NATIONAL AND KAPODISTRIAN UNIVERSITY OF ATHENS
DEPARTMENT OF INFORMATICS & TELECOMMUNICATIONS

# ABSTRACTS OF DOCTORAL DISSERTATIONS

Athens 2014

Volume 9

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

Volume 9

# 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 12/2012 to 12/2013.

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 innovative research results. We are happy to present the 9th 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 graduating Ph.D. students of the Department.

We would like to thank our Ph.D. graduates who contributed to this volume and hope that this has been a positive experience for them. Finally, we would like to thank PhD candidate Nikos Bogdos for his help in putting together this publication. The painting in the cover is called "*Pathways*" by artist *Διονύσης Καμπόλης*.

The Committee of Research and Development

A. Eleftheriadis

M. Koubarakis

E. Manolakos (Chair)

T. Theoharis

Athens, June 2014

# Table of Contents

## Doctoral Dissertations

# Robust Adaptive Machine Learning Algorithms for Distributed Signal Processing

Symeon Chouvardas*

Department of Informatics and Telecommunications

National and Kapodistrian University of Athens

15784, Ilissia, Athens, Greece

## Abstract

Distributed networks comprising a large number of nodes, *e.g.,* Wireless Sensor Networks, Personal Computers (PC's), laptops, smart phones, etc., which cooperate with each other in order to reach to a common goal, constitute a promising technology for several applications. Typical examples include: distributed environmental monitoring, acoustic source localization, power spectrum estimation, etc. Sophisticated cooperation mechanisms can significantly benefit the learning process, through which the nodes achieve their common objective.

In this dissertation, the problem of adaptive learning in distributed networks is studied, focusing on the task of distributed estimation. A set of nodes sense information related to certain parameters and the estimation of these parameters comprises the goal. Towards this direction, nodes exploit locally sensed measurements as well as information springing from interactions with other nodes of the network. Throughout this dissertation, the cooperation among the nodes follows the diffusion optimization rationale and the developed algorithms belong to the APSM algorithmic family.

First, robust APSM–based techniques are proposed. The goal is to "harmonize" the spatial information, received from the neighborhood, with the locally sensed one. This "harmonization" is achieved by projecting the information of the neighborhood onto a convex set, constructed via the locally sensed measurements. Next, the scenario,

in which a subset of the node set is malfunctioning and produces measurements heavily corrupted with noise, is considered. This problem is attacked by employing the Huber cost function, which is resilient to the presence of outliers. In the sequel, we study the issue of sparsity–aware adaptive distributed learning. The nodes of the network seek for an unknown sparse vector, which consists of a small number of non–zero coefficients. Weighted $\ell_1$–norm constraints are embedded, together with sparsity–promoting variable metric projections. Finally, we propose algorithms, which lead to a reduction of the communication demands, by forcing the estimates to lie within lower dimensional Krylov subspaces. The derived schemes serve a good trade-off between complexity/bandwidth demands and achieved performance.

**Subject Area:** Adaptive Learning, Distributed Signal Processing.

**Keywords:** Diffusion, Projections, APSM, hyperslabs.

# 1    Introduction

Distributed networks comprising a number of connected nodes, *e.g.,* Personal Computers (PC's), laptops, smart phones, surveillance cameras and microphones, wireless sensor networks etc., which exchange information in order to reach a common goal, are envisioned to play a central role in many applications. Typical examples of emergent applications involving distributed networks are: distributed environmental monitoring, acoustic source localization, power spectrum estimation, target tracking, surveillance, traffic control, patient monitoring and hospital surveillance, just to name a few [1,3,6,7,11]. All the previously mentioned applications share in common the fact that the nodes are deployed over a geographic region providing spatial diversity to the obtained measurements. Henceforth, the development of algorithms and node cooperation mechanisms, which exploit the information diversity over time and space, so that a common objective to be reached, becomes essential.

