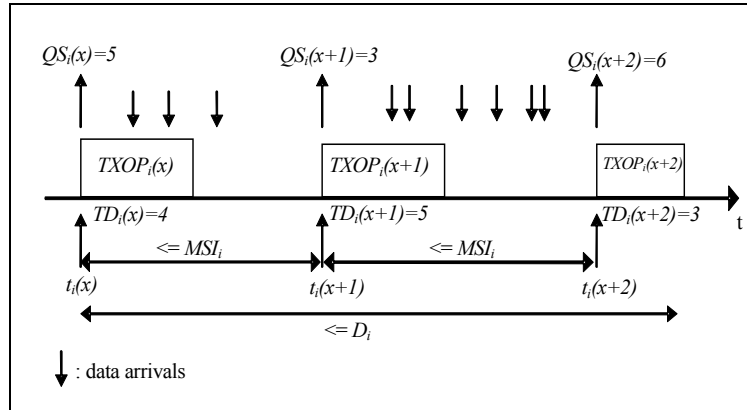


Fig. 3 TXOP assignment with ARROW

As can be observed in Fig. 3, for every $QSTA_i$, data arriving within the interval $[t_i(x), t_i(x+1)]$ can be transmitted no earlier than $TXOP_i(x+2)$ starting at $t_i(x+2)$. Therefore, in order not to exceed the delay deadline of MSDUs, assuming the worst

case that service intervals are equal to MSI_i and $TXOP_i(x+2)=MTD_i$, it should hold that:

$$\begin{aligned} D_i &\geq 2MSI_i + MTD_i \Leftrightarrow \\ \Leftrightarrow MSI_i &\leq \frac{D_i - MTD_i}{2} \end{aligned} \quad (6)$$

If the scheduler should also take into account possible retransmissions, relation (6) becomes:

$$MSI_i \leq \frac{D_i - MTD_i}{2 + m} \quad (7)$$

where m is the number of maximum retransmission attempts.

ARROW incorporates a traffic policing mechanism to ensure that the transmission requirements expressed through the Qs do not violate traffic characteristics expressed through the TSPECs. For that purpose, a *TXOP timer* is used, that implements the operation of a leaky bucket of time units. The TXOP timer value T_i for a $QSTA_i$ having n_i active TSs, increases with rate $r(T_i)$:

$$r(T_i) = \sum_{j=1}^{n_i} \left(\left(\frac{L_{ij}}{R_{ij}} + O \right) / \frac{L_{ij}}{\rho_{ij}} \right) \quad (8)$$

where O is the overhead due to PHY and MAC headers measured in seconds.

Equation (8) means that during the time interval needed for the generation of an MSDU of Nominal Size at mean data rate, the TXOP Timer should be increased by the time required for the transmission of this MSDU. The maximum TXOP Timer value $\max(T_i)$ equals the time required for the transmission of all maximum bursts:

$$\max(T_i) = \sum_{j=1}^{n_i} \left(\frac{MBS_{ij}}{R_{ij}} + O \right) \quad (9)$$

According to the operation of ARROW described below, no TXOP longer than the current value of T_i can be assigned to $QSTA_i$ at any time. After each TXOP assignment, the value of the respective TXOP timer is reduced accordingly.

The operation of ARROW can be divided in the following steps:

1. The scheduler waits for the channel to become idle.
2. When the channel becomes idle at a given moment t , the scheduler checks for QSTAs that:

a. can be polled without violating mSI , i.e., for a $QSTA_i$ that was last polled at time t_i , it should hold that:

$$t \geq t_i + mSI_i \quad (10)$$

and,

b. their TXOP timer value T is greater than the value of their mTD , to ensure enough time for the minimum TXOP duration.

3. If no QSTAs are found, the scheduler waits until (10) becomes true at least for one QSTA and returns to step 2.

4. In different case, the scheduler polls the QSTA with the earliest deadline. The deadline for a $QSTA_i$ is the latest time that this QSTA should be polled, i.e., $t_i + MSI_i$, where t_i is the time of the last poll for $QSTA_i$.

5. Assuming $QSTA_i$ having n_i active TSs is selected for polling, the scheduler calculates TD_i , as follows:

a. For every TS_{ij} of $QSTA_i$ ($j \in [1, n_i]$), the scheduler calculates TD_{ij} , as the maximum of (i) the time required to accommodate the pending traffic, as indicated by the queue size of that TS (QS_{ij}), plus any overheads (O), and, (ii) mTD_{ij} , to ensure that the assigned TXOP will have at least the minimum duration:

$$TD_{ij} = \max \left(\frac{QS_{ij}}{R_{ij}} + O, mTD_{ij} \right) \quad (11)$$

In the special case where QS_{ij} is equal to zero, TD_{ij} is set equal to the time for the transmission of a Null-Data MSDU. In this way, $QSTA_i$ is allowed to transmit a Null-Data MSDU, in order to update the queue size information for TS_{ij} . TD_i for $QSTA_i$ is calculated as the sum of all TD_{ij} :

$$TD_i = \sum_{j=1}^{n_i} TD_{ij} \quad (12)$$

b. Finally TD_i obtained from (12) is compared with the current TXOP Timer value T_i , to ensure conformance with the negotiated traffic profile:

$$TD_i = \min(TD_i, T_i) \quad (13)$$

6. After the scheduler assigns the TXOP, it reduces the respective TXOP timer value accordingly and returns to step 1:

$$T_i = T_i - TD_i \quad (14)$$

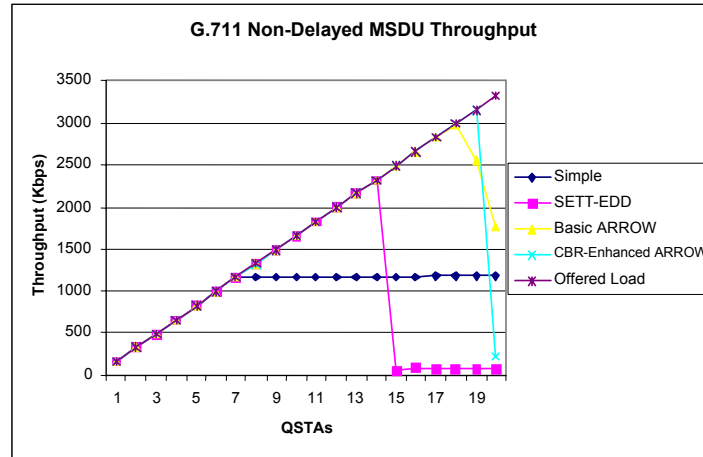
2.2 Simulation results

To measure the performance ARROW against Simple and SETT-EDD, a specialized 802.11e simulation tool developed by ATMEL Hellas was used [8]. The simulation scenarios considered an increasing number of QSTAs attached to a QAP. All QSTAs and the QAP were supporting the extended MAC layer specified in IEEE 802.11e and the PHY layer specified in IEEE 802.11g, with a transmission rate of 12Mbps. Each QSTA had two active sessions:

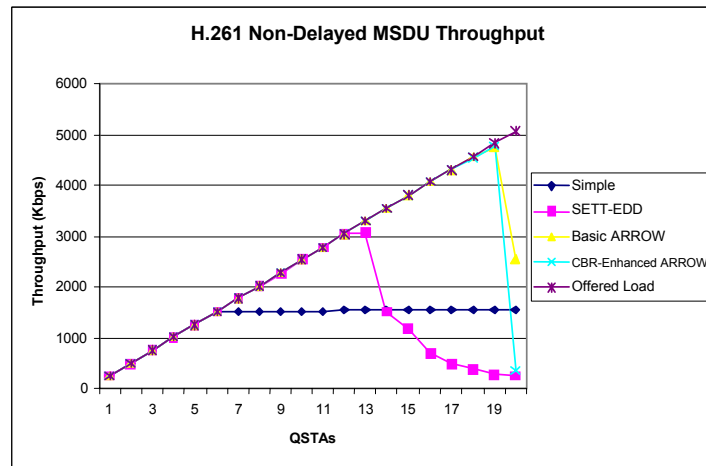
- a bi-directional G.711 voice session (CBR traffic), mapped into two TSs (one per direction), and,
- an uplink (from QSTA to QAP) H.261 video session at 256 Kbps (VBR traffic), mapped into one uplink TS.

Fig. 4 depicts throughput of non-delayed MSDUs for voice and video traffic. For voice traffic (Fig. 4a), basic ARROW accommodates up to 18 QSTAs, while SETT-EDD can manage up to 14 QSTAs and Simple up to only 7 QSTAs. Using the enhancement for CBR traffic, the number of QSTAs can be increased to 19 with CBR-enhanced ARROW, as a result of less required overhead. For video traffic (Fig. 4b), basic and CBR-enhanced ARROW outperform both SETT-EDD and Simple, accommodating up to 19 QSTAs, as opposed to 13 with SETT-EDD and 6 with Simple. The main reason for the considerably improved performance of basic ARROW is the accurate TXOP assignment it performs, by utilizing the queue size

information. This is also shown in detail using more metrics later in this section. From Fig 4a and 4b it is clear that CBR-enhanced ARROW can extend the admission capability of the system, as it can accommodate up to 19 QSTAs with voice and video TSs.



(a) G.711 Voice



(b) H.261 video

Figure 2. Throughput of Non-Delayed MSDUs

It is interesting to observe that throughput of SETT-EDD and ARROW (both basic and CBR-enhanced) reduces rapidly immediately after reaching its maximum value. The reason is that, due to the dynamic TXOP assignment performed by these algorithms, new TSs entering the system can participate equally to the channel assignment. This means that, after the overall input traffic exceeds a value that

corresponds to the maximum capability of the scheduler, none of the TSs (new or old) is serviced as required. The Simple Scheduler on the other hand, manages to provide a stable throughput regardless of the offered load, because static allocations for existing TSs are not affected by the traffic load increase. This effect highlights the need for an effective admission control scheme for SETT-EDD and ARROW that would prevent the offered load from exceeding the maximum scheduling capability.

References

1. D.Skyrianoglou, N.Passas, A. K. Salkintzis, E. Zervas "A Generic Adaptation Layer for Differentiated Services and Improved Performance in Wireless Networks", in Proc. PIMRC 2002, Lisbon, Portugal, September 2002.
2. D.Skyrianoglou, N.Passas, L.Merakos, "Improving QoS in a Multi-Hop Wireless Environment", in Proc. Med-Hoc-Net 2002, Sardegna, Italy, September 2002.
3. D.Skyrianoglou, N.Passas and S.Kampouridou, "A DiffServ-based Classification Scheme for Internet Traffic Over Wireless Links", in Proc. ICWLHN 2001, Singapore, December 2001.
4. D.Skyrianoglou, N.Passas and A.Salkintzis, "ARROW: An Efficient Traffic Scheduling Algorithm for IEEE 802.11e HCCA", IEEE Transactions on Wireless Communications, vol. 5, no.12, December 2006.
5. N. Passas, D. Skyrianoglou, and P. Mouziouras, "Prioritized Support of Different Traffic Classes in IEEE 802.11e Wireless LANs", Elsevier Computer Communications Journal, vol. 29, no. 15, September 2006.
6. D.Skyrianoglou, N.Passas and A.Salkintzis, "Traffic Scheduling for Multimedia QoS Over Wireless LANs", in Proc. ICC 2005, Seoul, Korea, May 2005.
7. D.Skyrianoglou, N.Passas and A.Salkintzis, "Traffic Scheduling in IEEE 802.11e Networks Based on Actual Requirements", in Proc. Mobile Venue 2004, Athens, Greece, May 2004.
8. A.Floros, D. Skyrianoglou, N. Passas, T. Karoubalis, "A Simulation Platform for QoS Performance Evaluation of IEEE 802.11e ", The Mediterranean Journal of Computers and Networks, vol. 2, no. 2, April 2006.
9. A. Salkintzis, N. Passas, and D. Skyrianoglou, "On the Support of Voice Call Continuity across UMTS and Wireless LANs", Wiley Wireless Communications and Mobile Computing (WCMC) Journal, vol. 8, issue 7, Sep. 2008.
10. A. Salkintzis, N. Passas, and D. Skyrianoglou, "Seamless Voice Call Continuity in 3G and WLANs", in Proc. 9th International Symposium on Wireless Personal Multimedia Communications (WPMC), San Diego, CA, September 2006.
11. A. Salkintzis, G. Dimitriadis, D. Skyrianoglou, N. Passas, F.-N. Pavlidou, "Seamless Continuity of Real-Time Video Across UMTS and WLAN Networks: Challenges and Performance Evaluation", Special issue on "Towards Seamless Interworking of WLAN and Cellular Networks", IEEE Wireless Communications Magazine, vol. 12, no. 3, Jun. 2005.
12. A Salkintzis, D. Skyrianoglou and N. Passas, "Seamless Multimedia QoS Across UMTS and WLAN Networks", in Proc. IEEE Vehicular Technology Conference (VTC) Spring 2005, Stockholm, Sweden, May 2005.
13. D. Skyrianoglou, N. Passas and A. Salkintzis, "Support of IP QoS over Wireless LANs", in Proc. VTC Spring '04, Milan, Italy, May 2004.
14. D. Skyrianoglou and N. Passas, "A Framework for Unified IP QoS Support Over UMTS and Wireless LANs", in Proc. European Wireless 2004, Barcelona, Spain, February 2004.
15. N. Passas, D. Skyrianoglou and A. Salkintzis, "Supporting UMTS QoS in WLANs" in Proc. Personal Wireless Communications (PWC) 2003, Venice, Italy, September 2003.

Transformational Government and Electronic Government Adoption Model

Teta Stamati*

Department of Informatics and Telecommunications
National and Kapodistrian University of Athens
teta@di.uoa.gr

Abstract. Despite the need expressed in the literature for shedding light upon the mechanisms that underpin the transformational process of Government, there is still research to be conducted regarding the critical factors that affect transformational government (t-Gov) and the citizens' adoption of government transformational services. To address this gap, this research reports on the findings of the use of the structured-case approach and suggests a holistic framework to investigate the success factors for t-Gov and a model that integrates the concepts that affect services adoption. The research reveals that t-Gov is not a state, but a process entailing experiential judgment. Existing acceptance theories, hence, need to be complemented by additional variables that affect citizens' adoption of transformational services.

Keywords: information systems, service science, electronic government, transformational government, adoption of electronic services, theory of reasoned action, interpretive and positivistic research.

1. Introduction

The successful delivery of public policy is increasingly dependent upon the effective use and application of new technologies and Information Systems (IS) [1]. However, significant issues are raised when policy conceptualizations travel through the many and often labyrinthine levels of public administration [1][2]. To address these issues and change the way citizens interact and communicate with each other, as well as to enhance the relationship between citizens and government, t-Gov comes to the fore [2][3][4][5][6][7][8].

The objective of the thesis is twofold. It concerns the development of the transformational process holistic model (*holistic t-Gov model*) and the construction of the adoption model (*t-Gov adoption model*) of relevant electronic services [1]. The approach of t-Gov based on three dimensions is proposed, namely organizational, social and technological. The three proposed dimensions are modeled as three submodels and their concepts and correlations are analysed. Based on qualitative research, the proposed submodels are evaluated using documentation techniques, focus groups and personal interviews and the hermeneutic conclusions are presented. The adoption model of t-Gov services is proposed and the research hypotheses are formulated. Based on qualitative research, the fundamental concepts are assessed using the aforementioned interpretive techniques and the proposed model and the conceptual hypotheses are constructed. The measurement of the model is conducted by creating the latent model. The statistic assessment is conducted in two phases, namely exploratory and confirmatory phase and it is based on methods and statistic indicators utilizing the statistic packages of SPSS and LISREL. The empirical assessment of the model is based on the positivistic approach. An assessment instrument has been developed for the data collection. The concepts of the model are operationalized to measurable variables and the measurement scales are created. The reliability and the validity of the measurement instrument are assessed. The data for the theory evaluation comes from the case study of the implementation of a unified platform for the e-services provision to citizens and businesses. The thesis reports on the use of the structured-case approach to investigate the success factors for a massive t-Gov initiative in Local Government Organisations (LGOs) to investigate the parameters that ensure the smooth use of the Local Government Application Framework (LGAF) [9]. The thesis outlines the contribution of the structured-case approach to build t-Gov theory following the interpretivist and positivist approach [6][7][8][10]. The fundament design

* Dissertation Advisor: Drakoulis Martakos, Assoc. Professor

principles of the platform are presented and the implementation objectives are documented based on the proposed models.

The structure of the article is as follows: after a brief review of t-Gov and the factors that affect its success, the proposed theoretical frameworks are presented. The research methods and context of the study are analysed. It follows the assessment of the proposed theory, the discussion of the research results as well as the presentation of the improved frameworks. The last section concludes the article.

2. Theoretical Frameworks

2.1 Holistic t-Gov Model

It has been stated that service dominates the global economy [5]. To respond to this tendency, organisations have been reconfiguring their business operations, incorporating innovative elements [5] as well as adopting new service philosophy [5]. The new philosophy to service provision has not only been adopted in the private, but in the public sector as well, where the successful delivery of public policy is increasingly dependent upon the effective use and application of new technologies and IS [1]. However, human, technical, and organisational issues seem to be arising when policy is translated to new and innovative services, which aim to transform the public administration, that is, transformational government services (t-Gov services). Regardless of the technologies selected and used, t-Gov services involve the adoption of best practices, principles and policies for their successful exploitation; on the other hand, their use creates unique opportunities, challenges and implications [5]. This calls for an analytical, interdisciplinary examination from both a theoretical and practical perspective regarding policy execution and materialization towards t-Gov [1].

The research initially addressed the aforementioned concerns and challenges by using the research-action approach according to the interpretivist approach [10][11][12][13][14][15][16]. The specific context is constituted by a large number of Government Organisations (GOs) which provide a significant number of governmental services to the citizens, visitors, enterprises based within their geographical limits and other governmental bodies. The proposed holistic approach presents the t-Gov process as a dynamic three dimensional system [1]. In this context, the concept of t-Gov is examined by three aspects namely, organizational, social, and technological. The modeling of the t-Gov process is integrated by the three corresponding sub-models t-Gov DO1, t-Gov DS2, and tGov DT3 which include meta-data and meta-information about the organizational, social and technological aspect respectively [1].

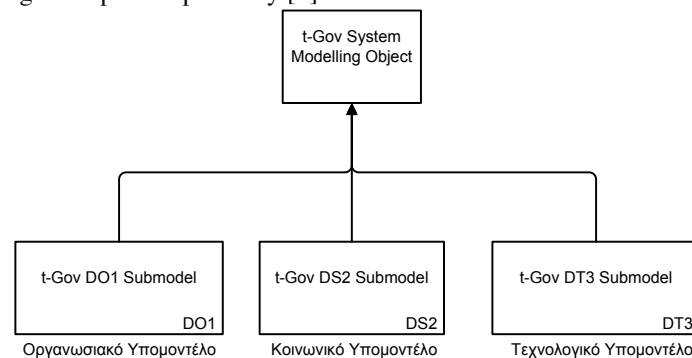


Figure 1. Holistic t-Gov Model

The following sub-models present the meta-data of each one of the three sub-systems that integrate the holistic t-Gov model.

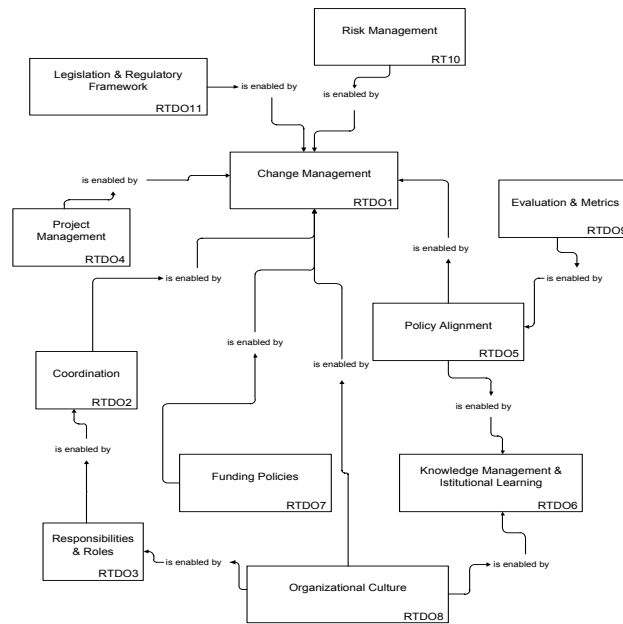


Figure 2. t-Gov DO1 Submodel

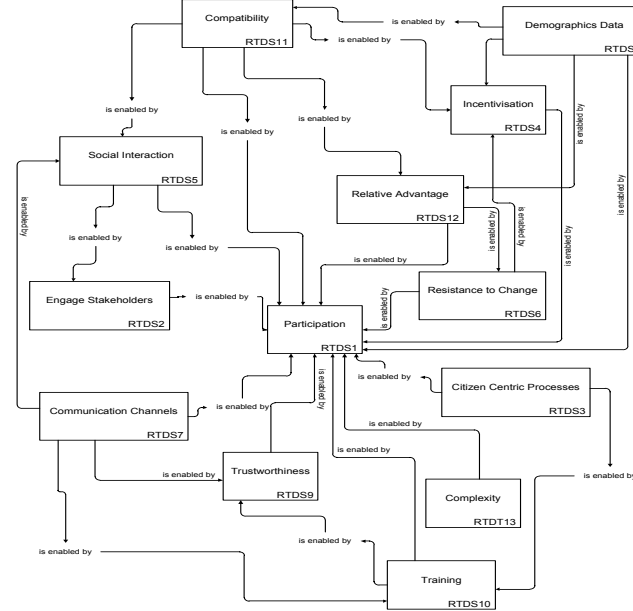


Figure 3. t-Gov DS2 Submodel

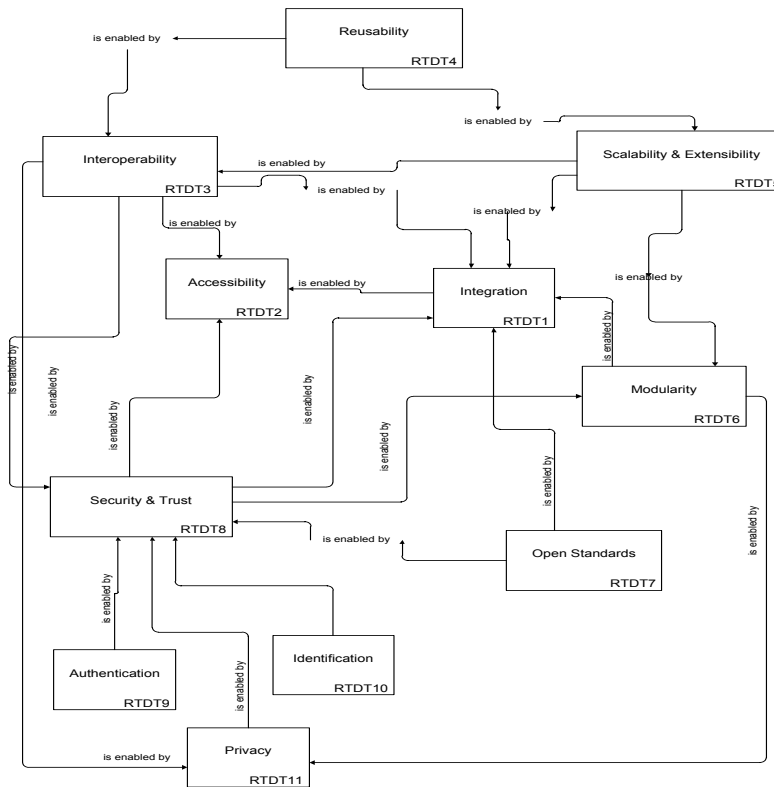


Figure 4. t-Gov DT3 Submodel

2.2 t-Gov Adoption Model

Past research on e-Gov has focused on implementation by using diffusion models. In particular, research has used Diffusion of Innovation (DOI) Theory [17]. Relevant studies [18][19][20][21] focusing on the role of administration size and professionalism on the adoption of computer technology [22]. Furthermore, literature has referred to the IS Success Model [23] and the Technology Acceptance Model (TAM) [24] as another means for discussing the particularities of the e-Gov implementation by measuring perceived usefulness (PU) and perceived ease of use (PEOU). TAM is based on the Theory of Reasoned Action (TRA). According to TRA, introduced by Martin Fishbein and Icek Ajzen [25], beliefs influence intentions, and intentions influence one's actions. TRA stresses that individual behaviour is driven by behavioural intentions, where behavioural intentions are a function of an individual's attitude toward the behaviour and subjective norms. Attitude toward the behaviour is defined as the individual's feelings about performing the behaviour. It is designated through an evaluation of one's beliefs regarding the consequences arising from a behaviour and an evaluation of the desirability of these consequences. Overall attitude can be assessed as the sum of the individual consequence multiplied by the desirability assessments, for all expected consequences of the behaviour. Subjective norm is defined as an individual's perception of whether people important to the individual think the behaviour should be performed. The contribution of the opinion of any given referent is weighted by the motivation that an individual has to comply with the wishes of that referent. Hence, overall subjective norm can be expressed as the sum of the individual perception multiplied by the motivation assessments, for all relevant referents.

TAM is one of the most well established theoretical frameworks that describe how users accept and use a technology [26]. The factors discussed by the TAM [27][28][29] have been utilised in various studies of acceptance of technology, IS, [30][31] and e-commerce [32][33][34]. Building on these TAM versions, the Unified Theory of Acceptance and Use of Technology (UTAUT) was introduced by [35], consisting of three factors namely performance expectancy, effort expectancy, and social influence and relevant studies have emerged [36][37]. However, Paul et al. [38] suggest that TAM is not conclusive and suffers from the absence of factors regarding social and human processes. Moreover, PEOU is not consistently linked to adoption [34][39][30][40]. Finally, TAM is criticised for representing subjective user assessments of a system [1][37].

Literature [37][41] suggests that since there are many similarities between e-commerce and e-Gov, TAM factors in e-commerce [32][33][34][42] could be used in the case of e-Gov [37].

However, the use of TAM has not been used extensively in the case of t-Gov, taking under consideration its nature [43][44]. Therefore, this study aims to understand the factors that affect citizens' adoption and on going usage of provided t-Gov services, and suggest a conceptual model explaining the dynamics of citizens and acceptance of the digital services. Considering the aforementioned adoption theories, the initial conceptual framework for t-Gov adoption extends the existing theoretical frameworks, as presented in the following diagram.

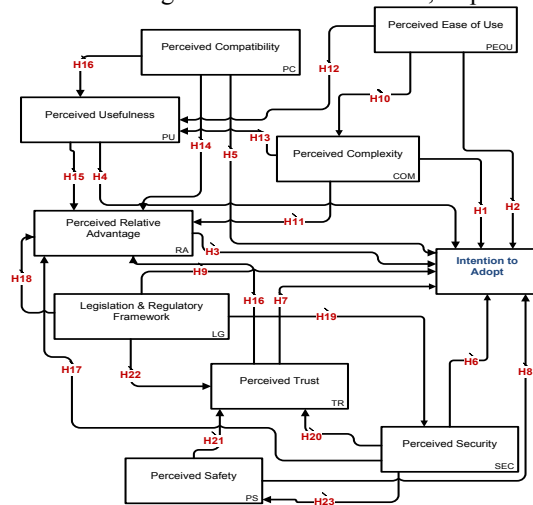


Figure 5. Conceptual Framework 1 – t-Gov Adoption Model

The Conceptual Framework 1 of the t-Gov Adoption Model consists of nine research constructs, namely:

PEOU: the degree to which a person believes that using a system would enhance his job performance

PU: the degree to which one person believes that using a particular system would be free of effort

COM: the degree to which a person believes that an innovation is being relatively difficult to use and understand. Complexity is comparable to the construct PEOU

PC: the degree to which a person believes that an innovation is seen to be compatible with existing values, beliefs, experiences and needs of adopters

RA: the degree to which a person believes that an innovation is seen as being superior to its predecessor

SEC: the degree to which a person believes that using a particular system would be secure

TR: the degree to which a person believes that using a particular system would be trustfulness

PS: the degree to which a person believes that using a particular system would be safe

LG: the degree to which a person believes that using a particular system would be according to the legislation

The resulting research hypotheses of the theoretical framework are presented below:

Hypothesis	Description
Hypothesis 1 (H1)	Higher levels of COM will reduce intention to adopt t-Gov services
Hypothesis 2 (H2)	Higher levels of PEOU of use will increase intention to adopt t-Gov services
Hypothesis 3 (H3)	Higher levels of RA will increase intention to adopt t-Gov services
Hypothesis 4 (H4)	Higher levels of PU will increase intention to adopt t-Gov services
Hypothesis 5 (H5)	Higher levels of PC will increase intention to adopt t-Gov services
Hypothesis 6 (H6)	Higher levels of SEC will increase intention to adopt t-Gov services
Hypothesis 7 (H7)	Higher levels of TR will increase intention to adopt t-Gov services
Hypothesis 8 (H8)	Higher levels of PS will increase intention to adopt t-Gov services
Hypothesis 9 (H9)	Higher levels of LG will increase intention to adopt t-Gov services
Hypothesis 10 (H10)	Higher levels of PEOU will reduce COM of t-Gov services
Hypothesis 11 (H11)	Higher levels of COM will reduce RA of t-Gov services
Hypothesis 12 (H12)	Higher levels of PEOU will be positively related to higher levels of PU of t-Gov services
Hypothesis 13 (H13)	Higher levels of COM will reduce PU of t-Gov services
Hypothesis 14 (H14)	Higher levels of PC will be positively related to higher levels of RA of t-Gov services

Hypothesis 15 (H15)	Higher levels of PU will be positively related to higher levels of RA of t-Gov services
Hypothesis 16 (H16)	Higher levels of TR will be positively related to higher levels of RA of t-Gov services
Hypothesis 17 (H17)	Higher levels of SEC will be positively related to higher levels of RA of t-Gov services
Hypothesis 18 (H18)	Higher levels of LG will be positively related to higher levels of RA of t-Gov services
Hypothesis 19 (H19)	Higher levels of LG will be positively related to higher levels of SEC of t-Gov services
Hypothesis 20 (H20)	Higher levels of SEC will be positively related to higher levels of TR of t-Gov services
Hypothesis 21 (H21)	Higher levels of PS will be positively related to higher levels of SEC of t-Gov services
Hypothesis 22 (H22)	Higher levels of LG will be positively related to higher levels of TR of t-Gov services
Hypothesis 23 (H23)	Higher levels of SEC will be positively related to higher levels of PS of t-Gov services

3. Research Method and Data Collection

Simon [45] stated that “there are always many ways to tackle a problem - some good some bad, but probably several good ways. There is no single perfect design. A research method for a given problem is not like the solution to problem in algebra. It is more like a recipe of beef stroganoff; there is no one best receipt”. One of the most important stages in this research was choosing the appropriate research philosophy, approach and method for the empirical inquiry. The study is conducted according both to the interpretivist school and the positivist school as figure 6 presents.

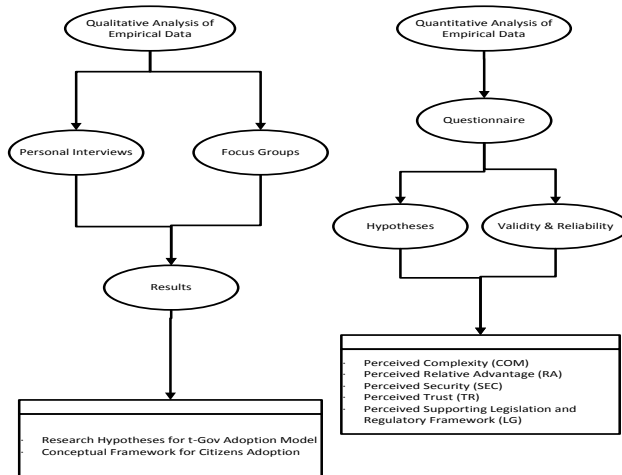


Figure 6. Research Method

The research approach throughout the study was to understand t-Gov process and to build new theory, rather than to test established theories. This was achieved by studying a number of existing theories and adoption perspectives as different theoretical lenses through which a complex phenomenon might be viewed. A methodological approach based on theory and experience about the adoption process of t-Gov was adopted and proceed to propose the holistic model for t-Gov that considers technical, human, social, legal and organizational parameters. The research that has been undertaken involves a series of case studies to a large number of GOs by means of the structured-case research method [46][47]. The application of this approach assured scientific rigor that was otherwise inadequate [7][8]. The approach provided a focused but flexible methodological approach to the field research process, through the following: outcomes integration allowed theory, knowledge and practice to emerge from the data collected; research guidance to follow and ensure accuracy; and ability to record the processes of knowledge and theory-building.