In this dissertation, the problem of distributed processing is studied with a focus on the distributed/decentralized estimation task. A number of nodes, which are spread over a geographic region, sense information related to certain parameters; the estimation of these parameters comprises our goal. The main idea behind distributed processing is that the nodes exchange information among them and make decisions/computations in a *collaborative* way instead of working individually, using solely the information that is locally sensed. It is by now well established, that the cooperation among the nodes leads to better results compared to the case where they act as individual learners, see for example [5, 10, 13]. The need to develop node cooperation

mechanisms is increased due to the presence of noise in the majority of applications. More specifically, the measurements observed at each node are corrupted by noise, and this fact adds further uncertainty on the obtained estimates of the unknown target parameters. This uncertainty can be reduced via the cooperation of the nodes.

In decentralized networks, the following issues have to be taken into consideration:

- *Performance:* A performance close to the optimal, that is the one associated with the centralized networks, which use all the available data, has to be achieved. In other words, despite the fact that direct communication among some of the nodes cannot be established, sophisticated cooperation mechanisms have to be developed, in order to "push" the performance to be as close as possible to the ideal scenario.

- *Robustness to possible failures:* As it has been already stated, a major drawback of the centralized topology is that if the FC fails then the network collapses. Decentralized networks have to be constructed so as to be robust against possible node failures.

- *Bandwidth and complexity constraints:* The amount of transmitted information has to be as small as possible, in order to keep the bandwidth low. Furthermore, since in decentralized networks a central processing unit with powerful computational capabilities is not present and usually cheap processing units comprise the nodes, low-complexity schemes have to be developed.

- *Adaptivity:* In many applications, such as, source localization, spectrum sensing, etc, the nodes of the network are tasked to estimate non–stationary parameters, *i.e.,* parameters which vary with time. Batch estimation algorithms, which use all the available training data simultaneously, cannot attack such problems. To this end, adaptive techniques have to be developed, where the data are observed sequentially, one per (discrete) time instance and operate in an online fashion for updating and improving the estimates.

The main objective of this dissertation is to develop algorithms in the context of *adaptive estimation* in distributed networks. The diffusion optimization rationale is adopted and the proposed algorithms belong to the Adaptive Projected Subgradient Method (APSM) algorithmic family.

# 2   Adaptive Robust Algorithms for Distributed Learning

As a first step, distributed algorithms, which follow the diffusion rationale and belong to the family of the Adaptive Projected Subgradient Method, are developed. The proposed algorithms adopt a novel combine–project–adapt cooperation protocol. The intermediate extra projection step of this protocol "harmonizes" the local information, which comprises the input/output measurements, with the information coming from the neighborhood, *i.e.,* the estimates obtained from the neighboring nodes. This is achieved by projecting the vector, occurring by combining the estimates of the neighbourhood, to a convex set, namely a hyperslab, which is constructed by exploiting locally sensed information. The steps of the algorithm can be summarised as follows:

1. **Combination Step:** The estimates from the nodes that belong to the neighbourhood are received and convexly combined with respect to the combination weights.

2. **Projection Step:** The resulting aggregate is first projected onto a properly constructed hyperslab.

3. **Adaptation Step:** The adaptation step is performed.

The following model is adopted. A network of $N$ nodes is considered and each node, $k$, at time $n$, has access to the measurements $d_{k,n} \in \mathbb{R}$, $\boldsymbol{u}_{k,n} \in \mathbb{R}^m$ generated by the linear system:

$$d_{k,n} = \boldsymbol{w}_*^T \boldsymbol{u}_{k,n} + v_{k,n}, \tag{1}$$

where $v_{k,n}$ is an additive noise process of zero mean and variance $\sigma_k^2$. The goal is the estimation of the $m \times 1$ vector $\boldsymbol{w}_*$.

As we have already mentioned, an APSM–based scheme, which employs projections onto hyperslabs, is developed. The scheme is brought in a distributed fashion by following the diffusion rationale. Moreover, here an extra step is added, that follows the combination stage and precedes the adaptation one. More specifically, the result of the combination step is projected onto the hyperslab $S'_{k,n}$, which is defined as

$$S'_{k,n} = \{\boldsymbol{w} \in \mathbb{R}^m : |d_{k,n} - \boldsymbol{w}^T \boldsymbol{u}_{k,n}| \leq \epsilon'_k\},$$

where $\epsilon'_k > \epsilon_k$ and $\epsilon_k$ is the user defined parameter associated with the hyperslabs, that will be used in the adaptation step at node $k$, *i.e.,*

$$S_{k,n} = \{\boldsymbol{w} \in \mathbb{R}^m : |d_{k,n} - \boldsymbol{w}^T \boldsymbol{u}_{k,n}| \leq \epsilon_k\}.$$

The algorithm comprises the following steps:

1. **Combination Step:** The estimates from the nodes that belong to $\mathcal{N}_k$ are received and convexly combined with respect to the combination weights $a_{k,l}$.

2. **Projection Step:** The resulting aggregate is first projected onto the hyperslab $S'_{k,n}$[1].

3. **Adaptation Step:** The adaptation step is performed.

$$\boldsymbol{\phi}_{k,n} = \sum_{l \in \mathcal{N}_k} a_{k,l} \boldsymbol{w}_{l,n}, \tag{2}$$

$$\boldsymbol{z}_{k,n} = P_{S'_{k,n}} \left( \boldsymbol{\phi}_{k,n} \right), \tag{3}$$

$$\boldsymbol{w}_{k,n+1} = \boldsymbol{z}_{k,n} + \mu_{k,n} \left( \sum_{j \in \mathcal{J}_n} \omega_{k,j} P_{S_{k,j}}(\boldsymbol{z}_{k,n}) - \boldsymbol{z}_{k,n} \right), \tag{4}$$

where $P_{S'_{k,n}}$ and $P_{S_{k,n}}$ are the projection operators onto the respective hyperslabs, $\sum_{j \in \mathcal{J}_n} \omega_{k,j} = 1$ and $\mathcal{J}_n := \overline{\max\{0, n-q+1\}, n}$. As it was experimentally verified, the proposed scheme exhibits an enhanced performance, both in terms of convergence speed as well as steady state error floor, compared to other state of the art algorithms, of similar complexity. Finally, it was proved that the algorithm enjoys a number of nice convergence properties such as monotonicity, strong convergence to a point and consensus.

# 3  Introducing Robustness to Cope with a Failure of Nodes

Consider a scenario, in which some of the nodes are damaged and the associated observations are very noisy. More specifically, it is assumed that that the noise is additive and white, albeit the standard deviation of the "damaged" nodes becomes larger, compared to the one of the "healthy" nodes. In such cases, the use of loss functions, suggested in the framework of robust statistics, are more appropriate to cope with outliers. A popular cost function of this family is the Huber cost function, *e.g.*, [8, 12].

In the current study we employ a slightly modified version of the Huber cost function, compared to the classical one. The difference is that in our

---

[1]The projection of a point onto a hyperslab is provided in Chapter 3.

context a 0–th level set is introduced, *i.e.,* a set of points in which the function scores a zero loss. Our goal is to find points, which lie in the previously mentioned 0–th level set. The geometry of the Huber function is illustrated in Fig. In contrast to the hyperslab case, the projection onto the 0–th level set of the Huber cost function, does not admit a closed form. For this reason, projections onto the halfspace, associated to the subgradient of the Huber loss function, take place. We can also include the extra projection step, described in the previous section, by introducing a modified version of the Huber cost function and following a similar rationale as in the hyperslab case. However, instead of projecting the result of the combination step onto an external hyberslab, we project it onto a halfspace that is generated by a properly modified cost function. The proposed algorithm comprises the following steps:

$$\phi_{k,n} = \sum_{l \in \mathcal{N}_k} a_{k,l} \boldsymbol{w}_{l,n}, \tag{5}$$

$$\boldsymbol{z}_{k,n} = P_{H'^-_{k,n}} \left( \boldsymbol{\phi}_{k,n} \right) ,, \tag{6}$$

$$\boldsymbol{w}_{k,n+1} = \boldsymbol{z}_{k,n} + \mu_{k,n} \left( \sum_{j \in \mathcal{J}_n} \omega_{k,j} P_{H^-_{k,j}} (\boldsymbol{z}_{k,n}) - \boldsymbol{z}_{k,n} \right), \tag{7}$$

where $P_{H^-_{k,j}}$ stands for the projection onto the halfspace associated to the Huber loss function and $P_{H'^-_{k,j}}$ is the previously described extra projection step. Under some mild assumptions, the developed algorithm enjoys monotonicity, asymptotic optimality, asymptotic consensus and strong convergence to a point that lies in the consensus subspace. Finally, numerical examples verified that the proposed scheme has an enhanced performance, compared to the other methodologies, in a network with malfunctioning nodes.