The research proposed theory, which has the form of a system of interconnected ideas that condense and organise knowledge [1][62]. The method explain, predict and provide understanding [8] determining the relationships between concepts in order to build a ‘web of meaning’ with respect to various issues of users’ adoption [46][47]. The development of conceptual frameworks namely, CF1, CF2... CFn was used to present the process of obtaining knowledge and theory building where CFn was the latest version of the theory built. The theory building process was interrelated with practice [46][47]. Applied research led to theory building, which led to further field research and theory building [1][8]. Thus, each research cycle led to changes to the existing CF. As part of the hermeneutic circle each new CF expressed the pre-understanding for the next cycle [63] following the natural human action of interpretation and world understanding. Essentially, a spiral towards understanding was enacted as current knowledge and theory foundations for yet another research cycle, which enhanced, revised and evaluated the research understanding. This was particularly appropriate for the present research, as t-Gov is an area distinguished by rapid changes, which suggests the need for theory and practice to become closely intertwined. The research methodological approach enabled theory to be developed that reflected the concerns, problems and issues facing t-Gov.

3.1 Interpretive Research

It was necessary to understand in depth the adoption process from the point of view of its meaning for the citizens as a social contract. Therefore, the research approach that was initially followed was described as being broadly interpretive. The main reasons behind this choice were the following: (a) interpretive studies attempt to understand phenomena through the meanings that people assign to them. In this research, the interpretivist approach allowed the empirically study of the factors that encourage or hinder the adoption of t-Gov services in a natural setting. These factors were influenced by many research issues and disciplines; such as organizational, managerial, technical and social; and (b) the unit of analysis in this research was the Government which is a complex social structure and is managed and controlled by different people sense-making: that is the t-Gov adoption process influences and is influenced by them. Many data collection methods, under the umbrella of the case study, have been used in the research, as the following sections explain.

Focus Groups

Focus groups referred to the form of group interview that capitalized on communication between research participants in order to generate data. Although group interviews are often used simply as a quick and convenient way to collect data from several people simultaneously, focus groups explicitly used group interaction as part of the method. The research conducted four focus group sessions designed to elicit perceptions from various stakeholder groups. A total of 87 citizens were participated in these sessions. Table 1 summarizes the types of participants at each focus group.

Focus Group	Duration	No of Participants	Role of Participants
FG1	6 h	60	employees of GOs, staff from private IT companies, citizens
FG2	3 h	15	employees from the Ministry of Interior, employees from the private IT company that was implementing LGAF
FG3	3 h	5	the team that was responsible for the implementation of LGAF
FG4	4h	7	potential adopters of LGAF

Table 1. Focus Groups

Focus groups brought together researchers who shared interests in common themes (e.g. public interest, technological and managerial issues, etc). The focus group sessions provided opportunities to explore shared beliefs and goals concerning t-Gov. The research included selected individuals at each focus group session to ensure content rich qualitative data from perspectives that would encompass the range of users and stakeholders beliefs and concerns.

In Depth Interviews: The research made use of unstructured or semi-structured set of issues and topics to guide the discussions during personal interviews. The objective of the exercises was to explore and uncover deep seated emotions, motivations and attitudes. The research attempted to deal with sensitive matters considering that the respondents were likely to give evasive or even misleading answers when directly questioned. The interviewers adhered to the following six fundamental rules [46]: avoid appearing superior or condescending and make use of only familiar words; put question indirectly and informatively; remain detached and objective; avoid questions that encourage 'yes' or 'no' answers; probe until all relevant details, emotions and attitudes are revealed; and provide an atmosphere that encourages the respondent to speak freely, yet keeping the conversation focused on the issues being researched. Totally, twelve (12) personal interviews took place.

Produced Research Model and Hypotheses: Participants interviewed during the interpretive techniques sessions completed a profile sheet which included quantitative and qualitative questions related to t-Gov success factors and t-Gov services adoption. The profile sheet asked respondents to assess, in a quantitative manner, adoption of t-Gov provided services. The participants used a Likert type scale [49] (from '1' to '5' in which '1' indicated 'strong adoption' and '5' indicated 'not adoption') to assess adoption of t-Gov services. Separate profile sheets were developed for each of the session in order to match the information needs with the various stakeholder groups. The profile provided with assessments about participant knowledge of t-gov services characteristics and attitudes, and qualitative information concerning expected user benefits, lessons learned and perceived barriers or threats to the adoption of t-Gov services. A database was created from these summaries and database management software was used to organise the data collected. A set of coding categories based on the actual data had been defined; the coding factors represented content found within the narrative summaries. Specific coding categories included categories for t-Gov services issues and information policy issues. Coding was used as a means of analyzing the data obtained from this data collection technique. Once analysed, the coding scheme provided a data reduction technique for project research. As a result of this analysis, researcher was able to query the database for specific incidents of particular factors without losing the ability to focus on the data content from a holistic perspective.

The aforementioned interpretive techniques regarding the produced conceptual framework of the adoption model, revealed the following significant results:

In the t-Gov adoption Model, the concepts of RA, PC and PU are loaded together. The constructs RA and PC have been loaded together also in other DOI research [37][50]. Moore and Benbasat conducted a thorough study using several judges and sorting rounds to develop reliable measures of diffusion of innovation constructs [18]. Although the items for RA and CT were identified separately by the judges and sorters, they all were loaded together. This may mean that, while conceptually different, they are being viewed identically by respondents, or that there is a causal relationship between the two [50]. For example, 'it is unlikely that respondents would perceive the various advantages of using t-Gov services, if their use were in fact not compatible with the respondents' experience or life style [50].

PU was also loaded with RA and PC. A similar argument to the one used to justify RA and PC loading together was used to explain that PU and RA were loading together. PU refers to the belief that a new technology will help one to accomplish a task, while RA refers to the belief that an innovation will allow one to complete a task more easily than he or she can currently.

Conceptually, these two constructs are very similar; they both refer to the use of an innovation to facilitate and ease the attainment of some goal. As RA and PU capture essentially the same concept, we decided to drop PU from further analysis. Similarly, the concepts of PEOU and COM are loaded together. The constructs of PEOU and COM were also loaded together [37]. Although the items for PEOU and COM were identified separately by the judges and sorters, they all were loaded together. This is because it is unlikely that respondents can easily perceive the provided services as ease of use if the governmental gate's use is complex. Finally, the concepts of TR and PS were loaded together. Again, although the items for TR and PS were identified separately, it is unlikely that respondents can perceive the provided services as trustworthy if Government does not provide mechanisms for safe transactions.

Considering the aforementioned issues, the model and the hypotheses tests were conducted with five independent variables – COM, RA, TR, LG and SEC as presented in the following diagram. The produced research hypotheses are H1, H3, H6, H7, H9, H11, H16, H17, H18, H19, H20 and H22 which is depicted in the produced conceptual framework as presented in the figure below.

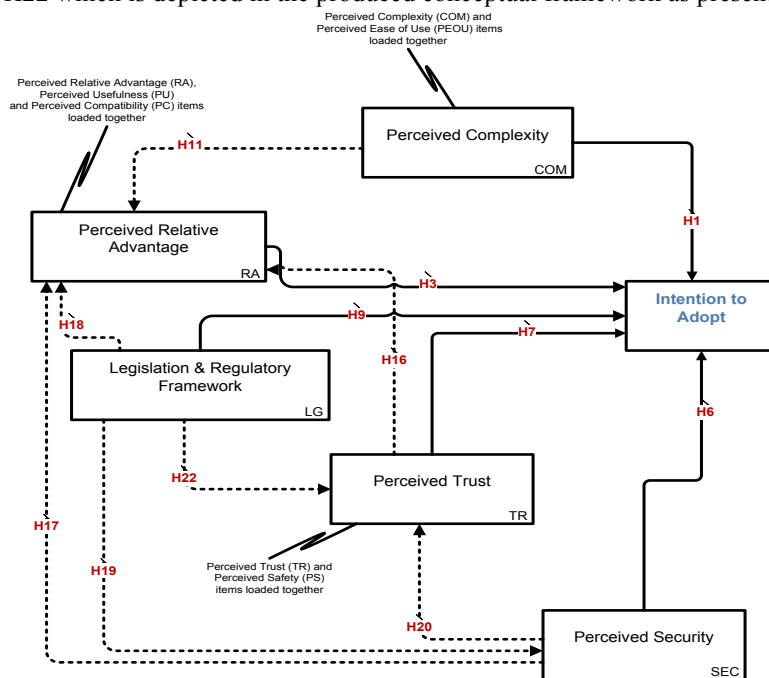


Figure 7. Conceptual Framework 2 – t-Gov Adoption Model

3.2 Positivist Research and Statistic Evaluation

An empirical study has been performed to test the model and the depicted relationships. The study was a laboratory experiment where various data used for the assessment of the conceptual frameworks. The data collection procedure and the development of measures used are shortly described in the following paragraphs.

Sample: Data were collected by administering an online questionnaire to a sample of 250 citizens. There was no information given about the actual purpose of the study. The subjects were first asked to answer to questionnaire items regarding the character of LGAF. Then they were requested to visit LGAF assuming that they were interested in a specific service provision. They were asked to look for the provision of e-services, find information about them and go through the procedure of completing in an electronic manner the transaction they were looking for. Actual running of the services was not required. After that, they were asked to indicate their responses to questionnaire items about using the provided e-service. A total of 227 responses were collected yielding an effective response rate of 90.8%. Questionnaires that were incomplete were discarded resulting in 207 usable responses.

Measures: Each of the five model constructs was operationalised using multiple items scales, based on Churchill's [51] paradigm. New scales were developed for some constructs. Items were generated based on the definition of the constructs and a review of the relevant literature so as to capture their conceptual meaning. For the rest of the constructs, items were borrowed from existing validated measures, as suggested by Straub [52], adapted with slight modifications where necessary to apply for electronic government context. All items, were measured using a 5-point Likert-type scale, ranging from '1' - strongly disagree to '5' - strongly agree, except items for

satisfaction from the overall interaction which were measured using a 5-point semantic differential scale. The measurement instrument was pretested with a sample of 12 people. The participants were presented with the list of items and a list of constructs, and were asked to assign each item to the construct that captured it best and to comment on the item's applicability to other constructs. Based on this test, the initial discriminant and convergent validity of the items was assessed to produce a refined set of measures which was used for data collection. The final version of the scales can be found in the following table.

<i>COM</i>	<i>Object</i>	<i>Source</i>
COM1	If something is complicated I do not deal with	Rogers (1995)
COM2	If something it is not easy to use I do not deal	Davis (1986)
COM3	The ease of use of new technological software plays important role in order to I utilize it	Davis (1986)
COM4	If I think that the an electronic service is simpler I will try to avoid the realisation of transaction with natural presence in the public authority	Carter & Bélanger (2004)
COM5	I easily can acquire the essential expertise in order to use the electronic provided services from the LGAF	New Object
COM6	I consider that the electronic transactions through LGAF is easy and evident process	New Object
<i>RA</i>	<i>Object</i>	<i>Source</i>
RA1	In general if something is useful in everyday routine or in work I will adopt it	Davis (1986)
RA2	If something is according to my the experiences my way of living I will adopt it	Rogers (1995)
RA3	I will do something if I believe that it offers me relative advantage	Rogers (1995)
RA4	I consider that LGAF will be useful regarding my transactions with the Public Administration	New Object
RA5	I will use LGAF if the way of running and completing a transaction is conformed with the way that I have learned to deal with the Public Authorities	New Object
<i>SEC</i>	<i>Object</i>	<i>Source</i>
SEC1	In general I have need to feel safety	Gefen (2000)
SEC2	I consider that the Government have the necessary mechanisms for citizens to feel safety	Gefen (2000)
SEC3	The security issue is significant for me	Cheung & Lee (2000)
SEC4	I feel secure to use the Internet for electronic transactions	Cheung & Lee (2000)
SEC5	I am sure that LGAF will provide secure transactions	New Object
<i>TR</i>	<i>Object</i>	<i>Source</i>
TR1	In general I trust people	Gefen (2000)
TR2	I think that in general people are reliable	Gefen (2000)
TR3	I trust Public Administration	Cheung & Lee (2000)
TR4	I trust the new technologies and the Internet	Cheung & Lee (2000)
TR5	LGAF offers trust mechanisms and thus I can entrust it for my electronic transactions with the Public Sector	New Object
TR6	I would execute an online transaction through LGAF that requires money exchange	New Object
<i>LG</i>	<i>Object</i>	<i>Source</i>
LG1	I consider important the existence of a transparent and unambiguous legislation and regulatory framework to support my transactions with the others	Carter & Bélanger (2005)
LG2	I feel confident that the legal framework will protect me in the internet	Carter & Bélanger (2005)
LG3	If the electronic transactions with the Government are imposed by the legislation framework, I will realize them through LGAF	New Object
LG4	I have been informed regarding the legislative framework that covers the electronic transactions in the internet	New Object
LG5	I believe that LGAF covers advanced legal issues in order to complete an electronic transaction with Public Administration	New Object

Table 2. List of items measures

Data Analysis: Data were analyzed with structural equation modeling techniques using SPSS and LISREL. Data analysis was based on the covariance matrix of the observed variables and was performed using maximum likelihood estimation method. The analysis was done in a two-stage procedure [53][54][55] in which the measurement model is first developed and estimated separately from the full structural equation model that models simultaneously measurement and structural relationships [54].

Measurement model: The measurement model, which is described by a set of structural equations representing the relationships between observed and latent variables, was assessed first. The model fit was adequate, with fit indices being within acceptable levels as presented in the following table [56][57][58].

Variable	KMO	Bartlett's Test of Sphericity	Eigenvalue	Correlation Matrix min value	Limits
COM	0.944	χ^2 : 1635.128 P<0,01	7.712	0.832	KMO>0,8 Eigenvalue>1 Correlation Matrix min value>0,3
RA	0.857	χ^2 : 739.336 P<0,01	3.445	0.304	
SEC	0.891	χ^2 : 818.808 P<0,01	3.840	0.538	
TR	0.883	χ^2 : 988.593 P<0,01	4.114	0.243	
LG	0.915	χ^2 : 1149.068 P<0,01	4.310	0.792	

Table 3. Variables fit indices

The goodness-of-fit indices also suggested evidence of convergent and discriminant validity as well as unidimensionality of the model constructs [57][58][59]. The measurement model was further assessed for construct reliability and validity through a Confirmatory Factor Analysis (CFA). All constructs demonstrated adequate reliability, with Cronbach's alpha values being .785 and above [58]. Reliability was assessed by computing the composite reliability of the constructs [34]. Composite reliability scores were .70 or higher, providing evidence of internal consistency [59][62][63]. Thus, all constructs were deemed reliable. Convergent validity was assessed by examining the ratio of factor loadings to their respective standard errors [56]. This ratio, represented by the t statistic value, must be greater than |2.00|, to indicate that each factor loading is greater than twice its associated standard error and should be significant for each factor loading. The model constructs satisfy both criteria for convergent validity. Each factor loading was more than double its standard error, with the lowest item t-value being 6.95. Furthermore, t-values (6.70-17.20) show that all items were loaded well on their assigned constructs. Discriminant validity was assessed with a chi-square difference test [56]. This involves setting the correlation between a pair of constructs to unity and comparing the chi-square of this model to the chi-square of the original unconstrained model. Discriminant validity between the two constructs in question was evidenced if the chi-square difference between the constrained and the unconstrained model was significant, smaller for the unconstrained model [53]. The chi-square difference was estimated for all construct pairs, providing evidence for the discriminant validity of the model constructs.

Structural model: Having tested the measurement model, the full structural model was estimated, to test the hypothesized relationships between the model constructs. The main model fit indices are as follows:

$\chi^2/d.f.$	1.40
RMSEA	0.05
CFI	0.98
NFI	0.94
NNFI	0.95
RMSR	0.078

Table 4. Model fit indices

The aforementioned indices were within acceptable levels [52][56][57]. The explanatory power of the proposed model was assessed by observing the squared multiple correlations of the endogenous constructs in the structural model estimation. The squared multiple correlations must be at least 0.10 in order for the latent construct to be judged adequate [60][61]. All model squared multiple correlations values satisfy this recommendation. The model with path coefficients for each endogenous construct is presented in the following figure.

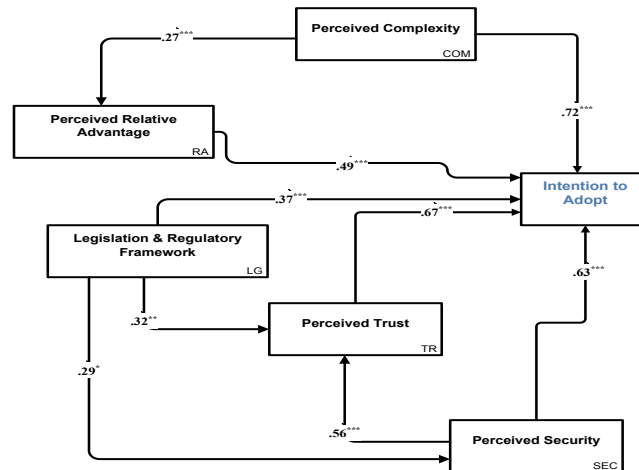


Figure 7. Conceptual Framework 3 – Proposed t-Gov Adoption Model – Structural Model Estimation

- * significant at the .05 level
 ** significant at the .01 level
 *** significant at the .001 level
 n.s. not significant at the .05 level

Hypothesis	Path Coefficient	Description
H1	.72	examined the effect of perceived complexity to intention to adopt
H3	.49	examined the effect of perceived relative advantage to intention to adopt
H6	.63	examined the effect of perceived security to intention to adopt
H7	.67	examined the effect of perceived trust to intention to adopt
H9	.37	examined the effect of the supporting legislation and regulatory framework to intention to adopt
H11	.27	examined the effect of perceived complexity to relative advantage
H19	.29	examined the effect of perceived supporting legislation and regulatory framework to security
H20	.56	examined the effect of perceived security to trust
H22	.32	examined the effect of perceived legislation and regulatory framework to trust

Table 5. Hypotheses testing results

Hypotheses 16, 17 and 18 were not supported as they were not found to be significant. Specifically, there is not significant impact between the concepts of perceived trust and relative advantage, perceived security and relative advantage, and finally between legislation and regulatory framework and relative advantage.

4. Concluding Discussion

This study presents a holistic approach for t-Gov and an integrated model for services adoption. The adoption model incorporates constructs from TAM and DOI and extends previous research. T-Gov initiatives have been identified as significant field for research in IS. Government services adoption raises important political, cultural, organisational, technological and social issues that must be considered carefully [64]. The key adoption factors of t-Gov are proposed which can be used as a tool to determine the roadmap for adoption of a t-Gov initiative. Further research should take place in order to explore the theory applicability to other environments. The next paragraphs reveal some of the research results.

Human and Social Constructs: *Compatibility* was found to have a significant relationship with use intentions in t-Gov. The theory assessment process strongly suggested that t-Gov services should be produced in a manner that is consistent with individuals' values, beliefs and experiences and provide information and work support in a manner that is consistent with what citizens are used on. Another significant concern was *Trustworthiness*. Citizens, who perceived the reliability and security of the internet to be low, presented obstacles when using t-Gov services [1]. There was a long debate between participants in the focus groups regarding the notion of initial trust to services

provided that refers to “trust in an unfamiliar trustee, a relationship in which the actors do not yet have credible, meaningful information about, or affective bonds with, each other” [1]. Regarding trustworthiness, citizens who perceived Government to be trustworthy consider the introduction of t-Gov system as a welcome initiative. Governmental-based trust was mainly associated with citizens’ perceptions of the governmental environment, such as the structures, regulations and legislation that make an individual feel safe and trustworthy [1]. Another important construct is the motivation or the perceived need for working ‘over the wire’. In demographic terms, the data analysis revealed that a percentage of 76% of the interviewees stated they intend to immediate use LGAF (early LGAF adopters) were people in young age, more educated (80% of them holding a University degree) and with relatively high incomes (40% of them had a net family income more than thirty thousand per year). This indicated that individual demographic characteristics were also influencing the adoption of provided services. The cases analysis proved that a group of individuals were more likely to keep using LGAF than others. Consequently, we examined two factors namely, the level of prior Internet usage and the citizens innovativeness. Individual innovativeness can be defined as ‘consumer acceptance’ of new ideas [1]. The findings supported that higher Internet usage led to LGAF adoption. Domain-specific innovativeness, i.e. innovation linked to certain domains was found to influence LGAF adoption. Finally, there was a group of users persuaded very quickly of the LGAF’s significant advantages compared to prior institutional systems. This proved that individual perceived relative advantage enforced the individual intention to use.

Organizational Constructs: The discussions concerned the coordination and ownership between and across GOs and departments, the political engagement regarding the delivery of technology supported services, the GO capacity including available resources (human, technical, etc.), change and risk management issues as well as the appropriate legal and legislation framework. The nature and mission of GOs were discussed and their relationship with the e-services provided. There was a clear concern regarding potential future developments and change [1]. Clear policies for GOs were seen to be critical. Key issues included sense of ownership and the required organisational transformation. A key concern was about ways to cope with organisational inertia. A particularly important area of risk was the access to governmental services and the issue of community inclusion. Furthermore, it emerged that measurement and evaluation techniques were necessary to realise the learning perspectives of t-Gov. To achieve successful transformational implementations it is necessary to establish coherent legitimacy and establish trust relationships between government and citizens. Since the legal framework regarding the provision of electronic services is ‘still in infancy’, a cohesive legal framework is required to speed the adoption of t-Gov. The research has revealed that four main sets of legislation are considered: personal data protection laws; privacy and security laws; information (provision) laws; and administrative laws.

Technical Constructs: Various technical parameters that might affect LGAF adoption and regular use were revealed. The supporting staff in GOs stressed the need for a less complex framework and more user-friendly in its user interface, and the forms and templates. The majority of interviewers and workshop participants were sceptical about the use of innovative technological tools, by aged users; the authors labelled this attribute ‘computer anxiety’. IT experts identified the need for flexible and scalable technology, privacy and security, shared services and common identity management, standards, coordination and integration between GOs operations and departments, identification and authentication. Regarding the notions of scalability and flexibility of governmental systems, the cases revealed that there is need to create flexible systems that can adapt and change on demand in accordance to the changing nature of t-Gov [1]. There was no definite agreement regarding what constitutes valid and appropriate access to information. Finally, issues of interoperability and standardisation arose, stemming from the way different GO’s departments can be managed, the technical tools needed for integration and the standardisation of certain data and services. To this extend, the notions of open standards and open source software were highlighted.

References

1. Stamati, T. and Martakos, D. Electronic Transformation of Local Government: an exploratory study, *International Journal of Electronic Government Research*, IGI Publishing, 7(1), 2010, 20-37.

2. Stamati, T., Karantjias, A., Martakos, D. Survey of citizens' perceptions in the adoption of National Governmental Portals, IGI Publishing, 2010, 213-235.
3. Stamati, T., Papadopoulos, T., Martakos, D. Transformational Services: a case study in the Greek public administration. In *Transformational Government through eGov: Socio-economic, Cultural, and Technological issues*, Emerald, 2010.
4. Stamati, T., Papadopoulos, T., Martakos, D.. Transformational Government citizens' services adoption: a conceptual framework, Volume 6846/2010, pp. 134-143, DOI: 10.1007/978-3-642-22878-0_12, Springer Berlin / Heidelberg, 2010.
5. Stamati, T., Karantjias, A. Inter-sector practices reform for e-Government integration efficacy, *Journal of Cases on Information Technology*, IGI Publishing, 13(3), 2010, 62-83.
6. Janssen, M., Shu, W. S. "Transformational government: basics and key issues". In *Proceedings of the 2nd International Conference on Theory and Practice of Electronic Governance*, ACM International Conference Proceeding Series 351 (Janowski, T. and Pardo, T. A, Eds), pp 117–122, ACM Publications, Cairo, Egypt, 2008.
7. Irani, Z., Elliman, T., Jackson, P. Electronic Transformation of government in the U.K. a research agenda. *European Journal of Information Systems*, 16, 2008, 327-335.
8. Irani, Z., Love, P., E., D., Elliman, T., Jones, S., Themistocleous, M. Evaluationg e-government: learning from the experiences of two UK local authorities, *Information Systems Journal*, 15, 2005, 61-82
9. LGAF project. Local Government Application Framework. 6th Framework Programme <http://wiki.kedke.org/wiki/>
10. Walsham, G. The emergence of interpretivism in IS research. *Information Systems Research* 6(4), 1995, 376-394.
11. Remenyi, D. *Doing research in business and management: an introduction to process and method*. London; Thousand Oaks, California, Sage Publications, 1998.
12. Denzin, N.K., Lincoln, Y.S. *Collecting and interpreting qualitative materials*. Thousand Oaks, Calif, Sage Publications, 1989.
13. Hussey, J., Hussey R. *Business research: a practical guide for undergraduate and postgraduate students*. Basingstoke: Macmillan Business, 1997.
14. Lee, A., Baskerville, R. Generalizing in information systems research. *Information Systems Research* 14(3), 2003, 221-243.
15. Myers, M.D. Qualitative research in information systems. *Management Information Systems Quarterly* 21(2), 1997, 241-242 .
16. Oates, B.J. *Researching information systems and computing*. London; Thousand Oaks, California: Sage Publications, 2006
17. Grönlund, Å. Ten years of eGovernment: the end of history and a new beginning. In Wimmer, M.A., Chappelet, J-L., Janssen, M., Scholl, H.J. (eds) (2010) *Electronic Government*. 9th IFIP WG 8.5 International Conference, EGOV 2010, Lausanne, Switzerland, August/September 2010. *Proceedings. LNCS 6228*, pp. 13-24, Springer. Best paper award
18. Rogers, E. M. *Diffusion of Innovations*. New York: Free Press, 1995.
19. Moon, M.J. The evolution of e-Government among municipalities: Rhetoric or reality? *Public Administration Review* 62(4), 2002, 424-433.
20. Moon, J., Norris, D. Does managerial orientation matter? The adoption of reinventing government and e-government at the municipal level. *Information Systems Journal* 15, 2005, 43– 60
21. Norris, D.F., Campillo, D. *Factors Affecting Innovation Adoption by City Governments: The Case of Leading Edge Information Technologies*, Maryland Institute for Policy Analysis and Research. University of Maryland, Baltimore, MD, USA, 2000
22. Angelopoulos, S., Kitsios, F., and Papadopoulos, T. Identifying Critical Success Factors in e-Government: A New service development approach. *Transforming Government: People, Process and Policy* 4(1), 2010, 95-118.
23. DeLone, W.H., McLean, E.R. Information systems success The quest for the dependent variable. *Information Systems Research* 3(1), 1992, 60–95
24. Davis, F. Perceived usefulness, perceived ease of use and user acceptance of information technology. *MIS Quarterly*, 13, 1989, 319–340.
25. Ajzen, I., Fishbein, M. Attitudes and normative beliefs as factors influencing intentions'. *Journal of Personality and Social Psychology*, 21, 1972, 1–9.
26. Davis, F. D., Bagozzi, R. P., Warshaw, P. R. User acceptance of computer technology: A comparison of two theoretical models. *Management Science*, 35(8), 982-1003 (1989).
27. Benbasat, I., Barki, H. Quo vadis TAM. *Journal of the Association for Information Systems*, 8(4), 2007, 211–218.
28. King, W. R., He, J.: A meta-analysis of the technology acceptance model. *Information & Management* 43, 740–755 (2006).
29. Schepers, J., Wetzels, M. A meta-analysis of the technology acceptance model: investigating subjective norm and moderation effects. *Information & Management*, 44, 2007, 90-103.
30. Venkatesh, V., Davis, F. A theoretical extension of the technology acceptance model: four longitudinal field studies. *Management Science* 46, 2000, 186–204.
31. Venkatesh, V., Morris, M. G. Why don't men ever stop to ask for directions? Gender, social influence, and their role in technology acceptance and usage behavior. *MIS Quarterly* 24(1), 2000, 115–139.
32. Gefen, D., Straub, D. The relative importance of perceived ease of use in IS adoption: a study of e-commerce adoption. *Journal of the Association for Information Systems*, 1, 2000, 1–28.

33. Moon, J., Kim, Y. Extending the TAM for a World-Wide-Web context. *Information & Management*, 38 (4), 2001, 217-230.
34. Gefen, D., Karahanna, E., Straub, D.: Trust and TAM in online shopping: an integrated model. *MIS Quarterly*, 27, 2003, 51–90.
35. Venkatesh, V., Morris, M. G., Davis, G. B., Davis, F. D. User acceptance of information technology: Toward a unified view. *MIS Quarterly*, 27(3), 2003, 425–478.
36. Shajari, M., Ismail, Z. A comprehensive adoption model of e-Government services in developing countries, *Advanced Management Science*, IEEE International Conference 2, 2010, 548-553.
37. Carter, L., Bélanger, F. The utilization of e-government services: citizen trust, innovation and acceptance factors. *Information Systems Journal* 15, 2005, 5-25.
38. Paul, L., John, I., Pierre, C. Why do people use information technology? A critical review of the technology acceptance model. *Association for Information Systems*, 40 (3), 20003, 191.
39. Ma, Q., Liu, L. The Technology Acceptance Model: A Meta-Analysis of Empirical Findings. *Journal of Organizational and End User Computing*, (16) 1, 2004, 59-72.
40. Taylor, S., Todd, P. A. Assessing IT Usage: The Role of Prior Experience. *MIS Quarterly*, 19(4), 1995, 561-570.
41. Barzilai-Nahon K., Scholl, J. Siblings of a Different Kind: E-Government and E-Commerce”, IFIP e-Government Conference, August, 2010.
42. Pavlou, P.: Consumer acceptance of electronic commerce: integrating trust and risk with the technology acceptance model. *International Journal of Electronic Commerce* 7, 69–103 (2003).
43. Sipior, J., Ward, B., Connolly, R. The digital divide and t-government in the United States: using the technology acceptance model to understand usage. *European Journal of Information Systems*, 19(1), 2010, 1-21.
44. Pilling, D., Boletzig, H.: Moving toward egovernment – effective strategies for increasing access and use of the internet among non-internet users in the U.S. and U.K. In *Proceedings of the 8th Annual International Conference on Digital Government Research: Bridging Disciplines & Domains Philadelphia*, (Cushing, JB and Pardo, TA, Eds) ACM International Conference Proceeding Series 228 Digital Government Research Center 2007, pp 35–46, ACM Publications (2007).
45. Simon, J.L. Basic research methods in social science; the art of empirical investigation. New York: Random House, 1969.
46. Carroll, J., Dawson, L.L., Swatman, P.A. Using Case Studies to Build Theory: Structure and Rigour. At *Proceedings of 9th Australasian Conference on Information Systems*, University of NSW, Sydney, Australia, 1998.
47. Carroll, J., Swatman, P. Structured-case: a methodological framework for building theory in information systems research. *European Journal of Information Systems*, 9, 2000, 235–242.
48. Dillon, W. R. Madden, T. J., Firtle, N. H. *Marketing Research in a Marketing Environment*, 3rd edition, 1994, 124-125.
49. Likert, R. A Technique for the Measurement of Attitudes, *Archives of Psychology*, 140, 1932.
50. Moore, G., Benbasat, I. Development of an instrument to measure the perceptions of adopting an information technology innovation. *Information Systems Research*, 2, 1991, 173–191.
51. Churchill, G. J. A Paradigm for Developing Better Measures of Marketing Constructs. *Journal of Marketing Research*, 16, 1979, 64-73.
52. Straub, D.W. Validating Instruments in MIS Research. *MIS Quarterly*, 13(2), 1989, 147-169.
53. Anderson, J.C., Gerbing, D.W. Structural Equation Modeling in Practice: A Review and Recommended Two-Step Approach. *Psychological Bulletin*, 103, 1988, 411-423.
54. Gerbing, D.W., Anderson, J.C. An Updated Paradigm for Scale Development Incorporating Unidimensionality and its Assessment. *Journal of Marketing Research*, 25, 1988, 186-192.
55. Kline, R.B. Principles and practice of structural equation modeling. The Guilford Press, New York, 1998.
56. Segars, A.H. Assessing the Unidimensionality of Measurement: a Paradigm and Illustration within the Context of Information Systems Research. *Omega*, 25(1), 1997, 107-121.
57. Gefen, D., Karahanna, E., Straub, D. Trust and TAM in online shopping: an integrated model. *MIS Quarterly*, 27, 2003, 51–90.
58. Straub, D., Boudreau M.C., Gefen, D. Validation Guidelines for IS Positivist Research. *Communications of the Association for Information Systems*, 13, 2004, 380-427.
59. J.F. Hair, R.E. Anderson, R. Tatham, W.C. Black, *Multivariate Data Analysis with Readings*, Macmillian, NY, 1992.
60. Fulk, J. Social construction of communication technology. *Academy of Management Journal*, 36, 1993, 921-950.
61. Fulk, J., Schmitz, J., Steinfield, C. W. A social influence model of technology use. Newbury Park, CA: Sage, J. Fulk & C. W. Steinfield Eds., *Organizations and Communication Technology*, 1990.
62. Neuman, W.L. *Social Research Methods: Qualitative and Quantitative Approaches*. Boston, MA, USA, 1991.
63. Gummerson, E. *Qualitative Methods in Management Research*. Newbury Park, CA, Sage, 1991.
64. Stamati, T., Kanellis, P., Martakos, D. Challenges of Complex Information Technology Projects: the MAC Initiative, *Journal of Cases on Information Technology*, IGI Publishing, 7(4), 2005, 41-58.