# 4 Sparsity–Aware Adaptive Distributed Learning

As a next step, an APSM–based sparsity–promoting adaptive algorithm for distributed learning in ad–hoc networks is developed. At each time instance and at each node of the network, a hyperslab is constructed based on the received measurements; this defines the region in which the solution is searched for. Sparsity encouraging variable metric projections onto these sets have been adopted. In addition, sparsity is also imposed by employing variable metric projections onto weighted $\ell_1$ balls. A combine adapt cooperation strategy is followed.

Let us introduce, here, the sparsity promoting *variable metric projection*, onto the respective hyperslabs, with respect to the matrix $\boldsymbol{G}_n$, defined as [14]:

$$\forall \boldsymbol{w} \in \mathbb{R}^m, \quad P_{S_n}^{(\boldsymbol{G}_n)}(\boldsymbol{w}) := \boldsymbol{w} + \beta_n \boldsymbol{G}_n^{-1} \boldsymbol{u}_n, \tag{8}$$

where

$$\beta_n = \begin{cases} \dfrac{d_n - \boldsymbol{u}_n^T \boldsymbol{w} + \epsilon}{\|\boldsymbol{u}_n\|_{\boldsymbol{G}_n^{-1}}^2}, & \text{if } d_n - \boldsymbol{u}_n^T \boldsymbol{w} < -\epsilon, \\ 0, & \text{if } |d_n - \boldsymbol{u}_n^T \boldsymbol{w}| \le \epsilon, \\ \dfrac{d_n - \boldsymbol{u}_n^T \boldsymbol{w} - \epsilon}{\|\boldsymbol{u}_n\|_{\boldsymbol{G}_n^{-1}}^2}, & \text{if } d_n - \boldsymbol{u}_n^T \boldsymbol{w} > \epsilon, \end{cases}$$

and $\|\boldsymbol{u}_n\|_{\boldsymbol{G}_n^{-1}}^2$ denotes the weighted norm, with definition $\|\boldsymbol{u}_n\|_{\boldsymbol{G}_n^{-1}}^2 := \boldsymbol{u}_n^T \boldsymbol{G}_n^{-1} \boldsymbol{u}_n$ (see Appendix C). Note that if $\boldsymbol{G}_n = \boldsymbol{I}_m$, then (8) is the Euclidean projection onto a hyperslab. The positive definite diagonal matrix $\boldsymbol{G}_n^{-1}$ is constructed following similar philosophy as in [2, 15]. The $i$-th coefficient of its diagonal equals to $g_{i,n}^{-1} = \frac{1-\overline{\alpha}}{m} + \overline{\alpha}\frac{|w_i^{(n)}|}{\|\boldsymbol{w}_n\|_1}$, where $\overline{\alpha} \in [0, 1)$ is a parameter, that determines the extend to which the sparsity level of the unknown vector will be taken into consideration, and $w_i^{(n)}$ denotes the $i$-th component of $\boldsymbol{w}_n$. In order to grasp the reasoning of the variable metric projections, consider the ideal situation, in which $\boldsymbol{G}_n^{-1}$ is generated by the unknown vector $\boldsymbol{w}_*$. It is easy to verify that $g_{i,n}^{-1} > g_{i',n}^{-1}$, if $i \in \text{supp}(\boldsymbol{w}_*)$, and $i' \notin \text{supp}(\boldsymbol{w}_*)$, where $\text{supp}(\cdot)$ stands for the support set of a vector, *i.e.,* the set of the non–zero coefficients. Hence, employing the variable metric projection, the amplitude of each coefficient of the vector used to construct $\boldsymbol{G}_n^{-1}$ determines the weight that will be assigned to the corresponding coefficient of the second term of the right hand side in (8). That is, components with smaller magnitude are multiplied with small coefficients of $\boldsymbol{G}_n^{-1}$. Loosely speaking, the variable metric projections accelerate the convergence speed when tracking a sparse vector, since by assigning different weights pushes the coefficients of the estimates with small amplitude to diminish faster. The geometric implication of it is that the projection is made to "lean" towards the direction of the more significant components of the currently available estimate.