Image processing methods and algorithms for accurate protein spot detection in 2-dimensional gel electrophoresis (2DGE)

Panagiotis Tsakanikas¹

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
tsakanik@di.uoa.gr

Abstract. The main goal of this dissertation is the development of methods to improve the accuracy and efficiency of the protein spot detection and quantification on 2DGE images. Image analysis is still considered as the bottleneck of the differential expression proteomics analysis workflow, due to the large variability in the protein spots' expression profiles where a lot of manual user work is needed in order to achieve acceptable results. The contributions of this dissertation are apparent to all the stages of a 2DGE image analysis: (i) development of a new and specialized method for image denoising based on multiresolution analysis and the Contourlet Transform which was evaluated with synthetic and real images and shown that it outperforms the existing denoising approaches in terms of noise suppression and optimizing the subsequent analysis results; (ii) development of a novel approach for delineating 2DGE image areas which -with high probability- include protein spots, based on Active Contours. The developed method has been evaluated using a large pool of synthetic and real gel images and shown that the extracted ROIs include the large majority of the true spots while it functions in a fully automatic way; (iii) novel hierarchical approach to protein spot segmentation based on machine learning techniques and Gaussian mixture models. The developed approach is applied on each previously extracted Region of Interest (ROI), aiming at removing local background and streaks, while estimating the number, location and borders of proteins spots contained. After an exhaustive evaluation the developed methodology proved to be more accurate and efficient than competing methods while it is grounded on the physical properties of protein spots and overlapping spots formation.

Keywords: Proteomics, two-dimensional gel electrophoresis, denoising, image segmentation, spot detection, spot quantification, spot modeling.

1 Introduction

During the last decade, the life and computer science communities are striving to build models in order to develop a global understanding of the living cell. This effort is deeply influenced by the development of the “omics” technologies (genomics, transcriptomics, proteomics, metabolomics, etc), which aim at establishing a holistic view on biological systems. *Proteomics* is the large-scale study of proteins, and in particular of their structures and functions [1,2]. Proteins are vital parts of living cells, as they are the main components of the physiological metabolic pathways. The proteome is the entire set of proteins [3] expressed by an organism or system (including the modifications made to a particular set of proteins). This varies with time and depends on the stresses that a cell or organism undergoes. So, proteomics is the study of the time varying proteome using the technologies of large-scale protein separation and identification. In other words it is the study of proteins, how they are modified, when and where they are expressed, how they are involved in signaling and metabolic pathways and how they interact with each other. Current research in proteomics requires that proteins in a biological sample be effectively resolved. To achieve this goal, proteins need to be separated first. This separation can be performed using two-

¹ Dissertation Advisor: Elias S. Manolakos, Associate Professor

dimensional gel electrophoresis (2DGE) and gives rise to protein spots of irregular shapes and sizes on the gel. Once proteins are separated and quantified, they can be identified. For this purpose, individual spots are cut out of the gel and cleaved into peptides with proteolytic enzymes. These peptides can then be identified using mass spectrometry methods combined with database search.

Two-dimensional gel electrophoresis is a powerful and widely used method for the analysis of complex protein mixtures extracted from cells, tissues, or other biological samples. This technique is used to separate the proteins in two steps according to two independent properties: isoelectric focusing (IEF), which separates proteins according to their isoelectric points (pI), and SDS-polyacrylamide gel electrophoresis (SDS-PAGE), which separates proteins according to their molecular weights (MW). In this way, complex mixtures, consisted of thousands of different proteins, can be resolved and the relative amount of each protein can then be determined. The resulting gel is then digitized and an image analysis pipeline is applied in order to find protein spot specific properties (i.e. location, quantity, etc). In general, a typical 2DGE image processing pipeline consists of the following stages [4]:

1. Image pre-processing (noise suppression, artifacts removal, streak removal and background correction).
2. Spot segmentation (spot detection). Delineate each individual spot area, outputting a list of spot centers, intensities and geometric features. The spot detection operations pipeline in most cases is:
 1. Detect the centers of as many spots as possible.
 2. Segment the gel into regions, each containing one of these spots.
3. Modeling and quantification (spot volume estimation). Model each extracted spot region by a parametric spot model in order to extract a characteristic vector for each spot for further data analysis, and to detect and separate co-migrated spots.
4. Corresponding protein spot matching across gels of different samples.
5. Identification of differential expression (using statistical methods).

1.1 Related Work and Current Limitations

Although the resolution of 2DGE seems impressive, it is still not sufficient compared to the enormous diversity of cellular proteins, and co-migrating proteins in the same spot are not uncommon [5]. Neighboring spots can obscure protein spot centers in these so-called complex regions, and their saturated nature can make the resolution of each individual protein impossible. Furthermore, spots tend to have symmetric diffusion in the pI dimension but often severe tails in the Mr dimension (streaks). This diffusion depends on the protein concentration, which is why streaks and smears occur with certain proteins. So, the “faint” protein spots require an expert eye to discriminate them from noise, so an efficient noise suppression method would be extremely useful. Also, the intensity of the background can vary across the image. Finally, in most cases there are incompletely separated (overlapping) spots (less-defined and/or separated) and several complex protein spot areas. Inefficiency due to the above limitations leads also to unmatched/undetected spots (leading to missing values), mismatched spots, and errors in quantification (several distinct spots may be erroneously detected as a single spot by the software and/or parts of a spot may be excluded from quantification). All of the above constitute some of the big challenges for the automation of spot detection in the bioinformatics pipeline. Throughout the years a lot of commercial and non-commercial methods for image analysis of gel images have been developed.

Most commercially available 2DGE image analysis tools use conventional spatial filters [7,11] to combat with noise, mainly due to the fact that they are conceptually simple and computationally efficient. However, spatial filtering introduces severe distortions at protein spot borders and alters considerably the intensity values of internal spot pixels [6]. To address spatial filtering limitations, multiresolution space-frequency domain techniques based on wavelets have been proposed [6]. It has been shown that the Wavelet Transform (WT) outperforms spatial filtering both in terms of

signal-to-noise ratio (SNR) performance and in terms of the resulting visual image quality [6]. Despite its several advantages, the Wavelet Transform has also some notable limitations [12]. The most relevant to 2DGE image denoising is its limited ability in capturing directional information, as needed to adequately represent the smoothness along spot boundaries.

Another pre-processing operation is the background subtraction used in order to eliminate meaningless changes in the gel background intensity level. A simple approach is to obtain the lightest and darkest point in the background and replace the whole background with the average intensity. Tyson and Haralick [8] find the local minima in the image, representing background depressions, and interpolate the background between these minima. Melanie II [7] subtracts the minimum intensity from all pixel values and then fits a third degree polynomial to the background image (with spots removed). Another technique is derived from 3-D mathematical morphology, where the operations of opening and closing a grayscale image by a structuring element is represented by sliding the structuring element respectively under and over the topographical image (intensity is regarded as height) [9]. Moreover, using a horizontal and/or vertical cylindrical structuring element, horizontal or vertical streaks may also be removed [10].

Previous work in 2DGE image segmentation includes several single-phase direct segmentation methods that will be reviewed here. They include methods using stepwise thresholding [14], second derivatives [15], the Watershed transform [7], and statistical spot modeling methods [13]. The stepwise thresholding approach is extremely sensitive to noise and artifacts, where additional criteria must hold in order to accept or reject the final connected areas. The second derivatives approach gives acceptable results only when proper noise suppression has been applied. Furthermore, it places the borders at the inside of the spots since the zero crossings of the second derivative are associated with the steepest part of the spot rather than its beginning. The Watershed transform based method has the major disadvantage of over-segmentation. Although this can be addressed using marker controlled watersheds, the selection of a good set of markers is not a trivial task. Finally, the approaches using statistical spot modeling are difficult to apply without prior knowledge of spot shapes and sizes and it is known that they perform poorly in areas with overlapping spots especially if these foreground areas are not accurately estimated (usually this estimation is performed by mathematical morphology). So, it is obvious then that the current protein spot detection approaches suffer from various disadvantages, such as sensitivity to noise and artefacts, spot border distortions, over-segmentation, poor performance in areas with overlapping spots [4] etc. Furthermore, they require careful post-processing and usually a lot of manual effort, to finally produce reliable detection results.

1.2 Dissertation Contributions

The goals of this dissertation have been the development of novel methods for 2DGE image analysis and especially for spot denoising, detection and quantification in order to improve the accuracy and efficiency of existing methods. Image analysis is still a bottleneck in expression proteomics workflows. Nowadays, there are several commercial software packages available such as PDQuest, ImageMaster, Progenesis, etc, that promise to be accurate and efficient but this is far from being a reality [16]. Motivated by the aforementioned limitations the main contributions are:

- ✓ *2DGE image denoising*

Since 2DGE gel images are inherently noisy due to dust and the imperfect image acquisition process, the first objective of this dissertation is the development of an effective denoising method, i.e. increasing the SNR without inserting significant distortions to the image. In this dissertation, a multiresolution image transform was employed, namely the Contourlet Transform (CT), which proved to be very effective for denoising 2DGE images and fit well the specific properties of gel images. The CT can approximate more accurately images with smooth contours and anisotropic characteristics. 2DGE images are anisotropic due to the large variety in shapes and orientation of the spots they contain. The developed method is fully automated and it is shown to outperform every

previously reported method [17,18]. We must note that the CT has not been used before for this type of images nor it has been coupled with the coefficient thresholding techniques that we have used.

✓ *Novel Active Contours based method for extracting foreground Regions of Interest (ROIs)*

A novel methodology [19] has been developed for delineating 2DGE image areas which, with high probability, include protein spots, based on Active Contours without Edges (ACWE). Moreover, a technique based on Contourlet Transform has been developed that leads to the automatic determination of the initial curve. This initialization method reduces the convergence time of the algorithm and improves its efficiency. Due to fact that the Contourlet Transform has properties that match particular characteristics of 2DGE images, a method based on it has been developed in order to enhance gel images and especially the faint spots. The method has been evaluated using a large pool of synthetic and real gel images. It has been shown that the extracted ROIs include the large majority of the true spots and are tight, i.e. they do not include large background areas. The evaluation has been performed using the popular commercial software package PDQuest and also relatively to the provided ground truth. Furthermore, our method does not require re-calibration of parameters every time a new image is processed and it can thus be fully automated.

✓ *Novel hierarchical approach for protein spot detection & quantification using machine learning methods*

This approach [20,21], unlike the traditional spot detection workflow where a gel image is directly segmented into spot regions following the spot modelling phase, it is applied on each ROI resulting from the previous image analysis step. First, it removes the local background pixels and streaks using 1-dimensional Gaussian mixture models applied on the intensity histograms of the extracted ROIs. Unlike mathematical morphology filtering it “kills” streak pixels without affecting the true spot pixels. A key idea of the developed method is that the informative image pixels are treated as sample data generators where machine learning methods are applied to the so generated data samples. A core technique used repeatedly in the developed methodology is Gaussian Mixture Modelling (GMM) [22] which is applied in an unsupervised manner [23]. Through an extensive evaluation, we have demonstrated that the developed methodology achieves trustworthy detection results and introduces much less spot artifacts. In addition, a comparison with the popular commercial software package PDQuest was conducted and shown that the developed methodology is more precise and more specific than PDQuest, while both methods achieve high sensitivity. Furthermore, it has been shown that it leads to more accurate spot quantification than PDQuest. Finally, the developed methodology can be fully automated and thus it is labor and error free from the user’s perspective, which is very important for high throughput proteomics projects.

2 Image Analysis

2.1 Image Denoising

The denoising methods commonly used so far, have the tendency to deform the protein spots on the gel to the extent that they create extraneous spots i.e. artifacts. This is a serious problem since insufficient or improper denoising affects the whole image processing pipeline from its early stages. So, it impacts negatively all the subsequent processes, such as spot detection, spot quantification, as well as spot matching across gels. In order to surpass those problems, a novel method for denoising 2DGE images has been developed, based on the Contourlet Transform [12]. The Contourlet Tranform (CT) is a multiresolution, flexible, directional image decomposition method based on contour segments. The main difference between the CT and the WT is that the CT allows for a different number of directions at each scale (can be any power of 2). So, CT can represent more efficiently smooth spot contours in a 2DGE image.

The denoising by multiresolution transforms involves the aforementioned analysis of signal followed by a coefficient thresholding method (also called *shrinkage*). For the developed CT-based

2DGE denoising methodology, two of the best performing shrinkage methods reported in the WT literature were adopted, namely the *BayesThres* [24] and *Bivariate* [25] methods.

As it is demonstrated the Contourlet Transform has properties that match well the characteristics of 2DGE images, and after a thorough evaluation with both synthetic and real gel images in terms of the achieved Signal to Noise Ratio (SNR), the distortions introduced and mainly via the benefits it offers to subsequent image analysis steps:

- 1) Protein spot detection - where by using the developed denoising methodology we avoid introducing a large number of artifacts and detect more faint spots.
- 2) Protein spot quantification - the estimated spot quantities are more close to the known ground truth and with less variance than when using wavelet-based denoising.

In conclusion, the developed denoising methodology is more effective than the currently used spatial filtering methods implemented by commercial software packages and also more effective than the more recently introduced Wavelet-based denoising methods.

Next we present the results of the developed method and compare them with the currently state-of-art method (wavelet denoising). We used PDQuest (version 8.0.1) to evaluate the different denoising approaches in terms of spot detection achieved after denoising using real images. We evaluated each method in terms of the TPs, FPs (artifacts) and false negatives (FNs) or missed spots. The results are summarized in Figure 1. We notice that regardless of the image used, the CT based denoising approaches result to a considerably smaller percentage of introduced extraneous spots (FPs), ranging from 4% to 8%, compared to the WT based denoising methods where the corresponding range was from almost 6% to 15%. The missed spots (FN) were also less when using the CT in all cases except for the CT-Bivariate and GelA case where they approached 4%. Overall, CT-Bayes denoising outperformed all other methods in terms of TPs, FNs and FPs and for both real images used (see Figure 1(c)).

2.2 Active Contours based method for extracting foreground Regions of Interest (ROIs)

Segmentation of 2DGE images requires partitioning them into areas of foreground (include protein spots) and background (no protein spots). We developed a new method based on Active Contours [26] that separates effectively those two areas in a way that: (i) reduces the number of missed faint spots, (ii) finds correct and tight borders for areas with spots, (iii) avoids over-segmentation.

Active contours (ACs) are a very powerful tool for image segmentation and object tracking. The key idea is the evolution of a curve, or curves, also called “snakes”, subject to constraints from the input image. Due to the properties of the specific application evolving curves should allow automatic topological changes of the curve. An AC approach that holds this property was introduced in [27] where the curve is modelled as a specific level set function of time in a higher dimensional surface. For more details on the developed methodology reader is referenced to [19].

Next we present some results of the developed methodology while for more details the reader is referenced to [19]. The results have also been compared to the ones obtained with PDQuest. The denotation followed is: spots that AC missed but were found by PDQuest (PDQ/nAC) partitioned into two subsets; existing spots that our AC base method missed (false negatives, FN) and extraneous spots (true negatives, TN = PDQuest artefacts) that AC correctly ignored. For those spots that AC detect but PDQuest missed (denoted as nPDQ/AC) we also consider two subsets. The ones that AC correctly found because they do exist (true positives, TP_2), and those that AC found but do not exist (false positive, FP = AC artefacts). The results obtained are summarized in Tables 1 and 2. Knowing that PDQuest is performs very well in segmentation we can conclude that our approach correctly reports the foreground regions in a wide collection of different of 2DGE images. We also see that we succeed in avoiding some PDQuest detected artefacts (TN) but also fail to detect very few spots that should have been detected (FN). On the other hand, we detect areas that contain spots missed by PDQuest (TP_2), but also introduce some artefacts (FP). In Table 1 we can see that the ratio of spots missed by our approach compared to PDQuest is pretty small (<3%, except for the Rj1

image) which also indicates that the proposed method results are highly reliable. In addition, as we can see from Table 2, that the proposed approach achieves sensitivity over ~91% for all images and a confidence above ~96%. These results indicate that ACs can be very effective in confining protein spots into tightly bounded spot areas. Accurate and correct spot areas segmentation is a prerequisite for spot detection and quantification. We have shown that the proposed AC based segmentation achieves comparable results with a mature tool for 2DGE image analysis (PDQquest) but with much less user intervention.

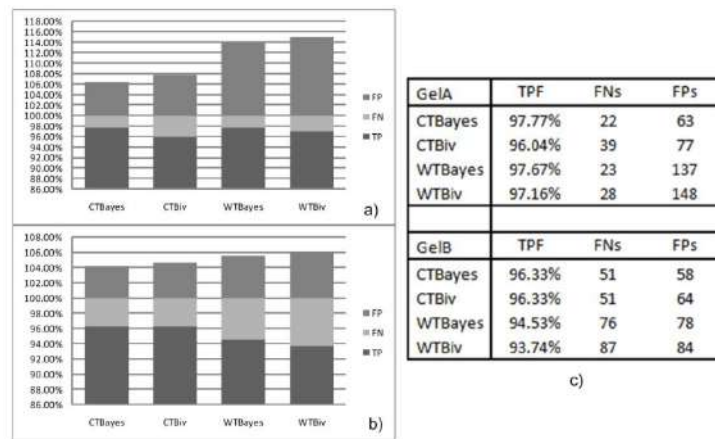


Figure 1 Spot detection results using the real images from [27] a) GelA and b) GelB respectively, c) the corresponding True Positive Fraction (TPF), False Positives (FPs) and False Negatives (FNs).

Image	PDQ	PDQ/AC	% PDQ/AC	PDQ/nAC	% PDQ/nAC	nPDQ/AC	%nPDQ/AC
1a	1112	1090	98,02%	22	1,98%	13	1,19%
2a	1315	1283	97,57%	32	2,43%	7	0,55%
MP1	262	256	97,71%	6	2,29%	13	5,08%
MP2	265	242	91,32%	23	8,68%	11	4,55%
MP3	227	223	98,24%	4	1,76%	24	10,76%
Rj1	146	123	84,25%	23	15,75%	4	3,25%
RGA	948	919	96,94%	29	3,06%	9	0,98%
RGB	1040	1018	97,88%	19	1,83%	40	3,93%

Table 1 Evaluation results. PDQquest spots in AC extracted foreground regions were all real spots ($PDQ/AC = TP_1$). We also report spots detected by PDQquest and not included in our foreground areas ($PDQ/nAC = (TN+FN)$) and spots in our foreground areas not detected by PDQquest ($nPDQ/AC = (TP_2+FP)$).

Image	PDQ	AC/PDQ (TP_1)	PDQ/nAC		nPDQ/AC		S	C
			FN	TN	FP	TP_2		
1a	1112	1090	5	17	5	8	99,55%	99,55%
2a	1315	1283	9	23	5	2	99,30%	99,61%
MP1	262	256	2	4	10	3	99,23%	96,28%
MP2	265	242	14	9	4	7	94,68%	98,42%
MP3	227	223	1	3	6	18	99,59%	97,57%
Rj1	146	123	11	12	1	3	91,97%	99,21%
RGA	948	919	20	9	6	3	97,88%	99,35%
RGB	1040	1018	12	7	13	27	98,86%	98,77%

Table 2 Evaluation results. Sensitivity is above 91% and Confidence above 96% for all images.

2.3 Protein Spot Detection & Quantification

Previously developed methods for protein spot detection suffer from various disadvantages, such as sensitivity to noise and artifacts, spot border distortions, over-segmentation and poor performance in areas with overlapping spots [4]. Furthermore, they require careful post-processing and a lot of manual effort, to finally produce reliable detection results [16]. To address these limitations, a novel approach for 2DGE automatic spot detection and quantification has been developed. Here, the informative ROI pixels are treated as sample data generators and Gaussian Mixture Modeling (GMM) is applied in an unsupervised manner [29,30], i.e. no pre-training is required and the whole process can be fully automated.

2.3.1. Local Background and Streaks removal

Although the ROI extraction step results in areas that include the vast majority of the protein spots present in a gel image, they may also include some local image background pixels and/or streak segments. In order to eliminate those pixels that may confuse spot modeling downstream from further processing, we classify the object pixels into 3 possible classes: class-1 represents the “strong” and/or saturated spot pixels, class-2 the “faint” spot pixels and tails of “strong” spots, and finally class-3 correspond to the pixels of the local background and/or streaks. The classification is performed on the histogram of pixel intensities using 1-dimensional Gaussian mixture modeling and a modified Expectation-Maximization (EM) algorithm [30]. The Minimum Message Length (MML) criterion (proposed in [30] for model selection) is applied to determine the optimal number of components in the mixture model that best fits the histogram data.

In summary, the use of the modified EM fits 3 pixel classes to the histogram while the MML criterion tests this fit and discards any class that is not essential for the histogram’s interpretation. So, if the histogram is well explained by 3 classes we can then threshold out the class-1 pixels that correspond to the background and/or streaks (light gray in Figure 3 (3)), else (if we end up with less than 3 classes) we keep the object intact. As one can notice, the developed method is performing local background and streak removal task only in areas and to the extent that it is really needed.

2.3.2 Initial Spot centers Estimation

At the next step we estimate the number of protein spots existing inside each ROI. To do so, a 5x5 (pixels) spatial filter is employed that finds the local minima (zero intensity corresponds to black pixels, maximum intensity to white) in the image (see green asterisks in Figure 3 (5)). Due to the pixel intensity saturation effect, it is possible that the filter identifies several closely located local minima. These are actually replicates of the same candidate spot centre and need to be grouped appropriately so as we do not end up with a very large and misleading number of candidate spot centers per object. This is accomplished by applying agglomerative Hierarchical Clustering (HC) [30] of the minima points using the Manhattan distance, the single linkage method for merging formed clusters. The extracted candidate spot centers are depicted by red circles in Figure 3(5).

2.3.3 2D Spot Modeling

The next and most distinguishing characteristic of our pipeline is the idea of using the pixel intensities as data generators. The total number of data points N generated by random sampling for each ROI is kept proportional to the number of estimated candidate spot centres, and not to the area of the ROI. The N points to be drawn are distributed among the pixels of the object according to their relative intensities (a “stronger” pixel “throws” more points in its neighbourhood). This is in accordance to the view that pixel intensity ideally represents the quantity of protein molecules concentrated at that particular gel location. Specifically, each pixel i with intensity I_i acts as the centre $\mu=(x_i, y_i)$ of a 2D Gaussian component $N(\mu, \Sigma)$ in a Gaussian Mixture Model [28] having as many components as the number of pixels and a predetermined fixed covariance matrix Σ .

As we can notice from Figure 3 (6) (light grey data points) the generated set of data points may include points that are far from all candidate spot centres. These outlier points (due to

remaining background pixels) can be identified since they have low likelihood for all components of the mixture (no component “feels strongly” about them) and are removed at this stage since they may adversely affect the subsequent step of spot detection and quantification.

The last step in the pipeline is the application of Gaussian Mixtures Models (GMM) [28] in 2-dimensions (see figure 3(7) & 3(8).

2.3.3 Results & Discussion

2.3.3.1 Spot detection & Quantification evaluation

In this Section the results of the developed methodology are presented. Figure 4 presents a summary of the results. It can be noticed that the developed method achieves a high TPF (over 91% for all images) and a high PPV (over 79% for all images), meaning that it is both sensitive and precise. This conclusion is also supported by the fact that the images exhibit very different characteristics.

This is because these images exhibit a larger dynamic range difference between spots and background. So, Active Contours based segmentation performs better at complex areas leaving out of the ROIs only a small portion of faint spots, unlike PDQuest which in order to achieve high detection efficiency needs a larger sensitivity parameter value which leads to a very large number of extraneous artefacts (greater than 300 artefacts).

Quantification of protein spots is also a very crucial step in 2-DE image analysis. Its accuracy and reliability greatly affects the subsequent differential spot expression analysis. Figure 5(a) presents the results obtained for the noise free synthetic image for which the ground truth (spot volumes) is known. This is a scatter plot of estimated vs. ideal spot volumes, as produced by PDQuest (data points denoted by circles) and the developed method (data points denoted with an x). The ideal regression line would be the diagonal line $y=\alpha x$ with $\alpha=1$. As it can be seen in Figure 5(a), the regression lines for each method (dotted line for PDQuest and solid line for developed method) deviate from the ideal. In the case of PDQuest, where $\alpha=1.906$, it is clear that there is an overestimation of the true spot volumes, which according to the table in Figure 5(b) (first column) has the mean value of 23.6 units. On the other hand, the developed method has $\alpha=0.8655$ i.e. it underestimates the true volume with a mean error of -4.6 units. Someone could argue that since the bias introduced by each method affects all the spots in a gel image it is of no great concern to differential expression analysis. However, regardless of the type of bias inserted by each method, the most important result is the error’s standard deviation (see Figure 5 (a) and (b)). As we can see, the developed methodology exhibits a low standard deviation of 0.42 while the STD of PDQuest estimated spot volumes for the same image is more than ten times higher (6.38). The corresponding Root Mean Square Error (RMSE) is 4.77 and 30.03 respectively. The much lower STD and RMSE values of quantification error suggest that when using the developed methodology the variance of the produced quantity estimates across spots with similar characteristic is much smaller. So, the developed approach is highly consistent, in the sense that a gradual decrease in spot maximal intensity (which also reflects in the spot’s area in this image) leads to a corresponding gradual decrease in volume estimates (see Figure 5(a)). This characteristic is what is the most important in practice for a reliable differential analysis based on spot volumes, irrespectively of any bias introduced by the method. Finally, at Figure 5(e) someone can see how noise addition and subsequent denoising affects spot detection performance for each method. It is clear that the developed method exhibits robustness to noise in detecting faint spots and avoiding introducing artefacts. As someone can notice, the developed methodology exhibits a high degree of robustness in detecting faint spots. Furthermore, it avoids introducing artifacts. It should also be mentioned that the developed method does not need any user intervention while PDQuest requires careful selection of its detection parameters every time the image is changing.

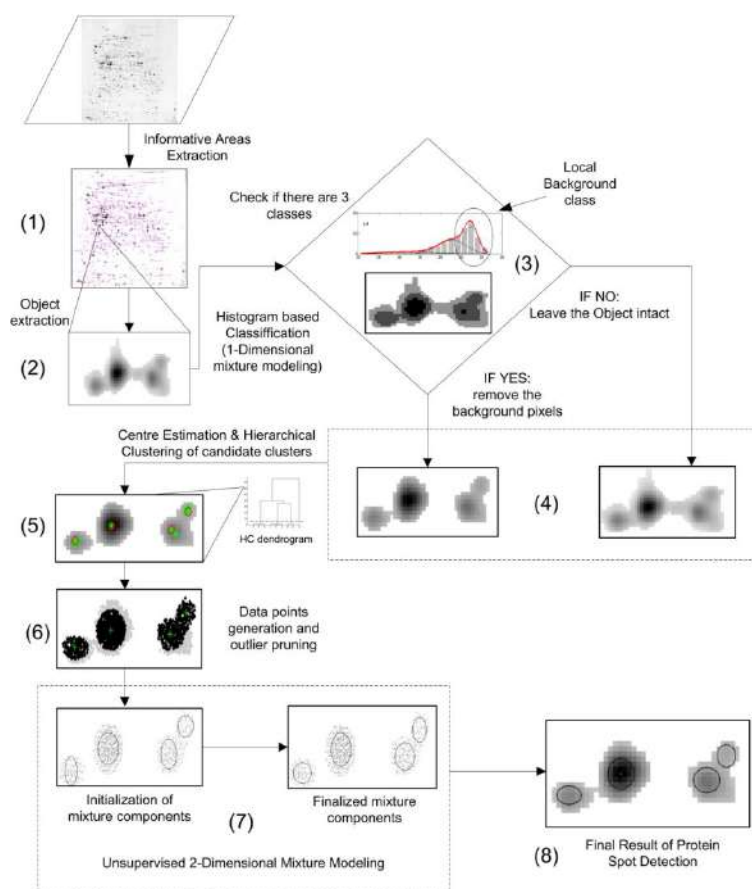


Figure 3: Overview of the developed spot detection method. Panels: **(1)** Extraction of areas containing protein spots using Active Contour based first level Segmentation, **(2)** Extracted image object, **(3)** Histogram based classification using 1-Dimensional mixture modelling (upper image shows the histogram of the object intensities, red line indicates the joint mixture density consisting of the individual densities represented with solid lines (1-D Gaussians); bottom image presents the resulting classified regions where the light gray area indicate the third class (local background) to be eliminated), **(4)** Resulting object after the histogram classification, **(5)** Centre estimation (green asterisks) and merging (red circles) of the neighbouring replicate centres with hierarchical clustering (dotted line at the dendrogram on the right), **(6)** Data generation and manipulation, **(7)** Unsupervised 2-Dimensional mixture modelling and final estimate, **(8)** Final detection result on the image object.

3 Conclusions

The goals of this dissertation have been the development of novel methods for 2DGE image analysis and especially for spot denoising, detection and quantification in order to improve the accuracy and efficiency of existing methods. Image analysis is still a bottleneck in expression proteomics workflows. Nowadays, there are several commercial software packages available such as PDQuest, ImageMaster, Progenesis etc, that promise to be accurate and efficient but this is far from being a reality [16]. In order to surpass current limitations we developed an end-to-end image processing pipeline which addresses effectively the current bottlenecks of the available 2DGE image analysis methods. We proved its robustness and efficiency using real and synthetic datasets while also evaluated its performance by means of what matters most for the proteomics scientists more than for engineering scientists. Finally, the developed methodology can be fully automated and free the user from the time consuming tasks of parameter fine tuning and spot editing.

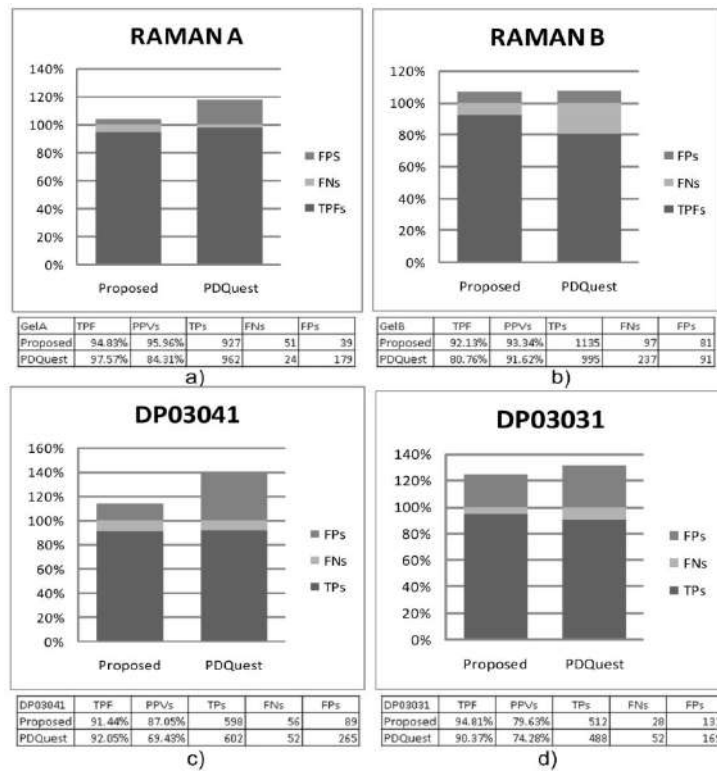


Figure 4 a) Comparison of results for the image *RamanaA* using the developed method and the PDQuest software package in terms of True Positives (TPs), False Negatives (FNs) and False Positives (FPs); b) for image *RamanB*; c) for image *DP03041*; d) for image *DP03031*.

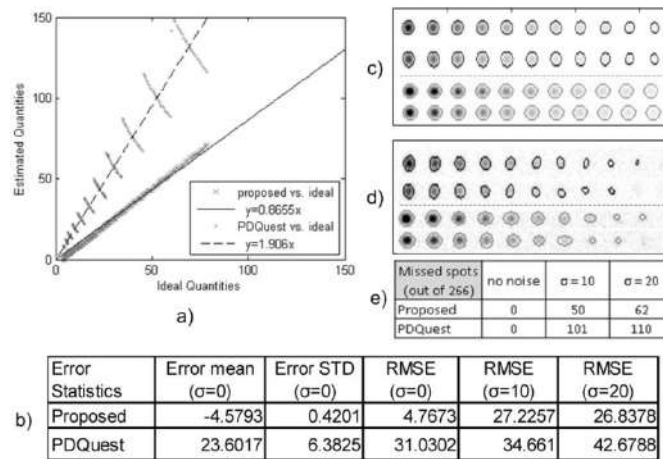


Figure 5: Comparative evaluation of the developed method to PDQuest for spot quantification performance. **a)** Estimated vs. ideal spot quantities using the synthetic image *Quant*. **b)** The table provides the quantity error means and standard deviations in the case of noise free image along with the corresponding root mean square error (RMSE) for all noise cases. **c)** Image patches with detected spots from the noise-free image (the top two spot rows correspond to results of the developed method, while the bottom two spot rows to PDQuest results). **d)** Image patches with detected spots from the noise corrupted image with $\sigma=10$ (top two spot rows correspond to results of our method while the bottom two spot rows to PDQuest results). **e)** Spot detection performance (missed spots) of each method after denoising, for high and very high noise levels.

4 References

- [1] Anderson NL, Anderson NG, "Proteome and proteomics: new technologies, new concepts, and new words", *Electrophoresis* 19 (11): 1853–61, 1998.
- [2] Blackstock WP, Weir MP, "Proteomics: quantitative and physical mapping of cellular proteins", *Trends Biotechnology*, 17 (3): 121–7, 1999.
- [3] Marc R. Wilkins, Christian Pasquali, Ron D. Appel, Keli Ou, Olivier Golaz, et al., "From Proteins to Proteomes: Large Scale Protein Identification by Two-Dimensional Electrophoresis and Amino Acid Analysis". *Nature Biotechnology* 14 (1): 61–65, 1996.
- [4] Dowsey, A. W., High-throughput image analysis for proteomics. PhD Thesis 2005, Department of Computing, Imperial College London.
- [5] Pietrogrande, M. C., Marchetti, N., Dondi, F., Righetti, P. G., "Spot overlapping in two-dimensional polyacrylamide gel electrophoresis separations: a statistical study of complex protein maps", *Electrophoresis*, 2002, 23, 283-291.
- [6] Kaczmarek, K., Walczak, B., de Jong, S., Vandeginste, B. G. M., Preprocessing of two-dimensional gel electrophoresis images. *Proteomics* 2004, 4, 2377-2389.
- [7] Appel, R. D., Vargas, J. R., Palagi, P. M., Walther, D., et al., "Melanie II – a third-generation software package for analysis of two-dimensional electrophoresis images: II. Algorithms", *Electrophoresis*, 1997, 18, 2735-2748.
- [8] Tyson, J. J., Haralick, R. H., "Computer analysis of two-dimensional gels by a general image processing system", *Electrophoresis*, 1986, 7, 107-113.
- [9] Sternberg, S. R., "Gray scale morphology", *Comp. Vis. Graph. Image Processing*, 1986, 333-355.
- [10] Skolnick, N. M., "Application of morphological transformations to the analysis of two-dimensional electrophoresis gels of biological materials", *Comput. Vis. Graph. Image Process.*, 1986, 35, 306-322.
- [11] M. Rogers, J. Graham, and R.P. Tonge, "Using statistical image models for objective evaluation of spot detection in two-dimensional gels", *Proteomics*, 2003, vol. 3, pp. 879- 886.
- [12] M.N. Do and M. Vetterli, "The contourlet transform: An efficient directional multiresolution image representation", *IEEE Trans. on Image Proces.*, 2005, vol. 14, pp. 2091-2106.
- [13] Conradsen, K., Pedersen, J., "Analysis of two-dimensional electrophoretic gels", *Biometrics*, 1992, 48, 1273-1287.
- [14] P. Cutler, G. Heald, I. R. White, J. Ruan, "A novel approach to spot detection for two-dimensional gel electrophoresis images using pixel value collection", *Proteomics*, vol. 3, pp. 392-401, 2003.
- [15] Lemkin, P. F., Lipkin, L. E., "GELLAB: a computer system for two dimensional gel electrophoresis analysis. III. Multiple two-dimensional gel analysis", *Comput. Biomed. Res.*, 1981, 14, 407-446.
- [16] Clark, B. N., and Gutstein, H. B., The myth of automated, high-throughput 2DGE image analysis, *Proteomics* 2008, 8, 1197-1203.
- [17] Tsakanikas, P., Manolakos, E. S., Improving 2-DE gel image denoising using Contourlets, *Proteomics* 2009, 9, 3877–3888.
- [18] Tsakanikas, P., Manolakos, E. S., Effective Denoising of 2D Gel Proteomics Images Using Contourlets *IEEE Proc. ICIP 2007*, San Antonio, Texas, USA, September 16-19, 2007, pp. VI: 269-272.
- [19] Tsakanikas, P., Manolakos, E. S., Active Contour Based Segmentation of 2DGE Proteomics Images, 16th European Signal Processing Conference (EUSIPCO-2008), Lausanne, Switzerland, August 25-29, pp. 83-87, 2008.
- [20] Tsakanikas, P., Manolakos, E. S., Protein Spot Detection and Quantification in 2-DE gel images using Machine Learning methods, *Proteomics* 2011, accepted, to appear.
- [21] Tsakanikas, P., Manolakos, E. S., "A fully automated 2-DE gel image analysis pipeline for high throughput proteomics", *International Conference on Acoustics, Speech and Signal Processing (ICASSP 2011)*, May 22-27, 2011, accepted, to appear.
- [22] McLachlan, G.J. and Peel, D. *Finite Mixture Models*, Wiley (2000).
- [23] Figueiredo M., and Jain A.K., Unsupervised learning of finite mixture models, *IEEE Transactions on Pattern Analysis and Machine Intelligence - PAMI*, vol. 24, no. 3, pp. 381-396, March 2002.

- [24] Abramovitch, F. F., Sapatinas, T., Silverman, B., Wavelet thresholding via a Bayesian approach, J. R. Stat. 1998, 60, 725-749.
- [25] Sendur, L., Selesnick, I. W., Bivariate shrinkage functions for wavelet-based denoising exploiting interscale dependency. IEEE Trans. on Signal Processing 2002, 50, 2744-2756.
- [26] T. F. Chan, L. A. Vese, "Active Contours without Edges", IEEE Trans. on Image Processing, vol. 10, pp. 266-277, Feb. 2001.
- [27] S. Osher, J. A. Sethian, "Front Propagating with curvature dependent speed: Algorithms based on Hamilton-Jacobi Formulation", J. Comput. Phys., vol. 79, pp. 12-49, 1998.
- [28] Raman, B., Cheung, A., Marten, M. R., Quantitative comparison and evaluation of two commercially available, two-dimensional electrophoresis image analysis software packages, Z3 and Melanie. Electrophoresis 2002, 23, 2194-2202.
- [29] McLachlan, G.J., Peel, D. Finite Mixture Models, Wiley 2000.
- [30] Figueiredo M., and Jain A.K., Unsupervised learning of finite mixture models, IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 24, no. 3, pp. 381-396, March 2002.
- [31] Theodoridis, S. and Koutroumbas, K., Pattern Recognition (Third Edition), Elsevier, 2006

Semantic Information Management for Pervasive Computing

Vassileios Tsetsos*

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
b.tsetsos@di.uoa.gr

Abstract. In recent years, knowledge and semantics management has been applied to many areas of Informatics and Telecommunications. The main reason is that such management enables the creation of more flexible and usable computing environments, with higher added value for the users. In this thesis various different cases of semantics information management are studied, that refer to different aspects of modern and future computing paradigms. A special focus is put on context information management, when it is described through knowledge representation techniques. Specifically, in the first part of the thesis several methods for semantic service discovery are presented. Moreover, an evaluation framework for service discovery engines is proposed, which satisfies the special requirements of this discovery process. This framework, as shown by the experimental evaluation performed, is appropriate for evaluating such engines, since it is fully compatible with their special characteristics. In the second part the combination of context-awareness and information dissemination in autonomic computing environments is explored. The main goal for such combination is to be able to achieve collaborative context-awareness. The main contribution of this thesis is the design and evaluation of a efficient scheme for collaborative context-awareness. Finally, in the third part, a framework for personalized services is presented along with two case studies: semantic location services and personalized interactive TV. The proposed framework relies on rules and ontologies for providing advanced services to TV subscribers. In this part several issues of architectural and technological nature are addressed that are closely relevant to the implementation of such services.

Keywords: Pervasive Computing, Semantic Information Management, Context-Awareness, Mobile Computing, Semantic Services, Service Personalization

1 Introduction

The main objective of this thesis is to study ways to develop intelligent computing environments and advanced services. The techniques presented in the following sections are related to various concepts which are briefly described as follows:

* Dissertation Advisor: Stathes Hadjiefthymiades, Assistant Professor

- “Pervasive Computing”: the term was first mentioned by Mark Weiser back to 1991. It refers to a new computing paradigm where computers, embedded in the physical user environment, are cooperating in a distributed way in order to provide an advanced user experience. Some key characteristics of this paradigm are a) system adaptation to changes in the user/application environment, b) user-centered and personalized applications, c) intelligent services.
- “Knowledge management and Semantic Web”: semantic information is formally structured metadata that clearly define concepts and entities (through relationships, constraints, properties etc). Management of semantic information (or knowledge) refers to its representation, persistence, and update processes. A formal, yet practical, way of semantics representation is ontology. According to R. Studer (1998) “Ontology is a formal, explicit specification of a shared conceptualization”. In practise, the most popular means for creating an ontology is the Semantic Web technologies (e.g., Web Ontology Language, Resource Description Framework Schema)
- “Context awareness”: it refers to a system’s ability to sense and react accordingly to changes in its operational parameters. A system may sense the context changes through sensors or other means. It can even rely on other systems to update it upon context changes, a technique called collaborative context awareness that was studied in one of the thesis’ topics. Another important aspect of context awareness is its context modeling. Several context models have been proposed in the literature. In this thesis we explore knowledge representation methods for describing context.

Our thesis is that all the aforementioned aspects of computing, and especially the semantic technologies, form the basis for creating advanced and intelligent systems with the following main features:

- Distributed computation in heterogeneous environments.
- Collaborative computation
- Intelligent service behaviour and application personalization

1.1 Organization and contribution of the doctoral dissertation

The thesis is divided in three parts that study different aspects of semantic information management in pervasive computing environments. Special emphasis is put on investigating issues related to the features mentioned in the previous section.

The first part of the thesis deals with service discovery and, particularly, its evaluation. Since service discovery is a core functionality of pervasive computing, we aimed to design a framework for assessing the effectiveness of service discovery systems.

Two topics were investigated:

- a) finding appropriate metrics for evaluating service discovery systems
- b) creating a methodology for evaluating discovery systems with graded relevance scale and without relevance judgments from experts.

The contribution of our research can be summarized to:

- a) two metrics that are used in generalized information retrieval systems were proposed. These are suitable for evaluating graded discovery results and are more accurate than the standard Precision and Recall metrics.

b) an evaluation method that combines various techniques from the information retrieval field was proposed. It does not require relevance judgments from experts since it automatically generates pseudo-relevance judgments. Then it exploits the aforementioned metrics in order to evaluate and compare the systems.

The second part of the thesis refers to collaborative context awareness through efficient message exchange in nomadic (ad hoc) environments (e.g., vehicular networks). In such environments, there is often the need for intelligent context aware services. An efficient way to achieve this is through knowledge representation and reasoning techniques. However, in such environments not all nodes do have sensors and their resources (e.g., energy, storage) are usually limited. Hence, the nodes have to collaborate efficiently in order to deliver the expected functionality. The contribution of the thesis in this area is a publish/subscribe scheme where the nodes exchange sensor-originating data [1]. The data model adheres to a context ontology, thus enabling support of knowledge management processes. The scheme makes minimal assumptions for the underlying network infrastructure and nodes and tries to be sensitive in terms of event detection (i.e., context changes). Its performance is shown to be superior than that of other schemes which can be applied to nomadic environments (i.e., periodic polling).

In the third part of the thesis we study issues related to personalization of context-aware services and applications, with the aid of semantics. Specifically, a generic framework is presented for developing such applications. Basic elements of such framework are models that represent all the involved entities, with more important being the user model. Modelling of such applications is based on ontologies and the actual personalization process is implemented through rules and the respective reasoning engines. The proposed framework was validated in two application domains: a) semantic location based services (navigation and other services in users with disabilities), and b) personalization of multimedia services in interactive TV environments (content and service provisioning based on the program semantics and the viewer's profile). In the area of semantic location based services, a user model is combined with a location (spatial) model so that the system can search for accessible paths, taking into consideration the user's disabilities or preferences. Usage of such extensive knowledge representation techniques in this application domain is considered as one of the thesis' key contributions [2][3].

The same applies to the second application domain, too [4] (i.e., semantic interactive TV). There have been proposed several similar systems in the literature [5], however none of them uses formal semantics technologies in order to provide the desired personalization. Such declarative way to affect the provisioning of value added services is, in our view, a key factor for their adoption and success.

2 Evaluation of Semantic Web Service Discovery

Service discovery and selection are central topics in distributed systems research. Semantic Web Services (SWS) are an evolution of Web Services that are based on metadata and allow for more expressive description of service capabilities, used both for service advertisements and requests. Such metadata is represented through well-

known knowledge representation tools, like ontologies and rules, implemented with Semantic Web technologies. The SWS paradigm involves, apart from the aforementioned description facilities, more sophisticated (mainly logic-based) matchmaking algorithms [6]. Typical information that the SWS discovery (SWSD) systems exploit are attributes such as: service Inputs, Outputs, Preconditions and Effects (a.k.a. IOPE attributes). However, no matter on which elements of a service description the matchmaking algorithm is applied to, the most important problem in matchmaking is that it is unrealistic to expect relevant advertisements and requests to be in perfect match. The problem is aggravated if we take into account that the service request may not fully capture the requestor's intention. The concept of the "Degree of Match" (DoM), a kind of a relevance scale, was introduced for dealing with this problem [6].

Similarly to other retrieval systems, such as Web search engines, SWS discovery systems should be evaluated in terms of performance and retrieval effectiveness. Many researchers have already undertaken performance assessment efforts for measuring retrieval times and the scalability of the available tools (an extensive review of SWSD engines can be found in [6]). However, to the best of our knowledge, only a few researchers have performed such experimental evaluations. There are several reasons for this situation, with the most important outlined below:

1. *Lack of established evaluation metrics.* The typical metrics used in Information Retrieval are not directly applicable and they do not fully take into account the semantics of the degrees of match. Moreover, the existing metrics do not assume service rankings with ties (aka weak or partial rankings) in their majority. Such rankings are typically returned by the discovery engines.
2. *Lack of (sufficiently large) test collections.* Current service test collections, either for plain WS or SWS have a small number of services and an even smaller number of queries and relevance judgments [7].
3. *Incomplete relevance judgments.* In general, SWSD should be evaluated with methods that assume incomplete or totally missing relevance judgments. This is a very important and realistic assumption in open environments like the Web.

2.1 Background on Evaluation of Service Discovery Processes

The entities involved in a service discovery process are the service advertisements (S_i) published in a service registry, the service request R posed by the user, and the matchmaking engine that is responsible for the actual service discovery. In essence, the matchmaking engine assigns a Degree of Match $e(R, S_i)$ to every service advertisement S_i . These values determine the ranking of the final advertisements for a specific request R . In order to evaluate the matchmaking engine effectiveness some expert mappings $r(R, S_i)$ (i.e., relevance judgments between R and each S_i) should be available/pre-specified. Hence, the vectors r and e are defined as:

$$r: Q \times S \rightarrow W, \quad e: Q \times S \rightarrow W$$

where Q is the set of all possible service requests, S the set of service advertisements and W the set of values denoting the degree of relevance (for r) or degree of match (for e) between a request from Q and a service from S . Both r and e may assume various types of values: Boolean ($W=\{0,1\}$), real numbers ($W=[0,1]$), fuzzy terms ($W=\{\text{"irrelevant"}, \text{"relevant"}, \dots\}$), etc. Given these informal definitions, the

evaluation of a matchmaking engine is the determination of how closely vector e (delivered by the engine) approximates vector r (specified by domain experts).

A Boolean evaluation scheme is the traditional scheme used in the relevant literature. In this case, standard measures such as precision and recall are used for measuring the system performance. However it has some considerable pitfalls:

- the graded results of matchmaking algorithm execution are transformed to Boolean values, and, thus, the matchmaking and service semantics is ignored,
- such transformation involves the definition of a threshold. The assignment of an optimal value to this evaluation parameter is not a trivial task and there are no formal and commonly agreed ways it can be done, and
- the Boolean relevance assessments are too coarse-grained and do not always reflect the real intention of the domain expert

Given these shortcomings, the Boolean evaluation scheme cannot accurately assess how close the discovered services are to the actual relevant services. A solution to these problems would be to use an evaluation scheme based on graded relevance. However, this implies that apart from the existence of graded relevance judgments, appropriate metrics are in place. In the following sections we review some metrics that have been proposed in the literature and propose necessary adjustments.

2.2 Related Work

Most of the initial approaches rely on the Boolean evaluation scheme. To our knowledge, the first attempt to apply the concept of graded relevance to the evaluation of service discovery is described in [8]. This work was followed by [9] that proposed new metrics and methods for evaluating SWSD. However, even if the authors exploit graded relevance and reach some very interesting conclusions, they overlook the fact that the rankings returned by the SWSD tools are partial. Moreover, the fact that their evaluation relies on manually created relevance judgments constitutes a limiting factor as was the case for [8].

Besides the standard retrieval evaluation metrics (Precision, Recall, F-measure etc.) other metrics have been proposed in the literature. A very popular metric is Average Precision (AveP) over all relevant items for a query/request. Other approaches for catering for incomplete relevance judgments are RankEff, Ap_all, InducedAP, SubCollectionAP and InferredAP. In [10] Kekäläinen and Järvelin introduced the concept of gain. Every level of the relevance scale holds a gain value which indicates the gain that the user receives from finding a service belonging to this relevance level. They also proposed Cumulated Gain (cg) at rank r which depicts the total gain that the user receives by exploring the resulted ranking till rank r . However, Cumulated Gain does not penalize late retrieval of relevant services. Hence, the authors also added a discount factor in the calculation which decreases the gain of services as the rank increases, resulting in Discounted Cumulated Gain (dcg). To be able to compare various DCG curves from different discovery engines they proposed the use of normalised dcg (nDCG). nDCG is computed by dividing the DCG value of the result ranking with the dcg of an optimal ranking, called idcg (ideal dcg).

Based on the concept of gain, various metrics have been proposed, which can be expressed in terms of Cumulated Gain, like Q-Measure [11] and Average Weighted

Precision (AWP) [9]. Sakai in [11] proposed a new metric for graded relevance called Q-Measure which integrates AveP and AWP. Küster and König-Ries in [9] proposed Average Weighted Discounted Precision (AWDP) to be used for the evaluation of SWS discovery. However, all the metrics described so far are capable of evaluating only full rankings of services, or items in general, i.e. rankings without ties. McSherry and Najork [12] proposed a method to extend metrics for full rankings to partial rankings. In their work they extended Precision, Recall, F1-measure, Reciprocal Rank, Average Precision and nDCG (we will call nDCGp' its version applied to condensed lists). From all of them only the last one supports graded relevance.

Finally, a basic issue in evaluation is the creation of a test collection with relevance judgments. This task is performed manually, as some experts must judge every item in the collection with respect to every query. The solution that is proposed in [13] is the use of pooling. This method of pooling can be used when the set of items is a finite set that does not change frequently. But this is not the case for dynamic environments like the World Wide Web (WWW), where both the set of items and queries are under frequent change and it would be useful to evaluate systems without the need of manually created relevance judgments. In that direction, various techniques have been proposed, that can be classified in two large categories, those who automatically create relevance judgments from the rankings returned [13] and those who evaluate systems without relevance judgments.

2.3 Proposed Metrics for evaluation of SWS discovery

In this section we investigate and propose some evaluation metrics that demonstrate the desired characteristics for SWSD evaluation, as reported in the previous section. The first proposal is based on generalized versions of Precision and Recall. The other metrics are adaptations of metrics already proposed by other researchers.

2.3.1 Generalized Metrics

In order to deal with the problems identified above, we can assume that a service discovery system is a generalized retrieval system. In [8] we proposed such an evaluation scheme and performed a preliminary evaluation. The main idea is that the domain experts assess the relevance of specific services against a given request through fuzzy linguistic terms. In order to be able to compare the degrees of match used by the engines with the corresponding expert relevance assessments we need to express them in a similar form. Moreover, new metrics are required in this case. The proposed measures are generalizations of the recall and precision measures, calculated from the two rankings of relevance assessments: those delivered by the engine and those performed by domain experts in a way similar to the Boolean case.

2.3.2 Metrics for Evaluating Partial Rankings

From the literature survey it became apparent that apart from the nDCGp' metric [12], Q-Measure and AWDP would also be capable of evaluating SWSD systems, if they supported partial rankings. Hence, we provided some extensions of these metrics

towards this direction. The first extension was the definition of AWDP for partial rankings. The second one was the Q-measure for partial rankings.

2.4 Automatic Generation of Relevance Judgments

As already mentioned, it is important to be able to assess the effectiveness of SWSD systems without manually created relevance judgments, since it is very difficult to obtain them from domain experts. In this section we study two social voting methods for creating pseudo-relevance judgments. The first of these methods is based on the Borda Count and the second is the Condorcet method. Condorcet method seems to have better properties for the task at hand. A brief description of the method follows. According to the Condorcet method, the voters (matchmaking engines) rank the candidates (services) in an order of preference. However, ties and incomplete rankings are allowed. The final ranking is calculated based on the pair-wise number of wins of each candidate (i.e., service).

Firstly, the method builds a comparison matrix CM of size $n \times n$, where n is the cardinality of the set of all services returned by the engines. The $CM[i,j]$ element denotes the “wins” of the service i over the service j . When a service has not been returned by an engine, then it is assumed that has been defeated by all other ranked services returned by the engine. Next, for each pair of services we find the final ranking. We compare pair-wise all services and we add one point for each pair-wise win, lose and tie in the respective columns. The rules for ordering the services are:

- Winner is the service with the most wins.
- If two services have the same number of wins, then winner is the service with the less defeats.
- If the number of wins equals the number of defeats, then the services are ranked equally.

Our final goal, however, is to align this ranking to the relevance scale we use, so that the final relevance judgments are generated. In order to do this we use the following formula:

$$f_i = \frac{m - rank_i}{m}$$

where m is the number of discrete relevance levels in the final ranking, and $rank_i$ is the position of the i^{th} element in the ranking.

2.5 Experimental Evaluation

In order to assess the applicability of the metrics we performed several experiments under variable settings. The main objective of the experiments was to explore the behavior of the metrics and techniques used.

Setup

In the evaluation experiment we involved the following matchmaking engines [6]:

- OWLS – MX, which provides five different matchmaking algorithms.

- OWLS-SLR, which performs logic-based DL reasoning and also implements two distance metrics.
- TUB OWL-S Matcher, which performs matchmaking based on DL subsumption over many service parameters.

Something very interesting is the type of ranking returned by each engine. Table 1 shows a categorization of the aforementioned engines according to the type of ranking they return (as assumed for the experiments).

Table 1. SWSD engine categorization based on the type of ranking returned

SWSD Engine	Ranking
OWLS-MX (M0)	Partial
OWLS-MX (M1-M4, M1mx2-M4mx2)	Full/Partial
OWLS-SLR (UC, ED)	Full/Partial
TUB	Partial

In order to measure the effectiveness of the engines that returned full rankings we used the nDCG', Q'-measure and AWDP' metrics (quotes mean that they are applied to condensed lists). We also used these latter metrics in order to measure the effectiveness of engines returning partial rankings so as to assess the error introduced by their misuse. In order to present aggregate results for the generalised precision we defined the Average Generalised Precision (Average P_G or APG) which is the sum of all P_G values for all relevant services divided by the number of relevant services.

For the experiments we relied on the TC3 service collection. The main characteristics of this collection are that it includes a relatively large number of services and service requests and the corresponding relevance judgments as decided by "human experts". Apart from TC3, we also used the automatically-generated relevance judgments generated with the Condorcet method. Both TC3 and pseudo-relevance judgments had the same relevance scale shown in Table 2.

Table 2. Relevance scale for TC3 and pseudo-relevance judgments

Rank	Range of rank	Gain
Highly Relevant	(0.75, 1]	3
Relevant	[0.5, 0.75)	2
Potentially Relevant	[0.25, 0.5)	1
Irrelevant	[0, 0.25)	0

Experiments and Results

Firstly, we study the correlation between the automatically generated relevance judgments (pseudo-relevance judgments) and the TC3 judgments.

The pseudo-relevance judgments were produced based on the Condorcet method. The Kendall's tau-b correlation between the results was computed between the condensed lists of pseudo-relevance judgments and the TC3 mappings. We used the Kendall's tau-b correlation coefficient because it takes into consideration ties among

the ranked items. For each query, only the services that had been discovered by all engines were used for calculating the correlations (even if they had been judged as totally irrelevant).

Finally, in order to check the “rationality” of pseudo-relevance judgments, we manually created the following variants of the TC3 relevance judgments:

- *TC3+30* : the value of 30% of the relevance judgments (selected randomly) was upgraded to the next higher level in the relevance scale (e.g., from “Potentially Relevant” to “Relevant”),
- *TC3-30* : the value of 30% of the relevance judgments (selected randomly) was downgraded to the next lower level in the relevance scale,
- *TC3r30* : the value of 30% of the relevance judgments for each request were modified randomly, either upgraded or downgraded by one level.

Each variant is supposed to “represent” a different expert behavior.

The results are shown in Fig. 1 (one can observe that the Condorcet judgments have the highest correlation with the original TC3 mappings).

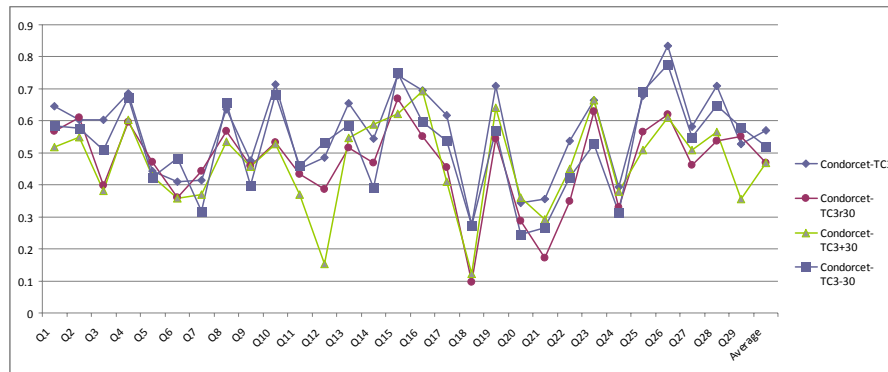


Fig. 1. Kendall's tau-b correlation: Condorcet and variants of TC3

Next, we show how much the SWSD evaluation depends on the selection of the correct metrics. Some aggregate results for better comparison are presented in Fig. 2 (for partial rankings/metrics and TC3 relevance judgments) and in Fig. 3 (same as previous but using the Condorcet relevance judgments). In all these figures, crossed lines denote differences between the rankings of the engines by the metrics. The correlation between the results of the various metrics is presented in Table 3.

Some observations that can be extracted from these figures are:

1. All metrics, except for the (average) $nDCG_p$, agree on the relative order of the engines' effectiveness (Fig. 2). $nDCG_p$ does not take into account the total number of relevant services (for which there exist relevance judgments), hence engines like the M2mx2 that return very few, but relevant, services per request, are ranked at a very high position. Hence, we can safely conclude that this metric is not very appropriate for the domain of SWSD evaluation.

- The AGP values are quite similar to AWDP(V) and Q-Measure(V) values, despite the fact that relies on a completely different approach, i.e., it is not based on the concept of *gain* (Fig. 2 and Table 3).

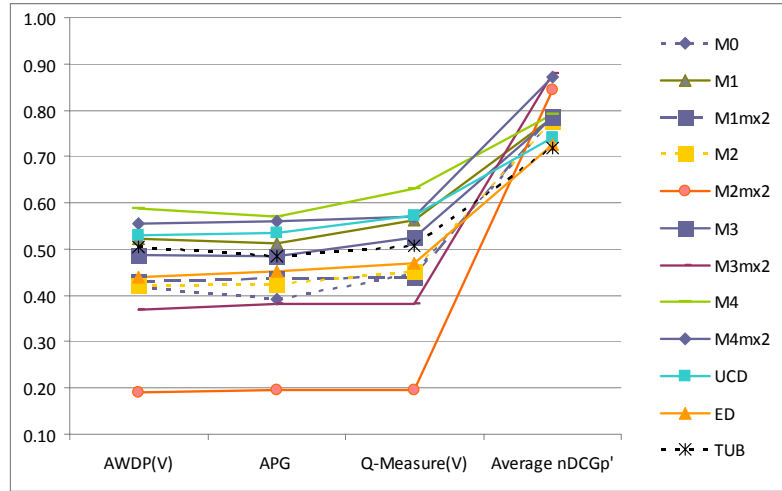


Fig. 2. Comparison of metrics and engines for partial service rankings (TC3)

Table 3. Pair-wise correlation of metrics (left:TC3, right: Condorcet)

		Kendall's tau			Kendall's tau
AWDP(V)	APG	0.96	AWDP(V)	APG	0.73
AWDP(V)	Q-Measure(V)	0.89	AWDP(V)	Q-Measure(V)	0.90
AWDP(V)	nDCGp'	-0.18	AWDP(V)	nDCGp'	0.29
APG	Q-Measure(V)	0.92	APG	Q-Measure(V)	0.63
APG	nDCGp'	-0.14	APG	nDCGp'	0.00
Q-Measure(V)	nDCGp'	-0.23	Q-Measure(V)	nDCGp'	0.40

The measurements show that the SWSD engines that participated in the creation of the pseudo-relevance judgments are favored in the experiments that use such judgments. However, we can always decide which are the best and worst services for a specific request. An observation is that the nDCGp' metric has the worst performance for the pseudo-judgments while the APG the best one.

As shown in Table 3, the APG metric has a high correlation coefficient with the other two “reliable” metrics, AWDP(V) and Q-measure(V). This fact, in combination with the observation of the previous subsection that it is less affected by the use of pseudo-judgments, constitutes it a rather promising metric. Another observation is that the absolute APG values are not affected significantly by the set of relevance judgments used (Fig. 2 and 3). Engines like M0 and TUB are slightly rewarded (since they participated in the creation of pseudo-judgments), but in general most of the engines have similar value ranges in the two experiments. In contrast, the AWDP(V) and Q-measure(V) metrics increase significantly the effectiveness values of the engines in Fig. 3. This is another reason why the generalized metrics seem more appropriate when the pseudo-judgments are used.