In the algorithm which is presented here, we go one step further, as far as sparsity is concerned. In a second stage, additional sparsity-related constraints, which are built around the weighted $\ell_1$ ball, are employed, [4]. A sparsity promoting adaptive scheme, based on set-theoretic estimation arguments, in which the constraints are weighted $\ell_1$ balls, was presented in [9]. Given a vector of weights $\boldsymbol{\psi}_n = [h_1^{(n)}, \dots, h_m^{(n)}]^T$, where $h_i^{(n)} > 0, \forall i = 1, \dots, m$, and a positive radius, $\delta$, the weighted $\ell_1$ ball is defined as: $B_{\ell_1}[\boldsymbol{\psi}_n, \delta] := \{\boldsymbol{w} \in \mathbb{R}^m : \sum_{i=1}^m h_i^{(n)} |w_i| \le \delta\}$. The projection onto $B_{\ell_1}[\boldsymbol{\psi}_n, \delta]$,
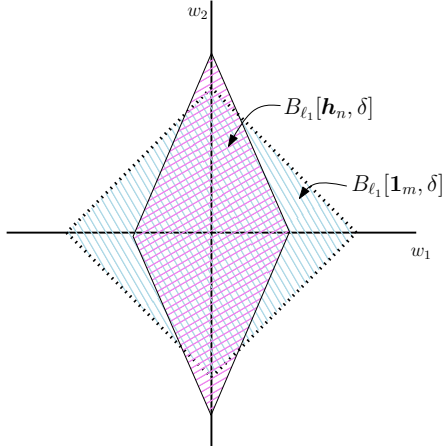
Figure 1: Illustration of a weighted $\ell_1$ ball (solid line magenta) and an unweighted $\ell_1$ ball (dashed line blue).

is given in [9, Theorem 1], and the geometry of these sets is illustrated in Fig. 1.

The steps of the algorithm are summarized in the sequel:

$$\boldsymbol{w}_{k,n+1} = P^{(\boldsymbol{G}_n)}_{B_{\ell_1}[\boldsymbol{\psi}_n,\delta]} \left( \boldsymbol{\phi}_{k,n} + \mu_{k,n} \left( \sum_{j \in \mathcal{J}} \omega_{k,j} P^{(\boldsymbol{G}_n)}_{S_{k,j}}(\boldsymbol{\phi}_{k,n}) - \boldsymbol{\phi}_{k,n} \right) \right), \quad (9)$$

The theoretical properties of the algorithm are studied and it is shown that under some mild assumptions, the scheme enjoys monotonicity, asymptotic optimality and strong convergence to a point that lies in the consensus subspace. Finally, numerical examples verify the enhanced performance obtained by the proposed scheme compared to other algorithms, which have been developed in the context of sparsity–aware adaptive learning.