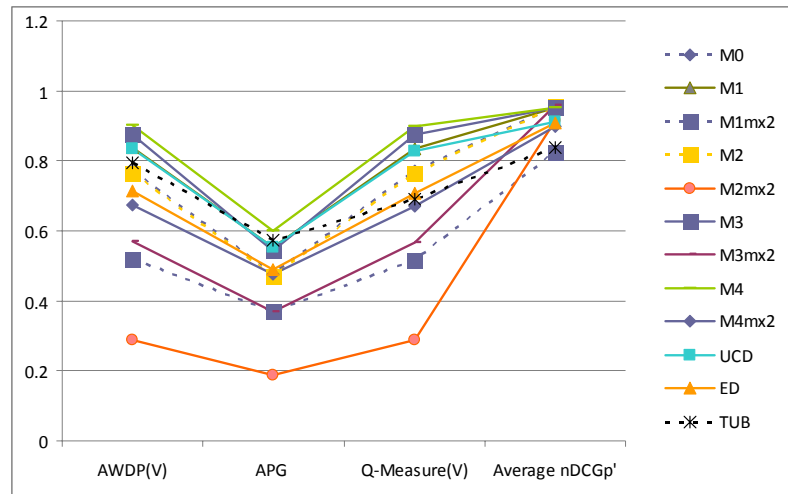


Fig. 3. Comparison of metrics and engines for partial service rankings (Condorcet)

Finally, as shown in the Fig. 2 the most effective service matchmaking engine is M4 (all metrics agree on that). Quite close performance have the UCD and M1 engines. In any case, M2mx2 engine has the worst service discovery performance. These results agree with the findings of the creators of the respective engines, as stated in the respective papers.

3 Conclusions

In the present thesis, several aspects of using semantic information management in building advanced applications and systems were studied. Some general conclusions follow:

Semantic Web Services

Semantic Web Services (SWS) are among the most usable and interesting results in the research field of the Semantic Web. Several researchers have dealt with SWS discovery and composition. However, there is a considerable lack of methods and tools for evaluating their efforts. In this work we proposed some methods in this respect, and we assessed their applicability experimentally. More research is necessary for establishing standard means of evaluation in this domain.

Collaborative context awareness

Context and situation awareness is an important aspect of modern and future computing systems. In this work we tried to increase the overall level of context awareness in a distributed environment, through a collaborative asynchronous scheme. The results were quite promising, since the performance of the system improved without decreasing the capability of the nodes to detect changes in their environment.

Semantic personalization

One of the most widely known forms of application intelligence is personalization. We proposed a reference framework that can facilitate the semantic description of entities and adaptation rules. The framework was validated in two application domains. In conclusion, semantic information management, where possible, can solve significant design problems of personalized systems.

References

1. Tsetsos, V. and Hadjiefthymiades, S. "An Innovative Architecture for Context Foraging", Eighth International ACM Workshop on Data Engineering for Wireless and Mobile Access (MobiDE, in conjunction with SIGMOD/PODS 2009), Providence, Rhode Island, (2009)
2. Tsetsos, V., Anagnostopoulos, C., Kikiras, P., and Hadjiefthymiades, S., "Semantically enriched navigation for indoor environments," *International Journal of Web and Grid Services*, vol. 2, no. 4, Inderscience Publishers, pp. 473--478, (2006)
3. Papataxiarhis, V., Riga, V., Nomikos, V., Sekkas, O., Kolomvatsos, K. Tsetsos, V. Papa-georgas, P. Xouris, V., Vourakis, S., Hadjiefthymiades, S., and Kouroupetroglou, G., "MNISIKLIS: Indoor LBS for All", *5th International Symposium on LBS & TeleCartography (LBS 2008)*, Salzburg, Austria, November, (2008).
4. Tsetsos, V., Papadimitriou, A., Anagnostopoulos, C. and Hadjiefthymiades, S. "Integrating Interactive TV Services and the Web through Semantics", *to appear in International Journal On Semantic Web and Information Systems, SI on "Semantic Media Adaptation & Personalization"*, IGI-Global, (2010)
5. Fernandez, B., Pazos Arias, Y., Lopez Nores, J.J., Gil Solla, A., & Ramos Cabrer, M. AVATAR: An Improved Solution for Personalized TV based on Semantic Inference, *IEEE Transactions on Consumer Electronics*. 52(1), pp. 223—231, (2006).
6. Tsetsos V., Anagnostopoulos C., and Hadjiefthymiades S., "Semantic Web Service Discovery: Methods, Algorithms and Tools", chapter in *"Semantic Web Services: Theory, Tools and Applications"* (Ed. Dr. Jorge Cardoso), IDEA Group Publishing, (2007)
7. Fan, J. and Kambhampati S. A Snapshot of Public Web Services, *SIGMOD Record*, 34(1), pp. 24--32, (2005).
8. Tsetsos, V., Anagnostopoulos, C., Hadjiefthymiades, S. "On the evaluation of Semantic Web Service matchmaking systems", *4th IEEE European Conference on Web Services (ECOWS)*, Zurich, Switzerland, (2006)
9. Küster, U., and König-Ries, B. "Evaluating Semantic Web Service Matchmaking Effectiveness Based on Graded Relevance," in *ISWC '08*, Karlsruhe, Germany, (2008).
10. Kekäläinen, K., and Jaana, J. "IR evaluation methods for highly relevant documents," in *SIGIR '00*, pp. 41—48, (2000)
11. Tetsuya Sakai, "New performance metrics based on multigrade relevance: Their application to question answering," in *NTCIR '04*, Tokyo, Japan, (2004).
12. McSherry, F., and Najork, M. "Computing Information Retrieval Performance Measures Efficiently in the Presence of Tied Scores," in *Lecture Notes in Computer Science*. Berlin Heidelberg: Springer-Verlag, (2008).
13. Nicholas C., & Cahan P. Soboroff I., Ranking retrieval systems without relevance judgments.: In *Proceedings of the 24th ACM SIGIR conference*, (2001).

Efficient algorithms for topology control and information dissemination/ retrieval in large scale Wireless Sensor Networks

Leonidas M. Tzevelekas¹

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
ltzev@di.uoa.gr

Abstract. Wireless Sensor Networks (WSNs) require radically new approaches for protocol/ algorithm design, with a focus towards energy efficiency at the node level. We propose two algorithms for energy-efficient, distributed clustering called Directed Budget Based (DBB) and Directed Budget Based with Random Delays (DBB-RD). Both algorithms improve clustering performance and overall network decomposition time when compared with state-of-the-art distributed clustering algorithms. Energy-efficient information dissemination in WSNs is proposed through modifying the movement of a simple Random Walk agent in a large scale Geometric Random Graph. The Random Walk with Jumps agent is compared against the Random Walk without backtracking agent. It is shown in simulations that the RW-J agent performs better than the RW agent in terms of Cover Time or Partial Cover Time. RW-J performs best when the underlying network topology has low connectivity, i.e. the graph has bad cuts. Information extraction from a large sensor field is a dual problem to information dissemination. We adopt the single mobile sink based information extraction methodology for collecting large amounts of data from a sensor field. The algorithms proposed are classified based on a. whether sensor nodes are allowed to transmit data over the wireless medium (single hop data forwarding) or not (almost zero hop data forwarding) and b. if there is detailed knowledge of sensor nodes locations in the field (deterministic algorithm) or not (randomized algorithm). All collection schemes are compared against required number of stops/ steps to completion and total physical distance covered by the mobile sink.

Keywords: Wireless Sensor Networks, distributed clustering, random walks, information dissemination, information harvesting from sensors, mobile sink based data harvesting

¹ Dissertation Advisor: Ioannis Stavrakakis, Professor

1 Dissertation Summary

Recent technological and scientific advances in the areas of solid state physics, integrated circuit design and telecommunications have led to enabling the functional design of innovative wireless networking models. These models comprise a large number of tiny sensor nodes working cooperatively towards building a wireless communication network. This doctoral dissertation focuses on the study/ proposal of energy-efficient algorithms for topology control and dissemination/ harvesting of information in large-scale Wireless Sensor Networks (WSNs). The doctoral dissertation comprises five chapters.

The first chapter of this dissertation is a detailed introduction in the subject of Wireless Sensor Networks. Specific details of recent technological advances towards making it technically feasible to produce large numbers of tiny sensor nodes are explained. A large ensemble of sensor nodes, which are embedded in the physical space, can produce a “smart environment”. The technical characteristics of “smart environments” are listed and the major models for Wireless Sensor Networks operation (as they have been proposed in literature) are explained. There is also a detailed report regarding the new technical challenges bound to be tackled by operation protocols designed for Wireless Sensor Networks. The most well known protocols in literature are explained at the end of the chapter (SPIN protocol, LEACH protocol, Directed Diffusion protocol, Publish/ Subscribe protocol).

Chapter 2 introduces two novel techniques for distributed clustering of sensor nodes in a large scale Wireless Sensor Network. The proposed Directed Budget Based (DBB) and Directed Budget Based with Random Delays (DBB-RD) algorithms have their basis on two previously published algorithms for distributed clustering of nodes in wireless networks, called Rapid and Persistent. The algorithms begin the distributed clustering process by distributing a set of coupons/ tokens offered to the initiator node evenly among the neighbors of that node; the process is then repeated in consecutive cycles of operation until the tokens are completely distributed or no more growth is possible in the network. The directed budget-based clustering algorithms called DBB and DBB-RD are proved through simulations to be energy-beneficial for the wireless sensor network due to both the reduced total number of exchanged messages in the network and the final cluster sizes achieved (close to the desired offered budget).

Another section of the dissertation describes an innovative technique for information diffusion in a large scale WSN. The described technique is based on random walks for information propagation inside the sensor network, which is modeled as a random geometric graph. The classic, well known, random walk involves the proliferation of the agent in the network by choosing uniformly at random among all next hop neighbors of the currently visited node. In contrast to this, the in-chapter-3-described technique involves the design of a “freezing” mechanism for the direction of movement of the random walk agent, such that the agent is allowed to be forwarded towards a specific direction in the network. The particular forwarding direction is retained by the random walk-with-jumps agent for as long as the agent will stay in the “freezing state”. It is shown through both simulations and

analysis that the incorporation of such a freezing mechanism into the otherwise pure random walk movement of the random walk agent will be beneficial for the overall covering process of the sensor network.

In Chapter 4, the random walk based movement is tested for contributing in the data harvesting process when the mobile sink based data harvesting is assumed. The movement of the mobile sink inside the sensor field is designed based on a) the wireless transmission of sensed data at 1-hop distance away from the producing node and b) the wireless transmission of sensed data at almost 0-hop distance away from the producing node. Furthermore, the movement of the mobile sink involves a deterministic variant (deterministically scheduling the changes in location of the mobile sink) and a random walk based variant (choosing randomly among all possible next mobile sink locations).

The final chapter of the dissertation summarizes contributions and results of previous chapters and furthermore touches on subjects of further work, which can be directly extracted from topics/ results presented in this dissertation.

2 Results and Discussion

2.1 Introduction

One of the main challenges associated with large-scale, unstructured and dynamic networking environments is that of *efficiently* reaching out to all or a portion of the network nodes (i.e. *disseminating information*), in order to provide, e.g., software updates or announcements of new services or queries. The high dynamicity and the sheer size of such networking topologies ask for the adoption of decentralized approaches to information dissemination [1], [2], [3], [4]. One of the simplest approaches employed for disseminating information in such environments, is the traditional flooding approach. Under flooding ([5], [6], [7], [8]), each time a node receives a message for the first time from some node, it forwards it to all its neighbors except from that node. Despite its simplicity and speed (typically achieving the shortest cover time possible, upper bounded by the network diameter), the associated large message overhead is a major drawback.

As flooding is considered not to be an option for large scale, wireless networking environments due to strict energy limitations of individual sensor nodes, approaches based on random walks are viewed as reasonable choices [9], [10], [11], [12]. Random walks possess several good characteristics such as simplicity, robustness against dynamic failures or changes to the network topology, and lack of need for knowledge of the network physical and topological characteristics. The Random Walk agent (RW agent) employed within a network of wireless sensors moves from neighbor node to neighbor node in a random manner, frequently revisiting previously covered nodes in a circular manner, even without backtracking (returning to the node it just came from is not allowed); these revisits constitute overhead and impact negatively on the cover time [8].

The Jumping Random Walk (J-RW) mechanism is proposed as an efficient alternative against the RW agent for information dissemination/ retrieval in large scale environments, like wireless sensor networks. The proposed scheme exploits the benefits of the RW mechanism (simple, decentralized, robust to topology changes) while providing a 'boost' in performance, i.e. accelerating the coverage process within the network. The latter is achieved by introducing a second state of operation to the RW agent in which the random movement paradigm is replaced by a non-random "directional" movement paradigm. It turns out that this improved significantly the cover time by "creating" virtual long links in topologies that lack them. It should be noted that the RW agent corresponds to a special case (parameter setting) of the proposed J-RW agent.

2.2 The RW agent in various topologies

A credible alternative to flooding for disseminating information in an unstructured environment, is the RW agent. In RW-based approaches, the initiator node employs an agent that will move randomly in the network, one hop/ node per time slot, informing (or querying) all the nodes in its path. Authors in [14] proposed a number of algorithms for RW-based searching in unstructured P2P networks, whereas probability-based information dissemination has been investigated for use in sensor networks [4], with data routing as the main consideration. Random walk in large-scale P2P nets has been shown to possess a number of good properties for searching and/or distributing of information within the network [15].

The overhead of RW-based solutions is considered to be much smaller than that of the flooding approaches, at the expense of a significant increase of cover time. *Cover time* is the expected time taken by a random walk to visit all nodes of a network. The generally relatively large (compared to flooding) cover time achieved under RW-based approaches depends on the network topology. For instance, it is $O(N \ln(N))$ for the fully connected graph (best-case scenario) and $O(N^3)$ for clique topologies (worst case scenario) [13], [14]. Random walks on random geometric graphs $G(N; r_c)$ have been shown to have optimal cover time $\Theta(N \ln(N))$ and optimal partial cover time $\Theta(N)$ with high probability given that the connectivity radius of each node r_c fullfills a certain threshold property, i.e. given that [16]

$$r_c^2 \geq \frac{c \ln(N)}{N}$$

For the rest of this work it is assumed that $G(V;E)$ is a connected network. Let N be the number of network nodes (equivalently, the size of set V). In such a network, a RW agent moving according to the previously described mechanism, will eventually visit or cover all network nodes after some time (cover time). Let $C_r(t)$ be the fraction of network nodes covered (or visited) after t time units or movements of the RW agent (i.e., the RW agent start moving at $t = 0$), for a particular realization (sample path) of the walk and for a given initiator node. $C_r(t)$ will be referred to hereafter as the coverage at time t . Clearly, $C_r(t)$ depends on the network size, the network topology, the initiator node and other factors. If T_r denotes the cover time then $C_r(T_r) = 1$; clearly $C_r(0) = 0$. As time increases, the RW agent is expected either to move to a

node that hasn't been covered previously (thus, $C_r(t)$ increases) or to move to an already covered node (thus, $C_r(t)$ remains the same). Therefore, $C_r(t)$ is a non-decreasing function ($C_r(t_1) \leq C_r(t_2)$, for $t_1 < t_2$).

The number of movements of RW-based solutions is much smaller than that under flooding approaches (where one movement of an agent corresponds to one message transmission), at the expense of a significant increase in cover time. For example, for the case of a fully connected network, a number of movements of the order of $N \ln(N)$, [17], is required under the RW mechanism, while under flooding approaches the corresponding number of movements (or messages) is of the order of N^2 . On the other hand, cover time under the RW mechanism is of the order of $N \ln(N)$, while under flooding it is upper bounded by the network diameter plus 1 (e.g., for the case of a complete graph cover time under flooding is 2).

2.3 The J-RW agent

2.3.1 Motivation

Fig. 1 illustrates a RW agent movement path initiated from the initiator node depicted inside the dotted ellipsis. The random walker spends some time revisiting nodes in the depicted "upper-left" network part, while nodes in other network parts are left unvisited. Suppose now that after a few time units long enough to "cover" a certain network part the RW agent moves to a "new" (most likely uncovered) network part ("bottom right" network part in Fig. 2). It is more likely than before to cover nodes that have not been visited previously by the agent, and therefore, accelerate the overall network cover process.

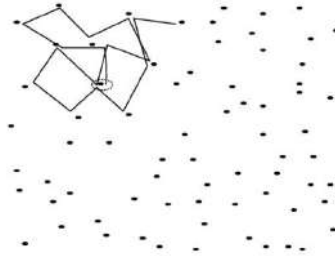


Fig. 1. RW Agent

One possible way for the agent to move away from a certain network region would be to carry out a number of consecutive directional movements, implementing a jump. This directional movement mechanism or jumping, initially proposed in [18], can be realized by switching occasionally away from the RW operation and engaging an operation implementing a directional movement. That is, such a RW agent (to be referred to as the Jumping Random Walk (J-RW)) operates under two states: State 0 under which it implements the typical RW mechanism without backtracking, and

State 1 during which the directional move is implemented; the time spent in state 1 (freezing state) will be referred to as the freezing (the direction) period.

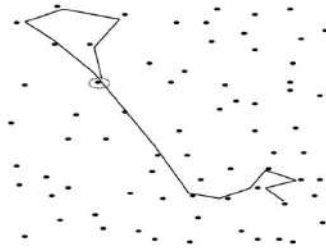


Fig. 2. J-RW Agent

The J-RW mechanism moves the agent - at the end of the freezing period - to networks that are expected (due to the directional freeze) to be geographically more distant than those reached by the RW agent after the same number of movements. That is, the introduction of the freezing state implements in essence jumps, defined as the physical distance between the nodes hosting the RW agent at the beginning and the end of the freezing period.

2.3.2 Description

The proposed J-RW mechanism is based on two underlying states. When in *State 0*, the J-RW agent operates as the already described RW agent. When in *State 1*, it implements a directional walk, by selecting as the next node to visit to be the neighbor of the current node that is the closest to the line connecting the current node and the node visited by the agent in the previous discrete time, in the direction away from the previously visited node. The directional walk may be easily implemented through a simple look up table involving the geographic locations of the neighbours of a node; this table determines the next node to forward the agent to under the directional walk, given that the agent came to this node from a given neighbour. The geographic information can be easily retrieved either at the time of deployment in the case of a static sensor field (with provisions for second, third, etc. choices when lower order choices are not available due to battery depletion), or after the deployment of the field with the help of a localization protocol run occasionally.

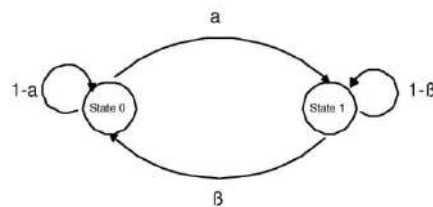


Fig. 3. Markov Chain mechanism for controlling J-RW agent movements

State transitions of the J-RW agent are assumed to occur at the discrete times according to a simple 2-state Markov chain, as shown in Fig. 3; let α (β) denote the transition probabilities from State 0 to State 1 (State 1 to State 0) and let $T_0 = 1/\alpha$ ($T_1 = 1/\beta$) denote the mean time (in discrete times of our reference time, or number of visits to nodes) that the agent spends in State 0 (State 1). Clearly, β (or, T_1) determines the length of time over which the directional walk is continuously in effect and, thus, the mean length of the induced jump. Similarly, α (or, T_0) determines the length of time over which the RW mechanism is continuously in effect. It should be noted that α and β should be carefully selected so that the mix of the two distinct operations is effectively balanced. It should be such that the implemented jump is sufficiently large to move the agent away from the current locality that is likely to be covered by the operation at State 0, and on the other hand, it should not be too large in which case it would leave uncovered large areas or require the random walk operation to operate long enough (at the increased cost of revisits) to cover the large areas between the start and the end of the jump. Similarly, α should be such that the time spent at State 0 be balanced so as to not over-cover or under-cover the current locality.

As previously for the RW mechanism, coverage and cover time under the J-RW mechanism may be defined in a similar manner, denoted by $C_j(t)$ and T_j respectively. $C_j(t)$ is a non-decreasing function of t taking values between 0 (for $t = 0$) and 1 (for $t > T_j$).

2.4 Coverage Analysis

2.4.1 Coverage under the RW and J-RW mechanisms

The main aim here is to derive an analytical expression for $C_r(t)$, which will serve as a tool for further understanding of random walk based information dissemination. Let's assume that the network topology is fully connected (i.e., all nodes are connected to all other nodes). This is actually the case for large values of r_c in geometric random graphs. For example, for nodes scattered in the $[0,1] \times [0,1]$ 2-dimensional plane, any value of $r_c > \sqrt{2}$ ensures that there is a link among any pair of nodes.

In such a network, each time the RW agent decides to move to a new neighbor node at time t (thus, arriving at time $t + 1$), coverage $C_r(t)$: (a) may increase ($C_r(t + 1) = C_r(t) + 1/N$), provided that the new node has not been covered previously; or (b) remain the same ($C_r(t + 1) = C_r(t)$), provided that the new node has already been covered. Note that at time t , in a fully connected network the RW agent may select one out of $N - 2$ network nodes (i.e., all network nodes excluding the one the agent came from and the one that is currently located at). Since $1/N$ corresponds to the coverage contribution of the node the agent came from and $1/N$ to the coverage contribution of the node that is currently located at, then $C_r(t) - 2/N$ is the coverage corresponding to the remaining $N - 2$ nodes and eventually, $(N - 2) \times (C_r(t) - 2/N)$

corresponds to the number of nodes that have already been visited by the agent (excluding the one the agent came from and the one that is currently located at). It is easily derived now that (on average) the increment of the coverage after a RW agent moves at time t , is given by the probability $1 - C_r(t)$ that it moves to a node not visited before multiplied by $1/N$ which is the contribution to coverage by each node that is visited for the first time. That is,

$$C_r(t+1) - C_r(t) = \frac{1}{N}(1 - C_r(t)) \quad (1)$$

Let t be continuous and let $\tilde{C}_r(t)$ denote the corresponding continuous and increasing function of $C_r(t)$. Based on previous equation we have

$$\frac{d\tilde{C}_r(t)}{dt} = \frac{1}{N}(1 - \tilde{C}_r(t)) \quad (2)$$

Equation (2) is a first class differential equation, and the solution is

$$\tilde{C}_r(t) = 1 - e^{-\frac{t}{N}} \quad (3)$$

Equation (3) was derived assuming a fully connected network. By reducing r_c in geometric random graphs, the number of neighbor nodes decreases and therefore a RW agent has fewer choices to move than before. Therefore, the fraction of nodes that have (not) been visited previously, is expected to deviate from $C_r(t)$. In order to account for the aforementioned decrease in the increase rate of $\tilde{C}_r(t)$ we introduce a positive constant k with $0 < k < 1$, such that

$$\tilde{C}_r(t) = 1 - e^{-\frac{k}{N}t} \quad (4)$$

The case of $k = 1$ corresponds to the fully connected network topology (i.e., large values of r_c) as it is concluded from Equation (3).

Coverage under the J-RW mechanism is related to r_c , α and β . However, an analytical expression for the coverage considering r_c , α and β is difficult to be derived and its further investigation will be based on simulations presented in the following section.

2.5 Simulation results and evaluation

There are multiple simulation runs executed under specific sets of parameters for the network and the investigated schemes. During each simulation run there is a large-scale node set up, with node population varying from 100 to 3000 nodes depending on the case. The nodes are placed at random locations on a square plain $[0,1] \times [0, 1]$. The random positions (x_u, y_u) of each node u in V are chosen within the set $[0,1]$ using the uniform probability distribution. Each node u is aware about its own position: (x_u, y_u) . Each node is connected to some other node if the euclidean distance among them is less or equal to r_c . Clearly, for $r_c > \sqrt{2}$ the resulting network is fully connected. Depending on N (the size of the network), the lower bound of r_c for which the topology remains connected varies (typically decreases as N increases). Four different values of r_c (0.05, 0.1 0.5 and 1.0) are considered in the sequel for those topologies of $N = 1000$. Note that all four values are less than $\sqrt{2}$.

2.5.1 The RW mechanism

Figure 5 presents simulation results for various topologies derived for $r_c = 0.05, 0.1, 0.5$ and 1.0 . The first observation is that for the appropriate value of k (i.e. $k = 0.3$ for $r_c = 0.05$, $k = 0.7$ for $r_c = 0.1$, $k = 0.9$ for $r_c = 0.5$ and $k = 1.0$ for $r_c = 1.0$), the analytical expression for coverage, given by Equation (3) approximates well the simulation results.

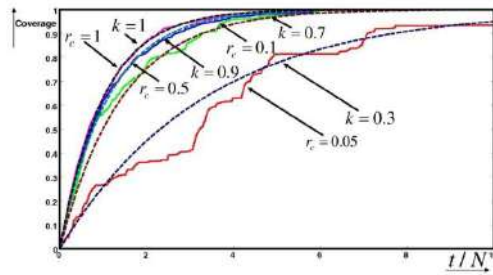


Fig. 4. Results for RW mechanism

2.5.2 The J-RW mechanism

Fig. 5-8 present simulation results under the J-RW mechanism for a network of 1000 nodes and various values of r_c and α . β has been kept constant and equal to 0.4, which means that as soon as State 1 is assumed (i.e., directional movement) the agent moves (on average) for 2-3 nodes towards a certain direction (more details are provided in the description of the J-RW mechanism in Section 3) before State 0 is assumed. In Fig. 5, coverage under the RW mechanism is clearly depicted and it is less than the coverage under J-RW for any value of α (e.g., 0.2, 0.4, 0.6 and 0.8). Note that the case depicted in Fig. 5 corresponds to a topology that is not highly connected ($r_c = 0.05$), thus even a relatively small value of $\beta = 0.4$ results in the J-RW doing significantly long jumps to get a performance improvement.

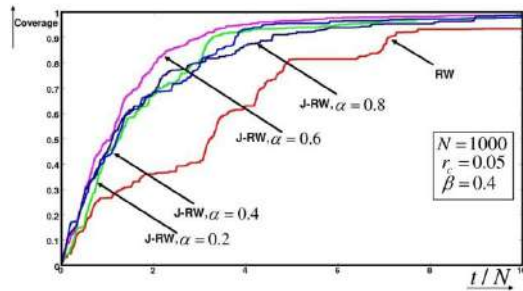


Fig. 5. Results for J-RW mechanism and low connectivity

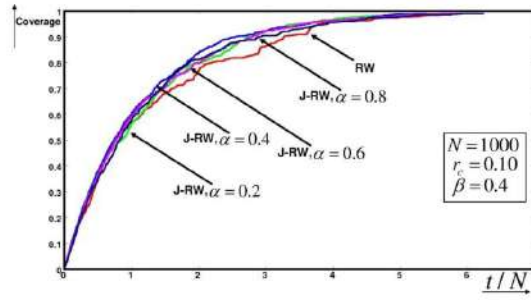


Fig. 6. Results for J-RW mechanism and medium connectivity

As the topology becomes more connected (r_c increases), the advantage of the J-RW mechanism is less obvious. For example, for the case depicted in Figure 6 ($r_c = 0.1$), coverage under the RW mechanism is still smaller than that under the J-RW mechanism (for any value of α), even though not that smaller as before, while for the case depicted in Figure 7 ($r_c = 0.5$), coverage under the RW mechanism is now larger than that under the J-RW mechanism (for any value of β). As r_c increases further, coverage under the RW mechanism is clearly higher than that under the J-RW mechanism for the specific combination of values of α and β . This is clearly depicted in Figure 8 for the case of $r_c = 1.0$ and can be attributed to the fact that the RW mechanism can now fully exploit the increased connectivity of the graph.

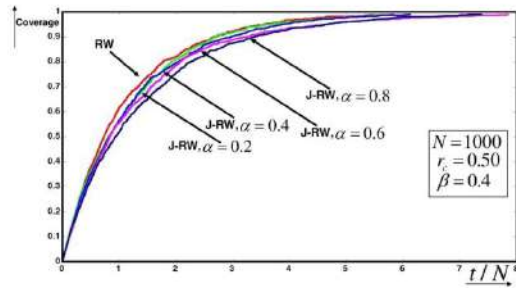


Fig. 7. Results for J-RW mechanism and higher connectivity

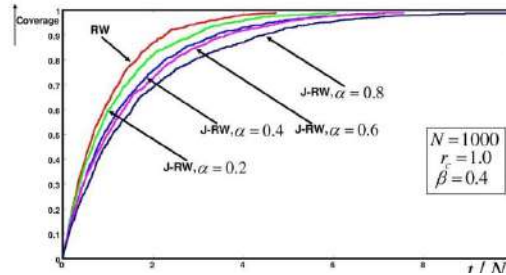


Fig. 8. Results for J-RW mechanism and highest connectivity

3 Conclusions

This dissertation presents results that are expected to foster technologies for next generation Wireless Sensor Networks. We develop energy efficient solutions for distributed clustering, for information dissemination and for information extraction in large scale Wireless Sensor Networks. We develop a methodology for distributed clustering in WSNs, called DBB, DBB-RD. We present a methodology for information dissemination that is energy efficient in terms of required number of steps of the RW agent to reach a given coverage level of nodes in the network. J-RW provides significantly higher coverage than RW because it is designed to avoid regions with bad cuts in the graph which trap the RW agent into revisiting already covered nodes. Finally we present a set of efficient algorithms for data collection in a Wireless Sensor Network when the single mobile-sink-based data harvesting methodology is adopted.

4 References

1. E. S. S. Dolev, J. Welch, Random walk for self-stabilizing group communication in ad hoc networks, IEEE Trans. on Mobile Computing 5 (7).
2. C. Avin, B. Krishnamachari, The power of choice in random walks: An empirical study, in: Proc. 9th ACM/IEEE International Symposium on Modeling, Analysis and Simulation of Wireless and Mobile Systems, (MSWiM), Malaga, Spain, 2006.
3. D. Braginsky, Rumor routing algorithm for sensor networks, in: Proceedings of the 1st ACM International Workshop on Wireless Sensor Networks and Applications, 2002, pp. 22-31.
4. U. Feige, A spectrum of time-space trade-offs for undirected s-t connectivity, Journal of Computer and System Sciences 54 (2) (1997) 305-316.
5. A. Segal, Distributed network protocols, IEEE Trans. on Information Theory IT (29) (1983) 23-35.

6. B. Williams, T. Camp, Comparison of broadcasting techniques for mobile ad hoc networks, in: MOBIHOC 2002, 2002, pp. 194-205.
7. K. Oikonomou, I. Stavrakakis, Performance analysis of probabilistic flooding using random graphs, in: The First International IEEE WoWMoM Workshop on Autonomic and Opportunistic Communications (AOC 2007), 2007.
8. D. Tsoumakos, N. Roussopoulos, Adaptive probabilistic search for peer-to-peer networks, in: 3rd IEEE International Conference on P2P Computing, 2003.
9. M. O.-K. M. Bani Yassein, S. Papanastasiou, On the performance of probabilistic flooding in mobile ad hoc networks, in: 11th International Conference on Parallel and Distributed Systems (ICPAD'05), 2005.
10. A. Abouzeid. Nabhendra Bisnik, Optimizing random walk search algorithms in p2p networks, *Comnet* 51 (6) (2006).
11. H. W. T. Lin, Search performance analysis in peer to peer networks, in: Third International Conference on Peer To Peer Computing (P2P'03), 2003, pp. 82-83
12. E. C. K. L. S. S. Q. Lv, P. Cao, Search and replication in unstructured peer-to-peer networks, in: ACM SIGMETRICS'02, 2002, pp. 258-259.
13. C. Avin, C. Brito, Efficient and robust query processing in dynamic environments using random walk techniques, in: Proc. IPSN 2004, Berkeley, California, 2004.
14. A. L.-O. R. Dorriv, P. Pralat, search algorithms for unstructured peer-to-peer networks, in: LCN 07: Proceedings of the 32nd IEEE Conference on Local Computer Networks, 2007, pp. 343-352.
15. M. Mihail. C. Gkantsidis, A. Saberi, Random walks in peer-to-peer networks, in: Proc. INFOCOM'04, Hong Kong, 2004.
16. C. Avin, G. Ercal, On the cover time of random geometric graphs, in: ICALP.(2005), pp. 677-689.
17. L. Lov'asz, Random walks on graphs: a survey, *Combinatorics*, Paul Erdos is Eighty, J. Bolyai Math. Soc., Vol. II 2 (1996) 353-397.
18. L. Tzevelekas, I. Stavrakakis, Improving Partial Cover of Random Walks in Large Scale Wireless Sensor Networks, 3rd IEEE WoWMoM Workshop on Autonomic and Opportunistic Communications (AOC'09), 15 June 2009, Kos, Greece.

Analysis and Retrieval of Mammographic Images

Stylianos Tzikopoulos*

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
stzikop@di.uoa.gr

Abstract. In this thesis two computer-aided diagnosis (CAD) systems are presented and implemented and their performance is evaluated. The first system proposed is a fully automated segmentation and classification scheme for mammograms based on breast density estimation and detection of asymmetry. First, image preprocessing and segmentation techniques are applied. Then, features for breast density categorization are extracted and Support Vector Machines (SVMs) are employed for classification, achieving accuracy of up to 85.7%. Most of these properties are used to extract a new set of statistical features for each breast, that are used to detect breast asymmetry between a pair of mammograms. The classifier adopted is an one-class SVM classifier, which resulted in a success rate of 84.47%. This composite methodology has been applied to the miniMIAS database. The results were evaluated by expert radiologists, are very promising and compared to other related works. The second system proposed is an experimental "morphological analysis" retrieval system for mammograms, using Relevance-Feedback techniques. The features adopted are first-order statistics of the Normalized Radial Distance, extracted from the annotated mass boundary. The system is evaluated on an extensive dataset of 2274 masses of the DDSM database, which involves 7 distinct classes. The experiments verify that the involvement of the radiologist as part of the retrieval process improves the results, reaching the precision rate of almost 90%. Therefore, Relevance-Feedback can be employed as a very useful complementary tool to a Computer Aided Diagnosis system.

Subject area: Digital Signal Processing, Medical Image Analysis, Pattern Recognition

Keywords: computer-aided diagnosis (CAD), content-based image retrieval (CBIR), relevance feedback (RF), mammography, classification

1 Introduction

Breast cancer, i.e., a malignant tumor developed from breast cells, is considered to be one of the major causes for the increase in mortality among women,

* Dissertation Advisor: Sergios Theodoridis, Professor

especially in developed countries. More specifically, breast cancer is the second most common type of cancer and the fifth most common cause of cancer death according to [1].

While mammography has been proven to be the most effective and reliable method for the early detection of breast cancer, as indicated by [2], the large number of mammograms, generated by population screening, must be interpreted and diagnosed by a relatively small number of radiologists. In addition, when observing a mammographic image, abnormalities are often embedded in and camouflaged by varying densities of breast tissue structures, resulting in high rates of missed breast cancer cases as mentioned by [3]. In order to reduce the increasing workload and improve the accuracy of interpreting mammograms, a variety of Computer-Aided Diagnosis (CAD) systems, that perform computerized mammographic analysis have been proposed, as stated by [4]. These systems are usually employed as a second reader, with the final decision regarding the presence of a cancer left to the radiologist. Thus, their role in modern medical practice is considered to be significant and important in the early detection of breast cancer.

The new CAD systems try to extract information not only from the external annotation given by the radiologist, but also from the images themselves. This way, they can provide to the user similar cases, by comparing the query image with the images of the available database. This methodology uses also the content-based image retrieval (CBIR systems). Over the recent years, these systems are gaining in importance [5, 6]. Such systems extract visual features from the "query" image, e.g. color, texture or shape and perform a comparison of it with the available images in a database, using specific similarity measures. The most similar images are returned to the user.

The scenario described above uses low-level features, which are not capable of capturing the image semantics, e.g. the high-level semantic concept that is meaningful for a user. This is known as the semantic gap. In order to address this gap, Relevance Feedback techniques (RF) have been developed since the early and mid-1990's [7]. In such a system, the user interacts with the search engine and marks the images that he perceives as relevant or non-relevant. Taking into account this feedback information, the engine "learns" and improved results are returned to the user during the next iteration.

Besides image retrieval, RF can also be employed to other systems, such as the retrieval of text documents, music or 3D objects. More recently it was used for medical image retrieval [8, 9]. In such a context, the aim of a retrieval system is to function in conjunction with a Computer Aided Diagnosis (CAD) system. The radiologists can be provided with relevant past cases -according to the query-, along with proven pathology and other information, making the diagnosis more reliable. RF seems an ideal scheme for the improvement of the performance of medical image retrieval systems, as it incorporates the radiologist's judgement, in order to capture some higher-level semantic concepts of the medical images. The judgement of such an expert is the result of a very complex and vague

procedure, combining a multitude of quantitative and qualitative facts, as well as the radiologist's experience, and therefore should be taken into consideration.

2 A fully automated scheme for mammographic segmentation and classification based on breast density and asymmetry

2.1 Overview of work

All of the CAD systems require, as a first stage, the segmentation of each mammogram into its representative anatomical regions, i.e., the breast border, the pectoral muscle and the nipple. The breast border extraction is a necessary and cumbersome step for typical CAD systems, as it must identify the breast region independently of the digitization system, the orientation of the breast in the image and the presence of noise. The pectoral muscle is a high-intensity, approximately triangular region across the upper posterior margin of the image, appearing in all the medio-lateral oblique (MLO) view mammograms. Automatic segmentation of the pectoral muscle can be useful in many ways [10]. One example is the reduction of the false positives in a mass detection procedure. In addition, the pectoral muscle must be excluded in an automated breast tissue density quantification method. The location of the nipple is also of great importance, as it is the only anatomical landmark of the breast, as mentioned by [11]. Most CAD systems use the nipple as a registration point for comparison, when trying to detect possible asymmetry between the two breasts of a patient. These automatic methods can also use the nipple as a starting point for cancer detection. Moreover, radiologists pay specific attention to the nipple, when examining a mammogram, according to [12].

Another important characteristic of a mammogram is the breast parenchymal density with regard to the prevalence of fibroglandular tissue in the breast as it appears on a mammogram. The relation between mammographic parenchymal density levels and high risk of breast cancer was first shown by [13], using four distinct classes for breast parenchymal density categorization. Thus, mammographic images with high breast density value should be examined more carefully by radiologists, for both physiological and imaging risk factors, creating a need for automatic breast parenchymal density estimation algorithms. In [14], such algorithms are presented and a new technique, introducing a histogram distance metric, achieves good results. Some existing algorithms, e.g., [15, 16], use the texture information of mammograms, in order to extract more features for breast density estimation.

Radiologists try also to detect possible asymmetry between the left and the right breast in a pair of mammograms, as it can provide clues about the presence of early signs of tumors such as parenchymal distortion. Many CAD systems analyze automatically the images of a mammogram pair and provide results for the detection of asymmetric abnormalities by applying some type of alignment and direct comparison, as implemented by [17]. In the works of [18, 19], directional

analysis methods are proposed, using Gabor wavelets, in order to detect possible asymmetry.

In this work, we propose a fully automated and complete segmentation methodology as the first stage of a multi-stage processing procedure for mammographic images [20–22]. Specifically, we have chosen to implement and apply the algorithm presented by [14] for breast boundary extraction, as the first step of the composite processing procedure; for the second step of pectoral muscle estimation, we enhanced the algorithm presented by [10] in order to achieve improved results; as a third step, we propose a new nipple detection technique, using the output of the breast boundary extraction procedure, when the nipple is in profile; that is, when it is projected on the background area of the mammogram, which is the recommended and usual case. The last algorithm, that is proposed in this work, besides locating the nipple point, can also serve as an improvement for the existing breast boundary algorithm, which misses the nipple if it is in profile. The improvement is obtained when updating the breast boundary, in order to include the detected nipple. Furthermore, as a fourth step, a new breast parenchymal density estimation algorithm is proposed, using segmentation of the inner-breast tissue, first-order statistics and fractal-based analysis of the mammographic image for the extraction of new statistical features, while the classification task is performed using Support Vector Machines (SVMs). Finally, a new algorithm is proposed for breast asymmetry detection, using the feature values already extracted from the breast parenchymal density estimation step, using an one-class SVM classifier. Both techniques achieve high success rates, often higher than the corresponding values of other algorithms in the relevant literature, while simpler and faster feature extraction methods have been employed. Our methodology has been tested on all the 322 mediolateral oblique view mammograms of the complete miniMIAS database, which is provided by [23], giving prominent results according to specific statistical measures and evaluation by expert radiologists, even in the case of such a difficult (very noisy) mammographic dataset.

2.2 Results and discussion

The complete system described was used for processing all the images of the miniMIAS database. All the intermediate results, i.e., breast boundary detection, pectoral muscle detection, nipple detection, asymmetry detection and breast density estimation, were examined in detail and evaluated by expert radiologists. It should be noted that the high level of noise, added to the images during the digitization process and the creation of the initial database images, makes the fully automated segmentation process a very challenging task.

The pre-processing techniques, which were selected to be applied in this work, were in general proved to be effective and successful, as the noise is correctly detected in most cases and sufficiently removed from the remaining stages of processing the images. The implemented breast boundary detection technique, which is based on a simple inference, gives satisfactory results. This is obvious

by a careful observation of the detected boundary of the images and also verified accordingly, as it is compared to the ground truth boundary using specific statistic measures, such as the Tanimoto Coefficient and the Dice Similarity Coefficient. The pectoral muscle estimate is accurate and further improved through the modification we propose, according to specific statistical measures extracted by the evaluation of the images from an expert radiologist. The new nipple detection technique tries to overcome the drawback of the breast boundary estimation method, i.e., not detecting the nipple, when this is in profile. In this way, it can serve as an improvement for the already established breast boundary, and in addition as a key point for further processing of the image, due to the importance of the nipple area in a mammographic image. Note that this technique can not be objectively compared to the algorithms proposed in previously published relevant literature, since the most similar one is the work by [12], which uses only a small subset of the miniMIAS database and has a different target than ours. The results were evaluated by expert radiologists and are promising enough to expect even better results, when applied to high quality digital mammograms.

The proposed algorithm for mammographic breast density estimation was tested on all the images of the miniMIAS database, fully annotated according to the 3 breast density classes. The results showed an accuracy of up to 85.71%, using the leave-one-out evaluation methodology. The achieved results are better compared to the relevant work of the bibliography [14, 16], although the latter one uses only a selected small portion of the miniMIAS database. The work of [15] achieves higher success rates, albeit it uses a different approach with higher-order textural features, which are computationally very expensive. The work we propose in this paper uses simple first-order statistical features and a new technique for the power spectrum estimation, making the whole process suitable for on-line training updates and real-time applications.

The asymmetry detection scheme was applied to all the images of the miniMIAS database, which is fully annotated, by characterizing each pair of mammograms as symmetric or asymmetric. The proposed methodology achieved a success rate of up to 84.47%. This success rate is similar to or even higher than the levels reported in the relevant literature, although it uses the complete set of images of the miniMIAS database, instead of a small subset, as the work of [18, 19]. Therefore, our experimental results can be considered more reliable and consistent. Furthermore, the use of the one-class classification algorithm turned out to be a simple yet effective way to overcome the problem of the imbalanced classes. The idea of the classification is to model as “target” the asymmetric cases and consider as “outliers” all the other cases, leading to an one-class scheme. The symmetric cases are not specifically modelled, but simply considered as non-asymmetric. In addition, note that our method is computationally much simpler and, more importantly, it is based on feature values that have already been computed and used. Thus, our method addresses the tasks of mammographic breast density estimation and asymmetry detection in an automatic, unified and generic way.

All the previously reported techniques can be combined and integrated to a clinical-level CAD system. All the algorithms are fully-automated and there is no need for external assistance. In addition, the processing time is not large enough, so each mammogram can be analyzed online; that is, on the fly as it is inserted the system. Moreover, the proposed scheme is considered to be robust against noise, as it has been verified by its application to the miniMIAS mammographic images database, in which the noise levels are very high and of varying nature.

3 Shape-based tumor retrieval in mammograms using relevance feedback techniques

3.1 Overview of work

As a second CAD system, a real application of a content-based medical image retrieval system is presented, while Relevance Feedback (RF) techniques are employed, in order to incorporate the radiologist to the retrieval process and further improve the results [24]. The system retrieves mammograms containing masses of the same morphology as the query image. The adopted features for the shape description are first-order statistics (mean value, standard deviation, mass circularity, entropy, area ratio parameter, zero-crossing count, roughness index) of the Normalized Radial Distance [25], extracted from the mass boundary. For the classification of the masses at step 0 of the RF procedure, a simple Euclidean minimum distance classifier [26] is used. On the next steps of the process, an SVM classifier is trained according to the feedback of the user. Note that we examine the performance of two different systems, in order to investigate the importance of the type of the patterns that the search engine returns to the user for labelling. In the simple SVM case [27], the system returns the most "confident" relevant patterns for labelling, e.g., the furthest patterns to the positive (relevant) side of the classifier. This can be easier for the user, but gives no useful information to the system, leading to slow convergence. However, in the active SVM case [28], the system returns the most "ambiguous" patterns for labelling, e.g. the patterns closest to the decision boundary, in order to improve the speed of the convergence. In order to evaluate the performance of the retrieval results at each round of the RF, the precision curve [29] is used. The precision at each round is defined as $pr = \frac{R}{N}$, where $N = 10$ is the total number of returned images to the user and R are the relevant images among them.

3.2 Results and discussion

For the evaluation of the Relevance Feedback scheme, a dataset of 2274 masses of the DDSM database [30], are used. Note that apart from the detailed boundary of each mass -used for the feature extraction step-, a classification of the shape of the masses in the following 7 distinct classes is also available: Irregular, Lobulated, Lymph Node, Oval, Round, Tubular or Other.

The experiments were carried out according to the following scenario:

- The user chooses a mass from the database as query image
- Repeat for steps 0 (no feedback yet) to 10 (user gave feedback 10 times)
 - The system returns to the user 10 images for evaluation and the precision is estimated
 - The system returns to the user 10 images to label
 - The user labels a subset of the images, as "relevant" or "non-relevant"
 - The system is re-trained, using the feedback of the user as new information

The above scenario is repeated for all the images of the database, in order to achieve more focused results. The system uses the simple SVM scheme [27], or the active SVM scheme [28]. In addition, the user is modeled as follows:

- The 'patient' user, that marks all the patterns returned by the system at each step as relevant or non-relevant, that can lead to imbalanced training sets.
- The less 'patient' user, that marks up to four relevant and four non-relevant patterns, among the patterns that the system returns at each step.
- The 'impatient' user, that marks up to three relevant and three non-relevant patterns, among the patterns that the system returns at each step.
- The 'lazy' user, that marks up to two relevant and two non-relevant patterns, among the patterns that the system returns at each step.

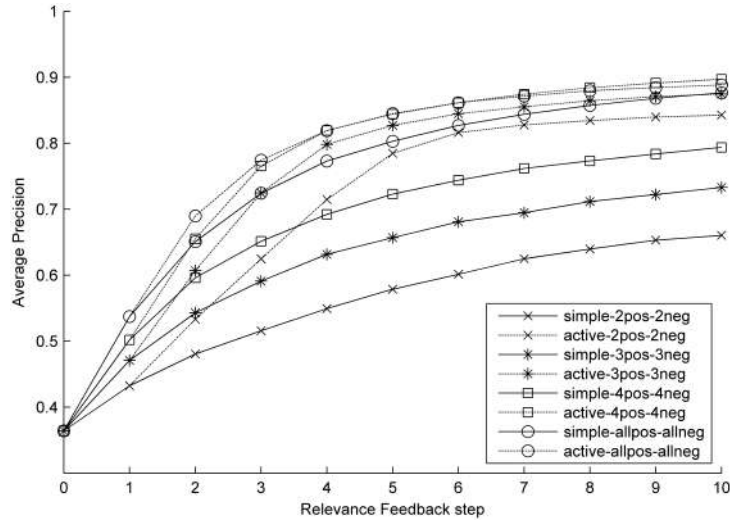


Fig. 1. Average precision at different steps of the RF procedure.

The average precision achieved at each iteration step for all the above configurations is shown in figure 1. Note that all the curves start from the same point

at step 0, as no information is given from the user. At step 1, the simple and active techniques of the same type of user achieve equal precision rate, as the available images at step 0 for each user type are the same for these two scenarios. However, at step 1 the user of the active scenario provides more informative feedback than the one of the simple scenario, leading to a quicker convergence of the classifier. This is the reason for the fact that active SVM outperforms the simple SVM at steps greater than or equal to 2, always for the same type of user. The maximum precision rate of 89.7% is observed for the case of active scenario that the user marks up to 4 relevant and 4 non-relevant patterns and not for the 'patient' user, because probably the latter one creates sometimes imbalanced training sets.

In this part of the thesis, Relevance Feedback has been employed as a complementary tool to a Computer Aided Diagnosis system, that retrieves masses with similar shape as the query one. The judgement of the radiologist is considered to be of high importance to such a sensitive system as a medical application, where the errors should be eliminated and therefore it is suggested to be taken into consideration. The results, which almost reach 90% precision rate, show that the retrieval process can be improved significantly, when the radiologist is incorporated in the retrieval process, even for a hard classification task of 7 classes, using features of first-order statistics.

The system converges much faster when the user is more actively involved in the process, by labeling more samples as "relevant" or "non-relevant". In addition, the active technique converges faster to better results than the simple one, while the average precision for each class follows the rules of the Relevance Feedback scheme. The mammographic dataset used for the evaluation is rather extensive, consisting of the large number of 2274 masses, categorized in 7 distinct classes; these facts ensure that the results presented are very useful, reliable and consistent.

The system is also available online for any user at [31].

4 Conclusions

The current thesis provides two CAD systems for the retrieval of mammographic images. The first one performs a -fully automated- segmentation of a mammogram. In addition, it estimated the breast density and detects possible asymmetry between a pair of images, corresponding to a pair of breasts. The system presented achieves high success rates, often higher than the works of the bibliography, although it uses simple -and inexpensive to compute- features. For this reason, the average processing time of an image is not large enough, so each mammogram can be analyzed online.

The other proposed CAD system examines the improvement achieved when RF techniques are used. More specifically, a "morphological analysis" retrieval system for mammograms is presented. It is evaluated on an extensive dataset of masses belonging to 7 distinct classes. The results, which almost reach 90% precision rate, prove that the retrieval process can be improved significantly,

when the radiologist is incorporated to the retrieval process, even for such a hard classification task.

References

1. R. Nishikawa, "Current status and future directions of computer-aided diagnosis in mammography," *Computerized Medical Imaging and Graphics*, vol. 31, no. 4-5, pp. 224-235, 2007.
2. M. Siddiqui, M. Anand, P. Mehrotra, R. Sarangi, and N. Mathur, "Biomonitoring of organochlorines in women with benign and malignant breast disease," *Environmental Research*, vol. 98, no. 2, pp. 250-257, 2005.
3. A. Wroblewska, P. Boninski, A. Przelaskowski, and M. Kazubek, "Segmentation and feature extraction for reliable classification of microcalcifications in digital mammograms," *Opto-Electronics Review*, vol. 11, no. 3, pp. 227-235, 2003.
4. R. Rangayyan, F. Ayres, and J. Leo Desautels, "A review of computer-aided diagnosis of breast cancer: toward the detection of subtle signs," *Journal of the Franklin Institute*, vol. 344, no. 3-4, pp. 312-348, 2007.
5. T. Gevers and A. Smeulders, "Content-based image retrieval: An overview." Prentice Hall, 2004.
6. R. Datta, J. Li, and J. Wang, "Content-based image retrieval: approaches and trends of the new age," in *Proceedings of the 7th ACM SIGMM international workshop on Multimedia information retrieval, November*. Citeseer, 2005, pp. 10-11.
7. X. Zhou and T. Huang, "Relevance feedback in image retrieval: A comprehensive review," *Multimedia systems*, vol. 8, no. 6, pp. 536-544, 2003.
8. W. Song, C. Huangshan, and T. Hua, "Analytic Implementation for Medical Image Retrieval Based On FCM Using Feature Fusion With Relevance Feedback," in *Proceedings of the 2nd International Conference on Bioinformatics and Biomedical Engineering, 2008. ICBBE 2008.*, 2008, pp. 2590 - 2595.
9. I. El-Naqa, Y. Yang, N. Galatsanos, R. Nishikawa, and M. Wernick, "A similarity learning approach to content-based image retrieval: application to digital mammography," *IEEE Transactions on Medical Imaging*, vol. 23, no. 10, pp. 1233-1244, 2004.
10. S. Kwok, R. Chandrasekhar, Y. Attikiouzel, and M. Rickard, "Automatic pectoral muscle segmentation on mediolateral oblique view mammograms," *Medical Imaging, IEEE Transactions on*, vol. 23, no. 9, pp. 1129-1140, 2004.
11. V. Andolina, S. Lill  , and K. Willison, *Mammographic imaging: a practical guide*. Lippincott Williams & Wilkins, 2001.
12. R. Chandrasekhar and Y. Attikiouzel, "A simple method for automatically locating the nipple on mammograms," *Medical Imaging, IEEE Transactions on*, vol. 16, no. 5, pp. 483-494, 1997.
13. J. Wolfe, "Risk for breast cancer development determined by mammographic parenchymal pattern." *Cancer*, vol. 37, no. 5, pp. 2486-92, 1976.
14. M. Masek, E. University of Western Australia School of Electrical, C. Engineering, and U. of Western Australia Centre for Intelligent Information Processing Systems, "Hierarchical Segmentation of Mammograms Based on Pixel Intensity," Ph.D. dissertation, 2004, <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.128.1245>.

15. A. Bosch, X. Munoz, A. Oliver, and J. Martí, "Modeling and Classifying Breast Tissue Density in Mammograms," in *Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition-Volume 2*. IEEE Computer Society Washington, DC, USA, 2006, pp. 1552–1558.
16. A. Oliver, J. Freixenet, A. Bosch, D. Raba, and R. Zwigelaar, "Automatic classification of breast tissue," in *Iberian Conference on Pattern Recognition and Image Analysis*. Springer, 2005, pp. 431–438.
17. F. Yin, M. Giger, K. Doi, C. Vyborny, and R. Schmidt, "Computerized detection of masses in digital mammograms: Automated alignment of breast images and its effect on bilateral-subtraction technique," *Medical Physics*, vol. 21, p. 445, 1994.
18. R. Ferrari, R. Rangayyan, J. Desautels, and A. Frere, "Analysis of asymmetry in mammograms via directional filtering with Gabor wavelets," *Medical Imaging, IEEE Transactions on*, vol. 20, no. 9, pp. 953–964, 2001.
19. R. Rangayyan, R. Ferrari, and A. Frere, "Analysis of bilateral asymmetry in mammograms using directional, morphological, and density features," *Journal of Electronic Imaging*, vol. 16, no. 1, pp. 13 003–13 003, 2007.
20. S. Tzikopoulos, H. Georgiou, M. Mavroforakis, N. Dimitropoulos, and S. Theodoridis, "A fully automated complete segmentation scheme for mammograms," in *Digital Signal Processing, 2009 16th International Conference on*. IEEE, 2009, pp. 1–6.
21. S. Tzikopoulos, H. Georgiou, M. Mavroforakis, and S. Theodoridis, "A fully automated scheme for breast density estimation and asymmetry detection of mammograms," in *17th European Signal Processing Conference (EUSIPCO)*, 2009.
22. S. Tzikopoulos, M. Mavroforakis, H. Georgiou, N. Dimitropoulos, and S. Theodoridis, "A fully automated scheme for mammographic segmentation and classification based on breast density and asymmetry," *Computer Methods and Programs in Biomedicine*, 2011.
23. J. Suckling, J. Parker, D. Dance, S. Astley, I. Hutt, C. Boggis, I. Ricketts, E. Stamatakis, N. Cerneaz, S. Kok *et al.*, "The Mammographic Image Analysis Society Digital Mammogram Database," in *Excerpta Medica. International Congress Series*, 1994, pp. 375–378.
24. S. Tzikopoulos, H. Georgiou, M. Mavroforakis, and S. Theodoridis, "Shape-Based Tumor Retrieval in Mammograms Using Relevance-Feedback Techniques," *Artificial Neural Networks-ICANN 2010*, pp. 251–260, 2010.
25. J. Kilday, F. Palmieri, and M. Fox, "Classifying mammographic lesions using computerized image analysis," *IEEE Transactions on Medical Imaging*, vol. 12, no. 4, pp. 664–669, 1993.
26. S. Theodoridis and K. Koutroumbas, *Pattern Recognition*, 4th ed. Academic Press, 2009.
27. H. Drucker, B. Shahrar, and D. Gibbon, "Relevance feedback using support vector machines," in *MACHINE LEARNING-INTERNATIONAL WORKSHOP THEN CONFERENCE*. Citeseer, 2001, pp. 122–129.
28. S. Tong and E. Chang, "Support vector machine active learning for image retrieval," in *Proceedings of the ninth ACM international conference on Multimedia*. ACM, 2001, pp. 107–118.
29. J. Luo and M. Nascimento, "Content-based sub-image retrieval using relevance feedback," p. 9, 2004.
30. M. Heath, K. Bowyer, D. Kopans, R. Moore, and P. Kegelmeyer, "The digital database for screening mammography," pp. 212–218, 2000.
31. Image Processing Techniques for Mammographic Images , <http://mammo.di.uoa.gr>, 2010.

Processing and Recognition of Handwritten Documents

Georgios Vamvakas *

¹ Department of Informatics and Telecommunications
National and Kapodistrian University of Athens

² Computational Intelligence Laboratory
Institute of Informatics and Telecommunications
National Centre for Scientific Research “Demokritos”

gbam@iit.demokritos.gr

Abstract. Nowadays, the accurate recognition of machine printed characters is considered largely a solved problem. A lot of commercial products are focused towards that direction, achieving high recognition rates. However, handwritten character recognition is comparatively difficult. So, the recognition of handwritten documents is still a subject of active research. In this thesis we studied the processing and focused on the recognition stages for handwritten optical character recognition. At the recognition stage a feature vector is extracted for all extracted characters in order to classify them to predefined classes using machine learning techniques. We studied several feature extraction techniques and developed methodologies that efficiently combine different types of features. Furthermore, a novel methodology that extracts features and classifies characters using a hierarchical scheme is proposed. This methodology, after being tested on well-known character databases, as well as on databases consisting of characters from historical documents and a database consisting of Greek contemporary handwritten characters, that were particularly created in this thesis, achieved recognition rates that are among the best one can find in the literature. This methodology was also applied to cursive handwritten words. The recognition rates in these experiments were also very high. Finally, an algorithm that automatically estimates the free parameters involved in character segmentation is also suggested. Character segmentation is very important because its result affects directly the recognition rates. Thus, the optimal segmentation is essential for a successful recognition.

Keywords: handwritten character recognition, feature extraction, hierarchical classification, machine learning techniques, character databases

1 Introduction

The large amount of documents that we have in our possession nowadays, due to the expansion of digital libraries, has pointed out the need for reliable and accurate Optical Character Recognition (OCR) systems for processing them. Although, the accurate recognition of contemporary machine printed characters is considered largely a solved problem, as mentioned above, handwritten character recognition is comparatively difficult, due to different handwriting styles, cursive handwriting and possible skew.

Another challenging task in OCR is the recognition of historical documents. Such documents are of great importance because they are a significant part of our cultural heritage. However, their low quality, the lack of standard alphabets and the presence of unknown fonts are major drawbacks in achieving high recognition

* Dissertation Advisors: ¹ Sergios Theodoridis, Professor – ² Dr. Basilis Gatos, Researcher

rates. In case of historical document processing in particular, an important area is word spotting. In many cases due to high levels of distortion, extremely poor quality, cursive handwriting etc. such documents can not be processed by an OCR system. In order to extract information from these documents page retrieval approaches, for searching or indexing, are adopted. However, this has to be done manually. In essence, this means that each occurrence of a word in a corpus must be annotated by hand. The goal of the word spotting idea is to greatly reduce the amount of annotation work that has to be performed.

According to the above, one can easily realize there are various issues that an OCR system have to deal with, such as character/word recognition and word spotting for either historical or contemporary documents. In this thesis, we suggest novel methodologies that attempt to deal with such issues, thus trying to assist document image processing. Moreover, we also propose an automatic unsupervised free parameter selection approach that optimizes the character segmentation algorithm adopted. This is essential because the segmentation step affects directly the recognition result.

2 Related Work

A widely used approach in OCR systems is to follow a two step schema: a) represent the image as a vector of features and b) classify the feature vector into classes. Selection of a feature extraction method is important in achieving high recognition performance. A feature extraction algorithm must be robust enough so that for a variety of instances of the same symbol, similar feature sets are generated, thereby making the subsequent classification task less difficult [2].

Feature extraction methods have been based mainly on three types of features [1, 2, 3 and 4]: a) statistical derived from statistical distribution of points b) structural and c) transformation-based or moment-based features. A survey on feature extraction methods can be found in [5]. Moreover, other approaches focus on measuring the similarity/dissimilarity between shapes by mapping one character onto another [6, 7].

All the above feature extraction techniques have been applied with great success to both historical and contemporary document recognition. However, there are also methodologies focused on the unique characteristics of the corresponding historical document they process, such as content and writing style [8, 9].

There have been quite a number of successes in determination of invariant features and a wide range of classification methods have been extensively researched. However, as mentioned in [10], most character recognition techniques use a “one model fits all” approach, i.e. a set of features and a classification method are developed and every test pattern is subjected to the same process regardless of the constraints present in the problem domain. It is shown that approaches which employ a hierarchical treatment of patterns can have considerable advantages compared to the “one model fits all” approaches, not only improving the recognition accuracy but also reducing the computational cost as well.

Most classification strategies in OCR deal with a large number of classes trying to find the best discrimination among them. However, such approaches are vulnerable to classification errors when patterns of similar shapes are present since they are not easily distinguished. In [11] a two-stage classification approach is presented to detect and solve possible conflicts between patterns with similar

shapes. During the first stage, a single classifier or ensemble of classifiers detect potential conflicts. The second processing stage becomes active only when a decision on the difficult cases must be taken. A comparative study between three different two-stage hierarchical learning architectures can be found in [12].

Word-spotting techniques for searching and indexing historical documents have been introduced. In [13], word images are grouped into clusters of similar words by using image matching to find similarity. Then, by annotating “interesting” clusters, an index that links words to the locations where they occur can be built automatically. In [14] and [15] holistic word recognition approaches for historical documents are presented. Their goal is to produce reasonable recognition accuracies which enable performing retrieval of handwritten pages from a user-supplied ASCII query.

3 Efficient Combination of Feature Extraction Techniques

In our approach [16], we employ four types of features. The first set of features is based on zones. The image is divided into horizontal and vertical zones, and for each zone we calculate the density of the character pixels (Fig. 1).

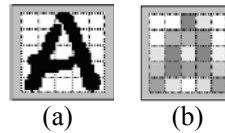


Figure 1. Feature extraction of a character image based on zones. (a) The normalized character image. (b) Features based on zones. Darker squares indicate higher density of character pixels.

In the second type of features, the area that is formed from the projections of the upper and lower as well as of the left and right character profiles is calculated. Firstly, the center mass (x_c, y_c) of the character image is found.

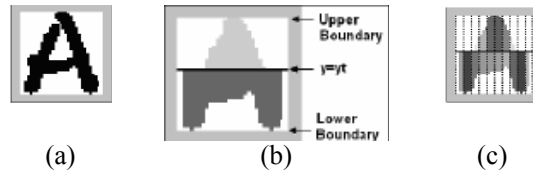


Figure 2. Feature extraction of a character image based on upper and lower character profile projections. (a) The normalized character image. (b) Upper and lower character profiles. (c) The extracted features. Darker squares indicate higher density of zone pixels.

Upper/lower profiles are computed by considering, for each image column, the distance between the horizontal line $y=y_c$ and the closest pixel to the upper/lower boundary of the character image (Fig. 2b). This ends up in two zones (upper, lower) depending on y_c . Then both zones are divided into vertical blocks. For all blocks formed we calculate the area of the upper/lower character profiles. Fig. 2c illustrates the features extracted from a character image using upper/lower character profiles. Similarly, we extract the features based on left/right character profiles.

The third feature set is based on the distances of the first image pixel detected from the upper and lower boundaries of the image, scanning along equally spaced vertical lines as well as from the left and right boundaries scanning along equally spaced horizontal lines (Fig. 3).

The forth set, calculates the profiles of the character from the upper, lower, left and right boundaries of the image, as shown in Fig. 4. The profile counts the number of pixels between the edges of the image and the contour of the character. These features are used because they describe well the external shape of the characters.

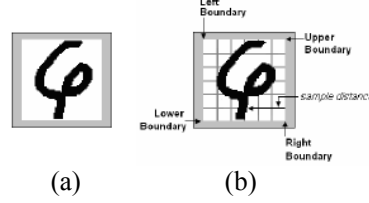


Figure 3. Feature extraction of a character image based on distances. (a) The normalized character image. (b) A sample distance from the right boundary.

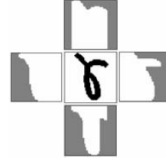


Figure 4. Features extraction of the character image based on distances.