# 5  Dimensionality Reduction in Distributed Adaptive Learning via Krylov Subspaces

In this section, the problem of dimensionality reduction in adaptive distributed learning is studied. As in the previous sections, the algorithm, to be presented here, is based on the APSM algorithmic family. At each time instant and at each node of the network, a hyperslab is constructed based on the received measurements and this defines the region in which the solution is searched for. Moreover, in order to reduce the number of transmitted coefficients, which is dictated by the dimension of the unknown vector, we seek
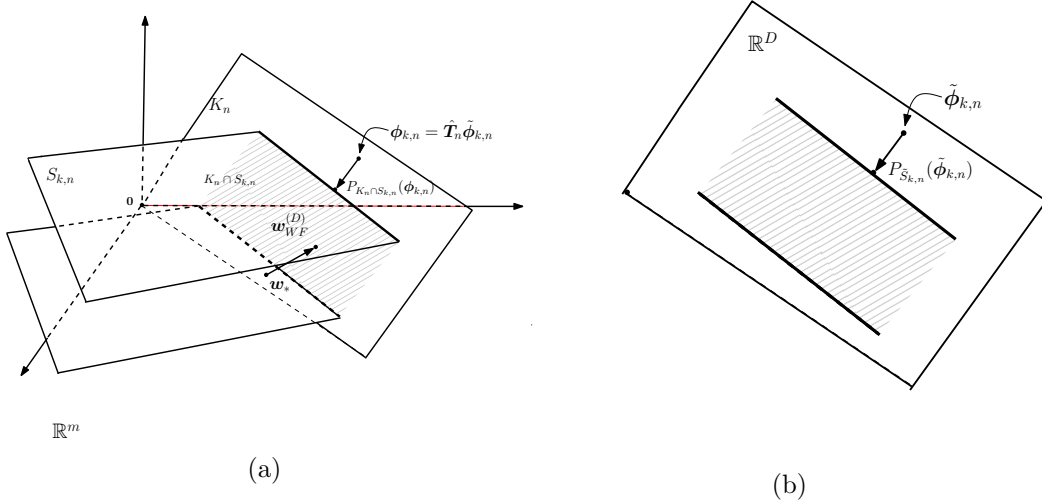
Figure 2: (a) Geometrical illustration of the algorithm for $q = 1$. The aggregate $\boldsymbol{\phi}_{k,n}$, which belongs in the subspace, is projected onto the intersection of the subspace and the hyperslab, generated by the measurement data. (b) The algorithmic scheme in the reduced dimension space, i.e., $\mathbb{R}^D$.

for possible solutions in a subspace of lower dimensionality; the technique will be developed around the Krylov subspace rationale. Our goal is to find a point that belongs to the intersection of this infinite number of hyperslabs and the respective Krylov subspaces. This is achieved via a sequence of projections onto the property sets as well as the Krylov subspaces. The proposed schemes are brought in a decentralized form by adopting the combine-adapt cooperation strategy among the nodes.

The steps of the algorithm can be encoded in the following formula:

$$\tilde{\boldsymbol{w}}_{k,n+1} = \tilde{\boldsymbol{\phi}}_{k,n} + \tilde{\mu}_{k,n} \left( \sum_{j \in \mathcal{J}} \omega_{k,j} P_{\tilde{S}_{k,j}}(\tilde{\boldsymbol{\phi}}_{k,n}) - \tilde{\boldsymbol{\phi}}_{k,n} \right), \tag{10}$$

where the vectors $tilde\boldsymbol{w}_{k,n+1}$, $\tilde{\boldsymbol{\phi}}_{k,n+1}$ belong to the reduced dimension space and the reduced dimension hyperslab is given by: $\tilde{S}_{k,n} := \{ \tilde{\boldsymbol{w}} \in \mathbb{R}^D : |d_{k,n} - \boldsymbol{u}_{k,n}^T \hat{\boldsymbol{T}}_n \tilde{\boldsymbol{w}}| \le \epsilon_k, \}$ where $\hat{\boldsymbol{T}}_n$ is a matrix, the columns of which, span the Krylov subspace. The geometrical interpretation of the algorithm is illustrated in Fig. 2

As in the previously derived schemes, the theoretical properties of the algorithm are studied and it is shown that the scheme enjoys monotonicity, asymptotic optimality and strong convergence to a point that lies in the intersection of the consensus subspace with the Krylov Subspace. Finally, numerical examples verify that the proposed scheme provides a good

trade-off between the number of transmitted coefficients and the respective performance.

# 6    List of Publications

**Journal Publications**

- S. Chouvardas, K. Slavakis, and S. Theodoridis. Adaptive robust distributed learning in diffusion sensor networks. *IEEE Transactions on Signal Processing*, 59(10):4692–4707, 2011.

- S. Chouvardas, K. Slavakis, Y. Kopsinis, and S. Theodoridis. A sparsity promoting adaptive algorithm for distributed learning. *IEEE Transactions on Signal Processing*, 60(10):5412 –5425, Oct. 2012.