Our methodology [16] for character recognition also considered a dimensionality reduction step, according to which the dimension of the feature space, engendered by the features extracted as described above, is lowered down to comprise only the features pertinent to the discrimination of characters into the given set of letters. In particular, we employed the Linear Discriminant Analysis (LDA) method, according to which the most significant linear features are those where the samples distribution has important overall variance while the samples per class distributions have small variance. Formally, this criterion is represented:

$$\text{LDA}(w) = \frac{w^T \text{Cov}(X) w}{w^T E_c [\text{Cov}(X | c)] w} \quad (1)$$

where w represents a linear combination of the original features, X the original feature vector, c the class, Cov is a the covariance matrix that has to be estimated from the samples and E_c is the expectation in respect to the classes. It turns out that finding the linear features that maximize the LDA criterion comes down to solving a generalized eigenvalue/eigenvector problem and keeping the eigenvectors that have greater eigenvalues. Moreover, the ratio of the sum of the eigenvalues kept to the overall eigenvalues sum provides as an index of quality of the feature subspace kept.

4 Hierarchical Character/Word Recognition

In this section a new feature extraction method followed by a hierarchical classification scheme is presented [17].

4.1 Feature Extraction

4.1.1 Characters

Let $im(x,y)$ be the character image array having 1s for foreground and 0s for background pixels and x_{max} and y_{max} be the width and the height of the character

image. Our feature extraction method relies on iterative subdivisions of the character image, so that the resulting sub-images at each iteration have balanced (approximately equal) numbers of foreground pixels, as far as this is possible. At the first iteration step (zero level of granularity, that is $L = 0$) the character image is subdivided into four rectangular sub-images using a vertical and a horizontal divider line as follows: Firstly, a vertical line is drawn that minimizes the absolute difference of the number of foreground pixels in the two sub-images to its left and to its right. Subsequently, a horizontal line is drawn that minimizes the absolute difference of the number of the foreground pixels in the two sub-images above and below. An important point is that the above dividing lines are determined taking into account sub-pixel accuracy. The pixel at the intersection of the two lines is referred to as the *division point* (DP). At further iteration steps (levels of granularity $L=1, 2, 3 \dots$), each sub-image obtained at the previous step is divided into four further sub-images using the same procedure as above (Fig.5).

Let L be the current level of granularity. At this level the number of the sub-images is $4^{(L+1)}$. For example, when $L = 0$ (Fig.5b) the number of sub-images is 4 and when $L = 1$ it is 16 (Fig.5c). The number of DPs at level L equals to 4^L . At level L , the co-ordinates (x_i, y_i) of all DPs are stored as features. So, for every L a $2 \cdot 4^L$ - dimensional feature vector is extracted.

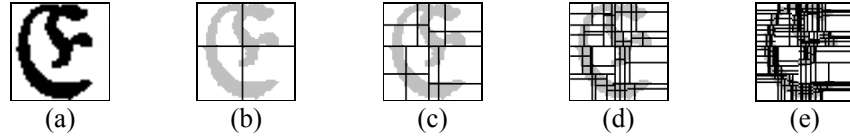


Figure 5. Character image and sub-images based on DP: (a) original image, (b), (c), (d), (e) subdivisions at levels 0, 1, 2 and 3 respectively.

After all feature vectors are extracted each feature is scaled to $[0, 1]$. Since each character is normalized to an $N \times N$ matrix all feature values f are in the range of $[1, N]$. Therefore, the value f_i of the i_{th} feature of every feature vector is normalized according to Eq.2.

$$f'_i = \frac{f_i}{N} \quad (2)$$

4.1.2 Words

In case of word recognition [18] the feature extraction technique applied for characters is also adopted. However, in order for features to be invariant of scaling the feature vector does not consist of the co-ordinates (x_i, y_i) of all DPs at a level L but of the pairs $(x_i - x_0, y_i - y_0)$, where x_i, y_i are the co-ordinates of the DP at L and x_0, y_0 are the co-ordinates of the initial DP (at level $L = 0$) of the word image. Furthermore, all features are normalized in the range of $[-1, 1]$. Since every word is normalized to an $N \times M$ matrix all feature values are scaled according to Eq. 3 and 4.

$$x'_i = \frac{x_i - x_0}{N} \quad (3)$$

$$y'_i = \frac{y_i - y_0}{M} \quad (4)$$

4.2 Hierarchical Classification

For the recognition procedure a hierarchical classification scheme is employed. Since characters/words with similar structure are often mutually confused when using a certain granularity feature representation, we propose to merge the corresponding classes at this level of classification. At a next step, we distinguish those character/word classes by employing a feature vector extracted at another level of granularity where the misclassifications between them are the least possible. The proposed classification scheme has a) a training and b) a recognition phase:

a. Training Phase

The training phase consists of three distinct steps: Step 1 is used to determine the level with the highest recognition rate for the initial classification, step 2 to merge mutually misclassified classes at the level found in step 1 and step 3 to find the level at which each group of merged classes is distinguished the best and to train a new classifier for each one at this level. These steps are described below:

Step 1: Starting from level 1 and gradually proceeding to higher levels of granularity, features are extracted, the confusion matrix is created and the overall recognition rate is calculated, until the recognition rate stops increasing. The level at which the highest recognition rate is achieved is considered to be the best performing granularity level (*BPGL*). Confusion matrices are created at each level from the training set using a *K*-fold cross-validation process. In our case *K* is set to 10.

Step 2: At *BPGL* where the maximum recognition rate is obtained the corresponding confusion matrix is scanned and classes with high misclassification rates are merged. Class merging is performed using the *disjoint grouping scheme* presented in [12]. Let the confusion matrix for *C* classes be $A_{i,j}$, where $A_{i,j}$ ($i, j = 1, 2 \dots C$) is the number of samples that belong to class *i* and are classified to class *j*. The similarity between classes *i* and *j* is defined according to Eq. 5.

$$N_{i,j} = A_{i,j} + A_{j,i}, (i < j) \quad (5)$$

Suppose we have two groups of classes G_p and G_q having *m* and *n* classes respectively. The similarity between these groups ($p < q$) is defined as:

$$S_{p,q} = \min_{i < j} N_{i,j}, (i = i_1 \dots i_m, j = j_1 \dots j_n) \quad (6)$$

Initially each class is a group. First two classes *i* and *j* with the highest $N_{i,j}$ value are found and merged into one group thus resulting in *C* – 1 groups. Next, the most similar groups according to Eq. 6 are merged into one. The procedure is iterated until all similarity values between groups are equal to zero in order to find all possible misclassifications.

Step 3: For each group of classes found in Step 2 the procedure described in Step 1 is performed again and the best distinguishing granularity level (*BDGL*) for its classes is found. Then, for every group another classifier is trained with features extracted at its *BDGL* in order to distinguish the merged classes at the next stage of the classification.

b. Recognition Phase

Each pattern of the test set is fed to the initial classifier with features extracted at *BPGL*. If the classifier decides that this pattern belongs to one of the non-group classes then its decision is taken into consideration and the unknown pattern is assumed to be classified. Else, if it is classified to one of the group classes then it is given to the group's corresponding classifier and this new classifier decides

about the recognition result. Note that if a sample is wrongly classified to a non-group class then at the next stage it will remain wrong. However, if it is misclassified to a group-class then it is possible to be correctly classified in the second stage.

5 Word Spotting

In this section a word spotting technique, based on the combination of results from different levels of the feature extraction method described in Section 4, is introduced [18]. Given a keyword that we want to match in a set of document images that have been segmented at word level, the matching algorithm is applied as follows:

Step 1: Create five lists R_i , $i = 1, 2, 3, 4$ and 5 , each one consisting of the *Euclidean Distances* between the keyword and every word of the set of documents, using feature vectors from granularity levels L_i , $i = 1, 2, 3, 4$ and 5 respectively, extracted according to the procedure described in Section 3.2.

Step 2: Normalize all distances in each R_i to $[0, 1]$ by dividing each one with the maximum distance in R_i .

Step 3: Merge all five lists in a list Q . Every word of the set of documents is represented in Q by five distances from the keyword. For each one we choose to keep the minimum distance and remove the others, resulting to Q' .

Step 4: Sort Q' in ascending order. Choose a threshold thr and keep only the first thr instances. List Q' now contains only the thr nearest words to the word we want to be matched.

6 Automatic Unsupervised Parameter Selection for Character Segmentation

Character segmentation is a difficult problem since low quality of document images and the wide variety of fonts can cause touching and broken characters. In most segmentation approaches a major problem is the selection of the free parameters that affect directly the segmentation results. The parameters are either user-specified and no training method is included [19, 20 and 21] or selected through a training procedure over a set of “optimal” parameter values that are usually manually selected based on some assumption regarding the training data [22]. However, ground truth or a priori knowledge of the fonts of the document image is not always available. To this end, we introduce a novel automatic unsupervised parameter selection methodology for character segmentation that is based on clustering [23]. The clustering is performed using features extracted from the segmented entities based on zones and from the area that is formed from the projections of the upper/lower and left/right profiles as described in Section 3. Optimization of an appropriate intra-class distance measure yields the optimal parameter vector.

Consider a character segmentation algorithm whose result depends on P parameters. Let $S_1, S_2 \dots S_v$ different parameter vectors (p-tuples) for different values of the parameters obtained using a standard selection method (e.g. random selection, selection through a grid). In our approach the well-known k-Means clustering algorithm is adopted due to its computational simplicity and the fact that, as all clustering techniques which use point representatives, is suitable for recovering compact clusters. If the expected number of different characters is in

the interval between k_1 and k_2 , then for every S_q we proceed to a k-Means clustering with k taking values from k_1 to k_2 .

Given a parameter vector S_q , in order to evaluate the performance of the clustering algorithm for every k between k_1 and k_2 , the mean squared distances from the centroids (within clusters sum of squares) is calculated as follows:

$$W_q(k) = \sum_{j=1,2,\dots,k} \frac{1}{n_{c_j}} \sum_{i \in C_j} d^2(x_i, \bar{x}_j) \quad (7)$$

where \bar{x}_j is the centroid of the cluster $C_j, j = 1, 2, \dots, k$, x_i is the i_{th} pattern inside cluster C_j , n_{c_j} is the cardinality of cluster C_j and d is the Euclidean Distance.

The value of $W_q(k)$ is low when the partition is good thus resulting to compact clusters. A measure of the quality of the segmentation result that corresponds to a parameter vector S_q is given as:

$$Q(S_q) = \frac{10^5}{\min_{k=k_1, \dots, k_2} (W_q(k))} \quad (8)$$

The optimal parameter vector S_{opt} is defined as:

$$S_{opt} = \arg \max_{S_q = S_1, S_2, \dots, S_v} (Q(S_q)) \quad (9)$$

7 Experimental Results

For our experiments the well-known CEDAR CD-ROM-1 [24], a database consisting of Greek handwritten characters (CIL-Database) [16], two databases [26] comprising samples of characters from old Greek Christian documents of the 17th century (HW and TW Databases) and a character database (TW-1 Database) [18] created by a part of a historical book from Eckartshausen which was published on 1788 and is owned by the Bavarian State Library [25]. The HW, TW and TW-1 databases were created using a semi-automatic procedure that relies on clustering as presented in [26]. For word recognition the IAM v3.0 database [29] was employed, as described in [39]. Finally, in order to demonstrate the results of the word spotting algorithm a set of historical handwritten images from George Washington's collection from the Library of Congress [27] was used.

Regarding the classification step, the Support Vector Machines (SVM) algorithm was adopted in conjunction with the Radial Basis Function (RBF) kernel [28].

In case of character recognition Tables 1, 2, 3, 4 and 5 show the experimental results for CIL, HW, TW, TW-1 and CEDAR databases, while Table 6 depicts the recognition accuracy for word recognition using the IAM database.

Table 1. Recognition Rates for CIL Database

CIL Database	
Zones	88.48%
Projections	87.75%
Distances	82.53%
Profiles	83.25%
Zones + Projections [30]	91.68%
Zones + Projections + Distances + Profiles with LDA [16]	92.05%
Hierarchical Classification v.1 [31]	93.21%
Hierarchical Classification v.3 [17]	95.63%

Table 2. Recognition rates for HW,TW and TW-1 Databases

	HW	TW	TW-1
Zones + Projections [23]	94.62%	95.44%	NA
Hierarchical Classification v.1 [32]	94.51%	97.71%	NA
Hierarchical Classification v.3 [17]	95.21%	98.24%	99.53 %

Table 3. Recognition rates for the CEDAR Database (52 Classes, a-z/A-Z)

CEDAR Database			
	Uppercase Characters	Lowercase Characters	Overall Recognition Rate
YAM[33]	NA	NA	75.70%
KIM [34]	NA	NA	73.25%
GAD[35]	79.23%	70.31%	74.77%
Hierarchical Classification v.3 [17]	86.17%	84.05%	85.11%

Table 4. Recognition rates for the CEDAR Database for uppercase only and lowercase only characters.

CEDAR Database						
	Uppercase Characters (26 Classes)			Lowercase Characters (26 Classes)		
	# Train Patterns	# Test Patterns	Recognition Rate	# Train Patterns	# Test Patterns	Recognition Rate
BLU[36]	7175	939	81.58%	18655	2240	71.52%
Hierarchical Classification v.3 [17]	11454	1367	95.90%	7691	816	93.50%

Table 5. Recognition rates for the CEDAR Database after merging lowercase and uppercase characters with similar shapes.

CEDAR Database				
	Number of Classes (all classes)	Recognition Rate	Number of Classes (after merging)	Recognition Rate
SIN [37]	52	NA	36	67%
CAM [38]	52	83.74%	39	84.52%
Hierarchical Classification v.3 [17]	52	85.11%	35	94.73%

Table 6. Recognition rates for the IAM Database.

IAM Database	
GAT [39]	87.68%
Hierarchical Classification v.3 (for words) [18]	90.56%

Regarding the experiments for word spotting two datasets from [25] (dataset-1) and from [27] (dataset-2) were used, for evaluating the proposed feature extraction technique, consisting of 13 and 10 document images respectively. Moreover, three words from each dataset were used as keywords: “Durchleucht”,

“nicht” and “Natur” that appear 10, 21, and 17 times respectively in dataset-1 and “public”, “appointments” and “government” that appear 9, 10 and 8 times respectively in dataset-2. Tables 7 and 8 present the F -measure rates using different values of threshold thr .

Table 7. F -Measure for dataset-1

Keyword	Threshold (thr)					
	5	10	15	20	25	30
Durchleucht	66.67%	50%	40%	40%	40%	40%
nicht	38.45%	64.50%	83.32%	97.55%	86.95%	82.35%
Natur	45.45%	74.07%	67.28%	70.27%	71.42%	72.38%

Table 8. F -Measure for dataset-2

Keyword	Threshold (thr)					
	5	10	15	20	25	30
appointments	75%	77.77%	80%	72.72%	74.99%	69.23%
public	80%	82.34%	73.67%	76.19%	69.56%	72%
government	42.85%	37.5%	33.34%	30%	36.36%	35.71%

8 Concluding Remarks

In this thesis novel methodologies that assist handwritten and historical document recognition are presented. In particular: An efficient feature extraction using different types of features followed by a dimensionality reduction step is proposed. Moreover, a novel feature extraction based on recursive subdivisions of the image is introduced. Even though the feature extraction method itself is quite efficient when a specific level of granularity is used, there is more to be gained in classification accuracy by exploiting the intrinsically recursive nature of the method. This is achieved by appropriately combining the results from different levels using a hierarchical approach. Several databases, historical or contemporary, were used to evaluate the performance of these methodologies. In all cases the experimentations depicted, regarding other state-of-the-art techniques that the recognition rates either for characters or words are among the highest one can find in the literature. Also, a new word-spotting algorithm is suggested that relies on the combination of features extracted at different levels of granularity. Finally, a methodology for automatic unsupervised parameter selection for character segmentation is proposed. The methodology is based on clustering; suggesting that the optimal segmentation output, relying on a set of parameters, should produce the best clustering. Experimental results, based on evaluation of segmentation using the ground truth, show that the proposed methodology is capable of finding the optimal or near optimal parameter set.

Figure 6 shows the recent achievements in character recognition. It is obvious that handwritten character recognition as well as historical character recognition that consist of symbols or ligatures that no longer exist in modern alphabets, are still active in research.

		Machine Printed			Handwritten		
		Single Font	Omni Font	Multi Font	Discrete	Cursive	Mixed
On-Line	Constrained				2	3	3
	Unconstrained				1	1	1
Off-Line	Noiseless	2	2	2	3	1	1
	Noisy	1	1	1	1	1	1

1

Need more research

2

Needs improvement

3

Well Done

Figure 6. Recent achievements in OCR

References

1. Luiz S. Oliveira, F. Bortolozzi, C.Y.Suen, "Automatic Recognition of Handwritten Numerical Strings: A Recognition and Verification Strategy", IEEE Transactions on Pattern Recognition and Machine Intelligence, 2001, Vol. 24, No. 11, pp. 1448-1456.
2. K. M. Mohiuddin and J. Mao, "A Comprehensive Study of Different Classifiers for Hand-printed Character Recognition", Pattern Recognition, Practice IV, 1994, pp. 437- 448.
3. L. A. Koerich, "Unconstrained Handwritten Character Recognition Using Different Classification Strategies", International Workshop on Artificial Neural Networks in Pattern Recognition (ANNPR), 2003.
4. N. Arica and F. Yarman-Vural, "An Overview of Character Recognition Focused on Off-line Handwriting", IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews, 2001, 31(2), pp. 216 - 233.
5. O. D. Trier, A. K. Jain, T. Taxt, "Features Extraction Methods for Character Recognition – A Survey ", Pattern Recognition, 1996, Vol.29, No.4, pp. 641-662.
6. S. Belongie, J. Malik, J. Puzicha, "Shape Matching and Object Recognition Using Shape Contexts", IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol.24, No. 4, pp. 509-522, 2002.
7. Anil K. Jain, Douglas Zongker, "Representation and Recognition of Handwritten Digits using Deformable Templates", IEEE Transactions on Pattern Analysis and Machine Intelligence, 1997, Vol. 19, No. 12, pp. 1386-1391.
8. V.G.Gezerlis and S.Theodoridis, "Optical Character Recognition for the Orthodox Hellenic Byzantine music notation", Pattern Recognition, 2002, Vol.35, pp. 895 – 914.
9. K. Ntzios, B. Gatos, I. Pratikakis, T. Konidakis and S.J. Perantonis, "An Old Greek Handwritten OCR System based on an Efficient Segmentation-free Approach", International Journal on Document Analysis and Recognition (IJ DAR), Special Issue on Historical Documents, 2007, Vol. 9, No. 2-4, pp. 179-192.
10. J. Park, V. Govindaraju, S. N. Shrihari, "OCR in Hierarchical Feature Space", IEEE Transactions on Pattern Analysis and Machine Intelligence, 2000, Vol. 22, No. 24, pp. 400-408.
11. L. Vuurpijl, L. Schomacker and M. van Erp, "Architecture for detecting and solving conflicts: two-stage classification and support classifiers", International Journal on Document Analysis and Recognition (IJ DAR), 2004. Vol. 5, No. 4, pp.213-223.
12. Tapan Kumar Bhowmik, Pradip Ghanty, Anandarup Roy and Swapam Kumar Parui, "SVM-Based Hierarchical Architectures for Handwritten Bangla Character Recognition", Int. Journal of Document Analysis and Recognition (IJ DAR), 2009, 12 (2), pp. 97-108.
13. T.M.Rath and R. Manmatha, "Word spotting for historical documents", International Journal on Document Analysis and Recognition (IJ DAR), 2006, Vol.9, No 2 – 4, pp. 139 – 152.
14. V. Lavrenko, T. M. Rath, R. Manmatha: "Holistic Word Recognition for Handwritten Historical Documents", Proceedings of the First International Workshop on Document Image Analysis for Libraries (DIAL'04), 2004, pp 278-287.
15. T. Adamek, N. E. O'Connor, A. F. Smeaton, "Word Matching Using Single-Closed Contours for Indexing Handwritten Historical Documents", International Journal on Document Analysis and Recognition (IJ DAR), Special Issue on Analysis of Historical Documents, 2006.
16. G. Vamvakas, B. Gatos, S. Petridis and N. Stamatopoulos, "An Efficient Feature Extraction and Dimensionality Reduction Scheme for Isolated Greek Handwritten Character Recognition", Proceedings of the 9th International Conference on Document Analysis and Recognition, Curitiba, Brazil, 2007, pp. 1073-1077.

17. G. Vamvakas, B. Gatos, S. J. Perantonis, "Handwritten Character Recognition through Two-Stage Foreground Sub-Sampling", *Pattern Recognition*, Vol.43, Issue 8, pp. 2807-2816, 2010.
18. G. Vamvakas, B. Gatos, S. J. Perantonis, "Efficient Character/ Word Recognition based on a Hierarchical Classification Scheme", *International Journal on Document Analysis and Recognition (IJDAR)*, suggested for the Special Issue: ICDAR'09, to appear.
19. Antonacopoulos, A., Karatzas, D., "Semantics-based content extraction in typewritten historical documents", in: *Eighth International Conference on Document Analysis and Recognition*, 48-53, 2005.
20. Liang, S., Shridhar, M., and Ahmadi, M., "Segmentation of touching characters in printed document recognition", *Pattern Recognition* 27 (6), 825-840, 1994.
21. Nikolaou, N., Makridis, M., Gatos, B., Stamatopoulos N., and Papamarkos, N., "Segmentation of historical machine-printed documents using Adaptive Run Length Smoothing and skeleton segmentation paths", *Image and Vision Computing*, doi:10.1016/j.imavis.2009.09.013, 2009.
22. Kavallieratou E., Stamatatos E., Fakotakis N., Kokkinakis G., "Handwritten character segmentation using transformation-based learning", *15th International Conference on Pattern Recognition*, vol. 2, pp. 634-637, 2000.
23. G.Vamvakas, N. Stamatopoulos, B.Gatos, S.J.Perantonis, "Automatic Unsupervised Parameter Selection for Character Segmentation", 9th IAPR International Workshop on Document Analysis Systems (DAS'10), pp 409-416, June 9-11, Boston, USA, 2010.
24. J.J. Hull, "A database for handwritten text recognition research", *IEEE Trans. Pattern Anal. Mach. Intell.* 16 (5) (1994) 550-554.
25. Carl von Eckartshausen, 1778, *Aufschlüsse zur Magie aus geprüften Erfahrungen über verborgene philosophische Wissenschaften und verdeckte Geheimnisse der Natur*, Bavarian State Library.
26. G. Vamvakas, B. Gatos, N. Stamatopoulos, S.J. Perantonis, "A Complete Optical Character Recognition Methodology for Historical Documents," 8th IAPR International Workshop on Document Analysis Systems (DAS'08), pp.525-532, Nara, Japan, September 2008.
27. <http://memory.loc.gov/ammem/gwhtml/gwhome.html>
28. Cortes C., and Vapnik, V, "Support-vector network", *Machine Learning*, vol. 20, pp. 273-297, 1997.
29. IAM Handwritten Database v3.0, <http://www.iam.unibe.ch/~fki/iamDB/>.
30. G. Vamvakas, B. Gatos, I. Pratikakis, N. Stamatopoulos, A. Roniotis, S.J. Perantonis, "Hybrid Off-Line OCR for Isolated Handwritten Greek Characters", 4th IASTED International Conference on Signal Processing, Pattern Recognition and Applications (SPPRA'07), pp. 197-202, Innsbruck, Austria, 2007.
31. G. Vamvakas, B. Gatos, S. J. Perantonis, "Hierarchical Classification of Handwritten Characters based on Novel Structural Features" (ICFHR'08), 11th International Conference on Frontiers in Handwriting Recognition, Montreal, Canada, August 2008.
32. G. Vamvakas, B. Gatos, S. J. Perantonis, "A Novel Feature Extraction and Classification Methodology for the Recognition of Historical Documents", 10th International Conference on Document Analysis and Recognition (ICDAR'09), pp 491-495, Barcelona, Spain, July 2009.
33. H. Yamada and Y. Nakano, "Cursive Handwritten Word Recognition Using Multiple Segmentation Determined by Contour Analysis", *IECIE Transactions on Information and System*, Vol. E79-D, pp. 464-470, 1996.
34. F. Kimura, N. Kayahara, Y. Miyake and M. Shridhar, "Machine and Human Recognition of Segmented Characters from Handwritten Words", *International Conference on Document Analysis and Recognition (ICDAR '97)*, Ulm, Germany, 1997, pp. 866-869.
35. P. D. Gader, M. Mohamed and J-H. Chiang, "Handwritten Word Recognition with Character and Inter-Character Neural Networks", *IEEE Transactions on System, Man, and Cybernetics-Part B: Cybernetics*, Vol. 27, 1997, pp. 158-164.
36. M. Blumenstein, X.Y. Liu, B. Verma, "A modified direction feature for cursive character recognition", *IEEE International Joint Conference on Neural Networks*, 2007, Vol.4, pp. 2983-2987.
37. S. Singh and M. Hewitt, "Cursive Digit and Character Recognition on Cedar Database". *International Conference on Pattern Recognition, (ICPR 2000)*, Barcelona, Spain. 2000, pp. 569-572.
38. F. Camastra and A. Vinciarelli, "Combining Neural Gas and Learning Vector Quantization for Cursive Character Recognition", *Neurocomputing*, vol. 51, 2003, pp. 147-159.
39. B. Gatos, I. Pratikakis, A.L. Kesidis and S.J. Perantonis, "Efficient Off-Line Cursive Handwritten Word Recognition", 10th International Workshop on Frontiers in Handwriting Recognition (IWFHR 2006), La Baule, France, October 2006, pp. 121-125.

Information Theory and Signal Processing for (nonlinear) Communication Channels

Konstantinos Xenoulis *

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
`kxen@di.uoa.gr`

Abstract. In this thesis, error decoding probability bounds and achievable rates for linear and nonlinear communications channels are presented. Gallager's upper bound as well as its variations through the Duman–Salehi bound are improved. The proposed technique relies on the new inverse exponential sum inequality and designates a new desirable characteristic for linear codes, that is directly connected with the concept of list decoding. The thesis also presents lower bounds on the capacity of nonlinear channels represented with Volterra series, combining the random coding technique with the theory of martingales. The proposed research follows the main ideas that dominate Shannon's basic work and properly utilizes exponential martingale inequalities in order to bound the probabilities of erroneous decoding regions. The specific analysis is also applied to cases where the noise statistical characteristics (mean value, deviation) remain unknown. The present work improves and extends the bound of Shulman–Feder for the family of binary, linear codes that are permutation invariant under list decoding. A new upper bound on list error decoding probability is presented that combines random coding techniques for non-random codes and decreases double exponentially with respect to the code's block length.

1 Introduction

Error probability evaluation and capacity are significant performance measures of coded information transmission over various communication channels. The high complexities involved in the calculation of error probability necessitates the introduction of efficient bounding techniques. Classical treatments [1] as well as modern approaches [2] provide tight bounds mostly for random and specific families of codes (turbo codes [3], LDPC codes [4]), since the latter are treated more easily than specific codes. Thus the existence of at least one optimum code within these families is assured, but the respective characteristics of the optimum code remain unknown. The development of new bounding techniques is crucial to the accommodation of optimum specific codes, which can achieve arbitrarily low error decoding probability with rates close to the channel's capacity.

* Dissertation Advisor: Nicholas Kalouptsidis, Professor

This summary is organized as follows; First, Section 2 introduces the basic error bounding techniques and the presents an improvement for discrete channels. Section 3 provides achievable rates for nonlinear channels under maximum likelihood and weakly typical set decoding. A new double exponential upper bound on list error decoding probability is presented in Section 4. It applies to specific codes while combines random coding techniques. Finally, concluding results are given in Section 5.

2 Error probability bounds

Let \mathcal{C} be a block code of length N and dimension k , over a field F with q elements. Let also $\mathbf{x}_i \in F^N, i = 0, \dots, q^k - 1$ and $S_d, d = d_{\min}, \dots, N$ denote respectively the codewords and the distance distribution of the code \mathcal{C} , with d_{\min} its minimum distance. For a vector $\mathbf{x} \in F^N$, $wt(\mathbf{x})$ denoted its corresponding Hamming weight. For an arbitrary set of messages \mathcal{M} with cardinality M , a message $m, 0 \leq m \leq M - 1$, is mapped to a codeword \mathbf{x}_m of the above code \mathcal{C} and is transmitted over a discrete communication channel with transition probability $P_N(\mathbf{y}|\mathbf{x}_m)$. \mathbf{y} is the received vector at the output of the channel, also of length N . The set of received vectors is denoted by \mathbf{Y} . Each received vector \mathbf{y} is decoded back onto the set of messages \mathcal{M} , according to the maximum likelihood (ML) rule. For the aforementioned transmission procedure, Gallager's upper bound [1] on the code's error decoding probability yields

$$P_{e|m} \leq \sum_{\mathbf{y} \in \mathbf{Y}} P_N(\mathbf{y}|\mathbf{x}_m) \mathcal{K}_m(\mathbf{y}, C, \lambda, \rho) \quad (1)$$

where

$$\mathcal{K}_m(\mathbf{y}, C, \lambda, \rho) = \left(\sum_{m' \neq m} \left(\frac{P_N(\mathbf{y}|\mathbf{x}_{m'})}{P_N(\mathbf{y}|\mathbf{x}_m)} \right)^\lambda \right)^\rho \quad (2)$$

and $\lambda, \rho \geq 0$. A modified version is provided by the DS2 technique [2, sec. 4.2.2]. Let $G_m(\mathbf{y})$ be an arbitrary nonnegative function over \mathbf{Y} , that may also depend on the transmitted message m . Then, for $\lambda \geq 0$ and $0 \leq \rho \leq 1$,

$$P_{e|m} \leq \left(\sum_{\mathbf{y} \in \mathbf{Y}} P_N(\mathbf{y}|\mathbf{x}_m) G_m(\mathbf{y}) \right)^{1-\rho} \cdot \left(\sum_{m' \neq m} \sum_{\mathbf{y} \in \mathbf{Y}} P_N(\mathbf{y}|\mathbf{x}_m) G_m(\mathbf{y})^{1-\frac{1}{\rho}} \left(\frac{P_N(\mathbf{y}|\mathbf{x}_{m'})}{P_N(\mathbf{y}|\mathbf{x}_m)} \right)^\lambda \right)^\rho. \quad (3)$$

The introduction of a tighter upper bound on the ML error decoding probability is made possible by the following inverse exponential sum inequality.

Theorem 1 (Inverse exponential sum inequality [5]). For positive numbers, $\alpha_1, \alpha_2, \dots, \alpha_N$ and $\beta_1, \beta_2, \dots, \beta_N$

$$\sum_{i=1}^N \beta_i \leq \left(\sum_{i=1}^N \alpha_i \right) \ln \frac{\sum_{i=1}^N \alpha_i e^{\frac{\beta_i}{\alpha_i}}}{\sum_{i=1}^N \alpha_i} \quad (4)$$

with equality if and only if $\frac{\beta_i}{\alpha_i} = \text{const.}$

The inverse exponential sum inequality of theorem 1 is used below in the error decoding probability analysis.

Theorem 2 ([5]). Consider the transmission of an arbitrary set of messages \mathcal{M} over a discrete channel, through the utilization of an (N, R) code \mathcal{C} . Let \mathbf{Y}_m^b denote the set of erroneous received vectors given that the message m is transmitted

$$\mathbf{Y}_m^b = \left\{ \mathbf{y} \in \mathbf{Y} : \exists m' \in \mathcal{M}, m' \neq m, P_N(\mathbf{y}|\mathbf{x}_{m'}) \geq P_N(\mathbf{y}|\mathbf{x}_m) \right\}$$

and

$$\mathcal{L}_m(C, \lambda, \rho) = \min_{\mathbf{y} \in \mathbf{Y}_m^b} \mathcal{K}_m(\mathbf{y}, C, \lambda, \rho).$$

Then the ML word error decoding probability for the specific code, given that the message m is transmitted, is upper bounded for all $\lambda, \rho \geq 0$ by

$$\begin{aligned} P_{e|m} &\leq \mathcal{L}_m(C, \lambda, \rho)^{-1} \left(\sum_{\mathbf{y} \in \mathbf{Y}} P_N(\mathbf{y}|\mathbf{x}_m) \mathcal{K}_m(\mathbf{y}, C, \lambda, \rho) \right) \\ &\leq \sum_{\mathbf{y} \in \mathbf{Y}} P_N(\mathbf{y}|\mathbf{x}_m) \mathcal{K}_m(\mathbf{y}, C, \lambda, \rho). \end{aligned} \quad (5)$$

Theorem 2 provides a bound on the ML word error decoding probability that is tighter than Gallager bound, as noted from the second inequality in (5). Moreover, the DS2 technique can be applied to the second term of the first inequality in (5) for all $\rho \leq 1$, thus leading to a tighter version of the DS2 bound.

Theorem 3 ([5]). Under the assumptions of theorem 2, the ML word error decoding probability is upper bounded for all $\lambda \geq 0, 0 \leq \rho \leq 1$ and any nonnegative function $G_m(\mathbf{y})$ by

$$\begin{aligned} P_{e|m} &\leq \mathcal{L}_m(C, \lambda, \rho)^{-1} \left(\sum_{\mathbf{y} \in \mathbf{Y}} P_N(\mathbf{y}|\mathbf{x}_m) G_m(\mathbf{y}) \right)^{1-\rho} \\ &\quad \left(\sum_{m' \neq m} \sum_{\mathbf{y} \in \mathbf{Y}} P_N(\mathbf{y}|\mathbf{x}_m) G_m(\mathbf{y})^{1-\frac{1}{\rho}} \left(\frac{P_N(\mathbf{y}|\mathbf{x}_{m'})}{P_N(\mathbf{y}|\mathbf{x}_m)} \right)^\lambda \right)^\rho. \end{aligned} \quad (6)$$

2.1 Special Cases of Theorems 2, 3

Discrete Channels and Coset Based Analysis Since the channel is memoryless and output symmetric, it holds

$$\left(\frac{P_N(\mathbf{y}|\mathbf{x}^\dagger)}{P_N(\mathbf{y}|\mathbf{x}_0)} \right)^\lambda = \left((q-1) \frac{1-p}{p} \right)^{\lambda(wt(\mathbf{y})-wt(\mathbf{y}-\mathbf{x}^\dagger))} \quad (7)$$

where \mathbf{x}_0 the all zero codeword. Given a specific $\mathbf{y}^* \in \mathbf{Y}_0^b$, in analogy to [6, eq.(33)], we define

$$\deg(\mathbf{y}^*|\mathbf{x}_0) = \left| \left\{ \mathbf{x}_m \in C, m \neq 0 : wt(\mathbf{y}^* - \mathbf{x}_m) \leq wt(\mathbf{y}^*) \right\} \right| \quad (8)$$

where $\deg(\mathbf{y}^*|\mathbf{x}_0)$ is the number of codewords whose Hamming distance from \mathbf{y}^* is lower than or equal to $wt(\mathbf{y}^*)$. The corresponding ratio (7) for each of the above codewords is greater or equal to 1 so that

$$\left(\sum_{m' \neq 0} \left(\frac{P_N(\mathbf{y}^*|\mathbf{x}_{m'})}{P_N(\mathbf{y}^*|\mathbf{x}_0)} \right)^\lambda \right)^\rho \geq \deg(\mathbf{y}^*|\mathbf{x}_0)^\rho. \quad (9)$$

In analogy again to [6, eq.(43)], for every $\mathbf{y}^* \in \mathbf{Y}_0^b$,

$$\deg(\mathbf{y}^*|\mathbf{x}_0) = \sum_{w=1}^{j_t} B_w^t - 1, \quad \left\lceil \frac{d_{\min}}{2} \right\rceil \leq j_t < N \quad (10)$$

where t denotes the coset of \mathbf{y}^* , $j_t = wt(\mathbf{y}^*)$ and B_w^t is the number of words of weight w in the coset t . In contrary to [6, eq.(43)], the term $B_{j_t}^t$ contributes to the sum in the right hand side of (10), since the inequality in (8) is not strict. Moreover, the absence of codeword \mathbf{x}_0 in the previous definition justifies reducing by 1 the aforementioned sum. Consequently, through (9) and (10),

$$\min_{\mathbf{y}^* \in \mathbf{Y}_0^b} \left(\sum_{m' \neq 0} \left(\frac{P_N(\mathbf{y}^*|\mathbf{x}_{m'})}{P_N(\mathbf{y}^*|\mathbf{x}_0)} \right)^\lambda \right)^\rho \geq \min_{t \in \mathcal{T}} \left(\min_{\lceil d_{\min}/2 \rceil \leq j_t < N} \sum_{w=1}^{j_t} B_w^t - 1 \right)^\rho \quad (11)$$

where \mathcal{T} is the set of all cosets t of the code \mathcal{C} .

Theorem 4 ([5]). *Under the assumptions of theorem 2, the ML word error decoding probability is upper bounded for all $\lambda \geq 0, 0 \leq \rho \leq 1$ and any nonnegative function $g(y)$ by*

$$P_{e|0} \leq \left(\min_{t \in \mathcal{T}} \left(\min_{\lceil d_{\min}/2 \rceil \leq j < N} \sum_{w=1}^j B_w^t - 1 \right)^\rho \right)^{-1} \left(\sum_y \Pr(y|0) g(y) \right)^{N(1-\rho)} \cdot \left(\sum_{d=d_{\min}}^N S_d \left(\sum_y \Pr(y|0) g(y)^{1-\frac{1}{\rho}} \right)^{N-d} \left(\sum_y \Pr(y|0)^{1-\lambda} \Pr(y|1)^\lambda g(y)^{1-\frac{1}{\rho}} \right)^d \right)^\rho. \quad (12)$$

Example 1. Consider the perfect Hamming code of length 7 with its coset weight distribution depicted in Table [7, p.170 ex. (1)]. Since the minimum distance of the code is 3, all cosets with minimum weight at least $\lceil 1.5 \rceil$ are examined. Then for $j_t = 2, \dots, 7$,

$$\min_{t \in \mathcal{T}} \left(\min_{2 \leq j_t < 7} \sum_{w=1}^{j_t} B_w^t - 1 \right)^\rho = (3 + 1 - 1)^\rho. \quad (13)$$

The minimum value is achieved for $j_t = 2$, since for $j_t > 2$, the sum over \mathcal{T} in the right hand side of (11) increases. Actually, since the minimum distance of the code is an odd number, there will always exist a term in the left hand side of (11) strictly greater than one.

3 Achievable rates for nonlinear channels

Random coding theorems and achievable rates for nonlinear additive noise channels are presented in this section both under ML and weakly typical decoding. Consider the transmission of an arbitrary set of messages \mathcal{M} with cardinality M over the nonlinear channel

$$\mathbf{y} = D\mathbf{x} + \boldsymbol{\nu}. \quad (14)$$

where \mathbf{x}, \mathbf{y} the corresponding input–output sequences of the channel and $\boldsymbol{\nu}$ the noise vector. The nonlinear behavior of the channel is represented by the Volterra system D applied to the channel’s input sequence $D\mathbf{x}$. The components of the latter vector satisfy

$$[D\mathbf{x}]_i = h_0 + \sum_{j=1}^d \sum_{i_1=0}^{\mu} \dots \sum_{i_j=0}^{\mu} h_j(i_1, \dots, i_j) x_{i-i_1} \dots x_{i-i_j} \quad (15)$$

where it holds

$$\|D\mathbf{x}\|_\infty \leq g_D(\|\mathbf{x}\|_\infty) \leq g_D(r) \quad (16)$$

and

$$g_D(x) = |h_0| + \sum_{j=1}^d \|h_j\| x^j, \quad x \geq 0, \quad \|h_j\| = \sum_{i_1=0}^{\mu} \dots \sum_{i_j=0}^{\mu} |h_j(i_1, \dots, i_j)|. \quad (17)$$

In the sequel we assume input causality i.e. $x_i = 0$ for all $i \leq 0$, and that the noise vector $\boldsymbol{\nu}$ is i.i.d gaussian with zero mean and variance σ_ν^2 . An ML error occurs if, given the transmitted message m and the received vector \mathbf{y} , another message $m' \neq m$ exists such that

$$\|\mathbf{y} - D\mathbf{x}_{m'}\|_2^2 \leq \|\mathbf{y} - D\mathbf{x}_m\|_2^2. \quad (18)$$

Under the random coding setup of Gallager [8, Chap. 5], the average ML error decoding probability $\bar{P}_{e,m}$, given the transmitted message m , satisfies

$$\bar{P}_{e,m} \leq ME \left[\exp \left(-\rho \frac{\|D\mathbf{x} - D\mathbf{x}'\|_2^2}{4\sigma_v^2} \right) \right]. \quad (19)$$

Suppose that x_j, x'_j are the j -th components of the corresponding random vectors \mathbf{x}, \mathbf{x}' . Let

$$\emptyset = F_0 \subset F_1 \subset \cdots \subset F_N, \quad F_i = \{x_1, \dots, x_i, x'_1, \dots, x'_i\} \quad (20)$$

and

$$\begin{aligned} Y_i &= X_i - X_{i-1}, \quad 1 \leq i \leq N, \quad X_i = E \left[\|D\mathbf{x} - D\mathbf{x}'\|_2^2 \mid F_i \right] \\ X_0 &= E \left[\|D\mathbf{x} - D\mathbf{x}'\|_2^2 \right], \quad X_N = \|D\mathbf{x} - D\mathbf{x}'\|_2^2. \end{aligned} \quad (21)$$

We refer to the sequence $\{Y_i\}_{i=1}^N$ as the martingale difference sequence [9] of the random variable X_N with respect to the joint filter $\{F_i\}_{i=0}^N$ in (20). The mean values appearing in (21) are with respect to all codewords the random variables \mathbf{x}, \mathbf{x}' can be assigned to. Under the previous setup, we note that

$$\sum_{i=1}^N Y_i = X_N - X_0 = \|D\mathbf{x} - D\mathbf{x}'\|_2^2 - E \left[\|D\mathbf{x} - D\mathbf{x}'\|_2^2 \right] \quad (22)$$

and thus (19) is equivalently expressed as

$$\bar{P}_{e,m} \leq M \exp \left(-\frac{\rho}{4\sigma_v^2} E \left[\|D\mathbf{x} - D\mathbf{x}'\|_2^2 \right] \right) E \left[\exp \left(-\frac{\rho}{4\sigma_v^2} \sum_{i=1}^N Y_i \right) \right]. \quad (23)$$

Due to the random coding setup and the independency of the ensemble's codewords, it holds

$$E \left[\|D\mathbf{x} - D\mathbf{x}'\|_2^2 \right] = 2 \left(\sum_{j=1}^N E \left[([Du]_j)^2 \right] - E \left[[Du]_j \right]^2 \right) = 2ND_v \quad (24)$$

where $D_v = E \left[([Du])^2 \right] - E \left[[Du] \right]^2$. Finally, combining (23) and (24), we obtain

$$\bar{P}_{e,m} \leq \exp \left(NR - \frac{\rho}{2\sigma_v^2} ND_v \right) E \left[\exp \left(-\frac{\rho}{4\sigma_v^2} \sum_{i=1}^N Y_i \right) \right]. \quad (25)$$

3.1 Random Coding Theorem

The development of exponential upper bounds for the mean value in the right hand side of (25) requires bounds on the conditional deviations $dev^+(Y_i)$ and conditional variances $var(Y_i|F_{i-1})$, where according to [9, pp. 24]

$$\begin{aligned} dev^+(Y_i) &= \max_{x_i, x'_i} Y_i \\ var(Y_i|F_{i-1}) &= E \left[(X_i - X_{i-1})^2 \mid F_{i-1} \right]. \end{aligned} \quad (26)$$

Appropriate bounds are derived in the lemma that follows.

Lemma 1 ([10]). *Under the assumptions that the components of all codewords \mathbf{x}_m , $0 \leq m \leq M-1$ are mutually independent, and r is chosen as in (3), the martingale differences Y_i (21) satisfy*

$$\text{dev}^+(-Y_i) \leq 4(\mu+1)g_D(r)^2, \quad \text{var}(-Y_i|F_{i-1}) \leq 16(\mu+1)^2g_D(r)^4. \quad (27)$$

The bounds provided by Lemma 1 lead to random coding upper bounds on the average ML error decoding probability. Tighter bounds can be obtained analytically for Volterra systems D of short memory.

Theorem 5 ([10]). *Consider the transmission of an arbitrary set of messages \mathcal{M} over a nonlinear Volterra additive gaussian noise channel (14). The components of the noise vector are i.i.d. random variables with 0 mean value and variance σ_v^2 . For each message m , $0 \leq m \leq M-1$, an N -length codeword \mathbf{x}_m is selected from the ensemble \mathcal{C} of (N, R) block codes with probability Q , independently from all other codewords, and is transmitted over the channel. If ML decoding is performed at the receiver and the assumptions of Lemma 1 about the codewords' components are valid, then the average error decoding probability \bar{P}_e is upper bounded as*

$$\bar{P}_e \leq e^{-N(E_c(Q, D, \sigma_v^2) - R)} \quad (28)$$

where

$$E_c(Q, D, \sigma_v^2) = \begin{cases} \frac{1}{2\sigma_v^2} \mathcal{D}_v - \left(\exp\left(\frac{\kappa}{4\sigma_v^2}\right) - 1 - \frac{\kappa}{4\sigma_v^2} \right), & \mathcal{D}_v > \frac{\kappa}{2} \left(\exp\left(\frac{\kappa}{4\sigma_v^2}\right) - 1 \right) \\ \frac{1}{\kappa} \left(-2\mathcal{D}_v + (\kappa + 2\mathcal{D}_v) \ln\left(1 + \frac{2\mathcal{D}_v}{\kappa}\right) \right), & \text{otherwise} \end{cases} \quad (29)$$

and $\kappa = 4(\mu+1)g_D(r)^2$.

Corollary 1. *All rates below $\max_Q E_c(Q, D, \sigma_v^2)$ (29) are achievable for transmission of information over a nonlinear additive gaussian noise channel under ML decoding.*

3.2 Weakly Typical Set Decoding for Nonlinear Systems

In this section, decoding rules for nonlinear channels are interpreted as concentration measures, and martingale theory is utilized. The analysis can be applied to cases where the channel's transition probability law is generally unknown or a suboptimum decoding algorithm is adopted. The nonlinear model (14) is undertaken using the correlation measure $W(\mathbf{x}, \mathbf{y}) = (D\mathbf{x})^T \mathbf{y} = (D\mathbf{x})^T (D\mathbf{x} + \boldsymbol{\nu})$. The input output pair (\mathbf{x}, \mathbf{y}) is called weakly ϵ -typical, if

$$W(\mathbf{x}, \mathbf{y}) \geq E_{\text{Pr}(\mathbf{x}, \mathbf{y})} [W(\mathbf{x}, \mathbf{y})] - N\epsilon. \quad (30)$$

Under the weakly typical decoding rule and the random coding setup, an error occurs given that message m is transmitted, either if codeword \mathbf{x}_m is selected

form the ensemble \mathcal{C} such that $(\mathbf{x}_m, \mathbf{y})$ does not satisfy (30) or if there exists another message $m' \neq m$ for which $\mathbf{x}_{m'}$ is selected independently of \mathbf{x}_m such that $(\mathbf{x}_{m'}, \mathbf{y})$ satisfies (30). Thus, the average error decoding probability, given that m is transmitted, equals

$$\bar{P}_{e,m} = \Pr \left((\mathbf{x}_m, \mathbf{y}) \text{ not } \epsilon\text{-typical} \bigcup_{m' \neq m} (\mathbf{x}_{m'}, \mathbf{y}) \epsilon\text{-typical} \right) \quad (31)$$

and is upper bounded due to the union bound and the Chernoff bound [8, eq. (5.4.11)] for $\lambda_1, \lambda > 0$ as

$$\begin{aligned} \bar{P}_{e,m} &\leq E [\exp (\lambda_1 (E [W(\mathbf{x}, \mathbf{y})] - N\epsilon - W(\mathbf{x}_m, \mathbf{y})))] \\ &+ \sum_{m' \neq m} E_{Q(\mathbf{x}_{m'})P_N(\mathbf{y})} [\exp (\lambda (W(\mathbf{x}_{m'}, \mathbf{y}) - E[W(\mathbf{x}, \mathbf{y})] + N\epsilon))]. \end{aligned} \quad (32)$$

The product probability $Q(\mathbf{x}_{m'})P_N(\mathbf{y})$ in the innermost term in the right hand side of (32) is a direct consequence of the random coding setup. Indeed, $\mathbf{x}_{m'}$ is independent of \mathbf{x}_m and consequently of \mathbf{y} . Noting that $\mathbf{x}'_m, \mathbf{x}_m$ are dummy variables in the above mean values, (32) satisfies

$$\begin{aligned} \bar{P}_{e,m} &\leq E [\exp (\lambda_1 (E [W(\mathbf{x}, \mathbf{y})] - N\epsilon - W(\mathbf{x}, \mathbf{y})))] \\ &+ ME_{Q(\mathbf{x}')P_N(\mathbf{y})} [\exp (\lambda (W(\mathbf{x}', \mathbf{y}) - E[W(\mathbf{x}, \mathbf{y})] + N\epsilon))]. \end{aligned} \quad (33)$$

The following lemma is crucial in the development of exponential martingale inequalities and is used in the proof of the random coding theorem under the weakly typical decoding rule.

Lemma 2 ([10]). *Suppose that all noise samples $\nu_i, i \in [1, N]$ are normally distributed $\mathcal{N}(0, \sigma_\nu^2)$. Then for any $\lambda > 0$ and $0 < k < 1$*

$$\begin{aligned} E_{x_i, x'_i} [\exp (\lambda Y'_i) | F'_{i-1}] &\leq \exp \left(\frac{\lambda^2}{2k} g_D(r)^2 \sigma_\nu^2 \right) \\ &\cdot \left(\frac{1}{2} \exp \left(-\frac{\lambda}{1-k} b' \right) + \frac{1}{2} \exp \left(\frac{\lambda}{1-k} b' \right) \right). \end{aligned} \quad (34)$$

Theorem 6 ([10]). *Let the transmission of an arbitrary set of messages \mathcal{M} over an additive noise nonlinear channel, under the same random encoding setup of Theorem 5. Let also the noise samples be i.i.d. and normally distributed $\mathcal{N}(0, \sigma_\nu^2)$, independent from the channel input. Then, for any $\epsilon, \epsilon_1 > 0$ arbitrarily small constants, the average error decoding probability \bar{P}_e is upper bounded as*

$$\bar{P}_e \leq \epsilon_1 + e^{-N(E'_c(Q, D, \sigma_\nu^2) - R)} \quad (35)$$

where

$$\begin{aligned} E'_c(Q, D, \sigma_\nu^2) &= \max_{0 < k < 1} \max_{0 < \lambda} \lambda (2D_v - \epsilon) - \frac{\lambda^2}{2k} g_D(r)^2 \sigma_\nu^2 - \\ &\ln \left(\frac{1}{2} \exp \left(-\frac{\lambda}{1-k} b' \right) + \frac{1}{2} \exp \left(\frac{\lambda}{1-k} b' \right) \right) > 0. \end{aligned} \quad (36)$$

Corollary 2. *Considering the transmission of information over a nonlinear Volterra additive gaussian noise channel, all rates below $\max_Q E'_c(Q, D, \sigma_\nu^2)$ (36) are achievable for the weakly typical set decoding rule.*

The tightness of the random coding exponent, given by Theorem 6, depends on lower bounds for the error decoding probability of the form provided in [11], for the specific functions $W(\mathbf{x}, \mathbf{y})$.

4 Random coding techniques for nonrandom codes

In this section, a new double exponential upper bound on the list error decoding probability of specific classes of codes over binary input symmetric output memoryless channels is derived.

When list decoding is performed at the output of the channel with list size L , the conditional error decoding probability of \mathcal{C} , given the transmission of message 0, satisfies

$$P_{e|0,\mathcal{C}}^\mathcal{L} = \sum_{\mathbf{y} \in \mathbf{Y}_{0,\mathcal{C}}^\mathcal{L}} P_N(\mathbf{y}|\mathbf{x}_0, \mathcal{C}) \quad (37)$$

where

$$\mathbf{Y}_{0,\mathcal{C}}^\mathcal{L} = \{\mathbf{y} \in \mathcal{J}^N : \exists \{l_i\}_{i=1}^L, l_i \neq 0 : P_N(\mathbf{y}|\mathbf{x}_{l_i}, \mathcal{C}) \geq P_N(\mathbf{y}|\mathbf{x}_0, \mathcal{C}), \forall i \in [1, L]\}. \quad (38)$$

Moreover, if for $\lambda, \rho \geq 0$ we set

$$\Omega_L(\mathbf{y}, \mathcal{C}, \lambda, \rho) = \frac{L^\rho P_N(\mathbf{y}|\mathbf{x}_0, \mathcal{C})^{\lambda\rho}}{\left(\sum_{m \neq 0} P_N(\mathbf{y}|\mathbf{x}_m, \mathcal{C})^\lambda\right)^\rho} \quad (39)$$

then, due to the definition in (38), the error decoding probability $P_{e|0,\mathcal{C}}^\mathcal{L}$ in (37) is upper bounded as

$$P_{e|0,\mathcal{C}}^\mathcal{L} \leq \sum_{\mathbf{y} \in \mathbf{Y}_{0,\mathcal{C}}^\mathcal{L}} P_N(\mathbf{y}|\mathbf{x}_0, \mathcal{C}) e^{1 - e^{\Omega_L(\mathbf{y}, \mathcal{C}, \lambda, \rho) - 1}}. \quad (40)$$

The current work is confined to the following L -list permutation invariant codes.

Definition 1 ([12]). *An (N, R) linear binary code \mathcal{C} with coset weight distribution matrix $\mathbf{\Gamma}$ is L -list permutation invariant if both the following properties are satisfied:*

\mathfrak{L}_1 : *there exists a $w_{opt} \geq \lceil d_{\min}/2 \rceil$ such that*

$$L = \min_{\kappa \in [1, K], \mathbf{\Gamma}_{\kappa, w_{opt}} \neq 0} \mathbf{\Gamma}_{\kappa, w_{opt}} - 1 > 0 \quad \text{and} \quad \max_{\kappa \in [1, K], w < w_{opt}} \mathbf{\Gamma}_{\kappa, w} < L + 1.$$

\mathfrak{L}_2 : *For all $\kappa \in [1, K]$, there exists a $w_\kappa^L > w_{opt}$ such that*

$$\mathbf{\Gamma}_{\kappa, w_{opt}+1} = \dots = \mathbf{\Gamma}_{\kappa, w_\kappa^L-1} = 0 \quad \text{and} \quad \mathbf{\Gamma}_{\kappa, w_\kappa^L} \geq L + 1$$

From an L -list permutation invariant code \mathcal{C} , we construct an ensemble of codes \mathcal{E} by considering all possible symbol position permutation $N \times N$ matrices \mathcal{P} . A position permutation matrix \mathcal{P} has a single 1 in every row and every column and is orthogonal, $\mathcal{P}\mathcal{P}^T = \mathcal{P}^T\mathcal{P} = I$. I is the $N \times N$ identity matrix and \mathcal{P}^T the transpose matrix of \mathcal{P} . The lemma provided below is crucial in the derivation of the double exponential bound on the list error decoding probability.

Lemma 3 ([12]). *For an (N, R) binary linear block code \mathcal{C} that is L -list permutation invariant, all codes in the permuted ensemble \mathcal{E} have the same error decoding region $\mathbf{Y}_0^\mathcal{C}$.*

Due to the channel symmetry, the average list error decoding probability $P_e^\mathcal{C}$ of any code \mathcal{C} , over all messages in \mathcal{M} , equals $P_{e|0, \mathcal{C}}^\mathcal{C}$ [13, Appendix C]. Thus, any bound on $P_{e|0, \mathcal{C}}^\mathcal{C}$ is also a bound on $P_e^\mathcal{C}$. Moreover, for a L -list permutation invariant code \mathcal{C} , in all codes of the ensemble \mathcal{E} , message 0 is encoded into the all-zeros vector \mathbf{x}_0 . Thus, $P_N(\mathbf{y}|\mathbf{x}_0, \mathcal{C}) = P_N(\mathbf{y}|\mathbf{x}_0)$. Consequently, if we take the average over \mathcal{E} on both sides of (40), then due the error decoding region invariance property stated in lemma 3, we have

$$P_{e|0}^\mathcal{C} \leq \sum_{\mathbf{y} \in \mathbf{Y}_0^\mathcal{C}} P_N(\mathbf{y}|\mathbf{x}_0) E \left[e^{1 - e^{\Omega_L(\mathbf{y}, \mathcal{C}, \lambda, \rho) - 1}} \right]. \quad (41)$$

Note that the function $\exp(1 - \exp(x - 1))$ is concave for $0 \leq x \leq 1$ since

$$\frac{de^{1 - e^{x-1}}}{dx^2} = e^{-1 - e^{-1+x} + x} (-e + e^x) \leq 0, \text{ for } 0 \leq x \leq 1.$$

Moreover, for any $\mathbf{y} \in \mathbf{Y}_0^\mathcal{C}$, $\Omega_L(\mathbf{y}, \mathcal{C}, \lambda, \rho) \leq 1$. Therefore, application of Jensen's inequality to the right hand side of (41) gives

$$P_e^\mathcal{C} \leq \sum_{\mathbf{y} \in \mathbf{Y}_0^\mathcal{C}} P_N(\mathbf{y}|\mathbf{x}_0) e^{1 - e^{E[\Omega_L(\mathbf{y}, \mathcal{C}, \lambda, \rho)] - 1}}. \quad (42)$$

The following technical lemma is useful in the derivation of a closed form upper bound on $P_e^\mathcal{C}$.

Lemma 4 ([12]). *The mean value of the double exponent in (42) is lower bounded for all $\rho' \geq 0$ as*

$$\begin{aligned} E[\Omega_L(\mathbf{y}, \mathcal{C}, \lambda, \rho)] &\geq L^{\rho' P} P_N(\mathbf{y}|\mathbf{x}_0)^{\frac{\rho'}{1+\rho'}} \\ &\cdot \left[(M-1)^{\rho' P} \left(\sum_{l=0}^N b_l \left(\frac{v_l}{b_l} \right)^Q \right)^{\rho' \frac{P}{Q}} \left(\sum_{\mathbf{x} \in \mathcal{I}^N} 2^{-N} P_N(\mathbf{y}|\mathbf{x})^{\frac{1}{1+\rho'}} \right)^{\rho'} \right]^{-1} \end{aligned} \quad (43)$$

where

$$b_l = \frac{\binom{N}{l}}{2^N}, \quad v_l = \frac{S_l}{M-1}, \quad 0 \leq l \leq N, \quad \frac{1}{P} + \frac{1}{Q} = 1, \quad P, Q \geq 1. \quad (44)$$

Combining lemma 4 with (42) and passing from $\mathbf{Y}_0^{\mathcal{L}}$ to the set of all received vectors \mathbf{Y} , we get the following theorem.

Theorem 7 ([12]). *Consider an (N, R) binary linear block code \mathcal{C} which is L -list permutation invariant with distance spectrum S_l , $0 \leq l \leq N$ and coset weight distribution matrix $\mathbf{\Gamma}$. \mathcal{C} is utilized in the transmission of an arbitrary set of messages \mathcal{M} , with cardinality $M = 2^{NR}$, over a binary input, symmetric output discrete memoryless channel. If $p/(q-1)$ is the error transition probability of the channel, then the average list error decoding probability, over all messages in \mathcal{M} , $P_e^{\mathcal{L}}$ of \mathcal{C} is upper bounded for all $\rho' \geq 0$ as*

$$P_e^{\mathcal{L}} \leq \sum_{h=0}^N \binom{N}{h} (1-p)^{N-h} \left(\frac{p}{q-1}\right)^h \sum_{k=\delta(q-2)h}^h \binom{h}{k} (q-2)^{h-k} \exp(1) \cdot \exp \left(- \exp \left(\frac{2^{\rho'N} (M-1)^{-\rho'P} L^{\rho'P} \left((1-p)^{N-h} \left(\frac{p}{q-1}\right)^h \right)^{\frac{\rho'}{1+\rho'}}}{\left(\sum_{l=0}^N b_l \left(\frac{p_l}{b_l}\right)^Q \right)^{\rho' \frac{P}{Q}} \mathcal{K}_1(p, q, \rho')^{N-h+k} \mathcal{K}_2(p, q, \rho')^{h-k}} - 1 \right) \right) \quad (45)$$

where

$$\mathcal{K}_1(p, q, \rho') = \left((1-p)^{\frac{1}{1+\rho'}} + \left(\frac{p}{q-1}\right)^{\frac{1}{1+\rho'}} \right)^{\rho'}, \quad \mathcal{K}_2(p, q, \rho') = \left(\frac{p}{q-1}\right)^{\frac{\rho'}{1+\rho'}} \\ \frac{1}{P} + \frac{1}{Q} = 1, \quad P, Q \geq 1 \quad \text{and} \quad \delta(q-2) = \begin{cases} 1, & q = 2 \\ 0, & \text{otherwise} \end{cases} \quad (46)$$

We note that the upper bound of theorem 7 fails to reproduce the random coding exponent for an L -list permutation invariant code \mathcal{C} , as in [14, Th.1]. Additionally, it does not admit a closed form expression for continuous output channel. Nevertheless, since

$$e^{1-e^{x-1}} \leq \frac{1}{x}, \quad x > 0 \quad (47)$$

(45) is tighter than the generalized version of Shulman-Feder bound in [15, eq.(A17)], [13, Cor.8]. Moreover, for L -list permutation invariant codes, application of (47) in (45) provides a new version of the generalized SFB, which nicely complements the one presented in [15, eq.(A17)].

5 Conclusions

This thesis deals with issues regarding reliable and efficient information transmission over linear and nonlinear communication channels. For discrete linear

symmetric channels, improved upper bounds are developed under maximum likelihood decoding. Furthermore, double exponential upper bounds on the list decoding error probability of specific codes are presented that combine random coding techniques. Finally, the thesis presents achievable rates for nonlinear channels both under maximum likelihood and weakly typical set decoding, utilizing properly the theory of martingales.

References

1. R. Gallager, "A simple derivation of the coding theorem and some applications," *IEEE Trans. Inf. Theory*, vol. 11, pp. 3–18, Jan. 1965.
2. I. Sason and S. Shamai, "Performance analysis of linear codes under maximum-likelihood decoding: A tutorial," *Foundations and Trends in Communications and Information Theory*, vol. 3, 2006.
3. C. Berrou, A. Glavieux, and P. Thitimajshima, "Near Shannon limit error-correcting coding and decoding," in *Proc. IEEE (ICC'93) Geneva, Switzerland*, May 1993, pp. 1064–1070.
4. R. Gallager, *Low-Density Parity-Check Codes*. Cambridge, MA: MIT Press, 1963.
5. K. Xenoulis and N. Kalouptsidis, "Improvement of Gallager upper bound and its variations for discrete channels," *IEEE Trans. Inf. Theory*, vol. 55, no. 9, pp. 4204–4210, Sep. 2009.
6. A. Cohen and N. Merhav, "Lower bounds on the error probability of block codes based on improvements on de Caen's inequality," *IEEE Trans. Inf. Theory*, vol. 52, no. 2, pp. 290–310, 2004.
7. F. MacWilliams and N. Sloane, *The Theory of Error-Correcting Codes*. North Holland, 1983.
8. R. Gallager, *Information Theory and Reliable Communication*. John Wiley and Sons, New York, 1968.
9. C. McDiarmid, "Concentration," *Probabilistic Methods for Algorithmic Discrete Mathematics*, pp. 195–248, 1998.
10. K. Xenoulis and N. Kalouptsidis, "Random coding theorems for nonlinear channels," *IEEE Trans. Inf. Theory*, accepted under changes, 2010.
11. S. Verdú and T. S. Han, "A general formula for channel capacity," *IEEE Trans. Inf. Theory*, vol. 40, no. 4, pp. 1147–1157, 1994.
12. K. Xenoulis and N. Kalouptsidis, "Double exponential performance bounds for permutation invariant binary linear block codes over symmetric channels," *IEEE Trans. Inf. Theory*, submitted, Dec. 2009.
13. E. Hof, I. Sason, and S. Shamai, "Performance bounds for erasure, list and feedback schemes with linear block codes," Aug. 2010, to appear *IEEE Trans. Inf. Theory*.
14. N. Shulman and M. Feder, "Random coding techniques for nonrandom codes," *IEEE Trans. Inf. Theory*, vol. 45, no. 6, pp. 2101–2104, 1999.
15. S. Shamai and I. Sason, "Variations on the Gallager bounds, connections, and applications (ComSoc & IT joint paper award 2003)," *IEEE Trans. Inf. Theory*, vol. 48, no. 12, pp. 3029–3051, 2002.