- S. Chouvardas, K. Slavakis, and S. Theodoridis. Trading off complexity with communication costs in distributed adaptive learning via Krylov subspaces for dimensionality reduction. *IEEE Journal of Selected Topics in Signal Processing*, 7(2):257–273, 2013.

- Symeon Chouvardas, Konstantinos Slavakis, Sergios Theodoridis, and Isao Yamada. Stochastic analysis of hyperslab–based adaptive projected subgradient method under bounded noise. *To appear in IEEE Signal Processing Letters*, 2013.

**Conference Publications**

- S. Chouvardas, K. Slavakis, and S. Theodoridis. A novel adaptive algorithm for diffusion networks using projections onto hyperslabs. In *CIP*, pages 393–398. IEEE, 2010 **(best student paper award)**.

- Symeon Chouvardas, Konstantinos Slavakis, and Sergios Theodoridis. Trading off communications bandwidth with accuracy in adaptive diffusion networks. In *ICASSP*, pages 2048–2051. IEEE, 2011.

- Symeon Chouvardas, Konstantinos Slavakis, Yannis Kopsinis, and Sergios Theodoridis. Sparsity-promoting adaptive algorithm for distributed learning in diffusion networks. In *Proceedings of the 20th European Signal Processing Conference (EUSIPCO)*, pages 1084–1088. IEEE, 2012.

- Symeon Chouvardas, Gerasimos Mileounis, Nicholaos Kalouptsidis, and Sergios Theodoridis. A greedy sparsity–promoting LMS for distributed

adaptive learning in diffusion networks. In *Proc. of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2013.

- Symeon Chouvardas, Gerasimos Mileounis, Nicholaos Kalouptsidis, and Sergios Theodoridis. Training-Based and Blind Algorithms for Sparsity-Aware Distributed Learning. In *Proceedings of the European Signal Processing Conference (EUSIPCO)*, 2013.

# References

[1] Ian F Akyildiz, Weilian Su, Yogesh Sankarasubramaniam, and Erdal Cayirci. A survey on sensor networks. *IEEE Communications Magazine*, 40(8):102–114, 2002.

[2] Jacob Benesty and Steven L. Gay. An improved pnlms algorithm. In *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 2, pages 1881–1884, may 2002.

[3] Doron Blatt and Alfred O Hero. Energy-based sensor network source localization via projection onto convex sets. *Signal Processing, IEEE Transactions on*, 54(9):3614–3619, 2006.

[4] E.J. Candes, M.B. Wakin, and S.P. Boyd. Enhancing sparsity by reweighted $\ell_1$ minimization. *Journal of Fourier Analysis and Applications*, 14(5):877–905, 2008.

[5] F.S. Cattivelli and A.H. Sayed. Diffusion lms strategies for distributed estimation. *IEEE Transactions on Signal Processing*, 58(3):1035–1048, 2010.

[6] Chee-Yee Chong and Srikanta P Kumar. Sensor networks: evolution, opportunities, and challenges. *Proceedings of the IEEE*, 91(8):1247–1256, 2003.

[7] Deborah Estrin, Lewis Girod, Greg Pottie, and Mani Srivastava. Instrumenting the world with wireless sensor networks. In *Acoustics, Speech, and Signal Processing, 2001. Proceedings.(ICASSP'01). 2001 IEEE International Conference on*, volume 4, pages 2033–2036. IEEE, 2001.

[8] PJ Huber and EM Ronchetti. *Robust Statistics*. New York: Willey, 2009.

[9] Y. Kopsinis, K. Slavakis, and S. Theodoridis. Online sparse system identification and signal reconstruction using projections onto weighted $\ell_1$ balls. *IEEE Transactions on Signal Processing*, 59(3):936 –952, 2011.

[10] Jinchao Li and Ali H Sayed. Modeling bee swarming behavior through diffusion adaptation with asymmetric information sharing. *EURASIP Journal on Advances in Signal Processing*, 18, 2012.

[11] Joel B Predd, SB Kulkarni, and H Vincent Poor. Distributed learning in wireless sensor networks. *Signal Processing Magazine, IEEE*, 23(4):56–69, 2006.

[12] M. Rabbat and R. Nowak. Distributed optimization in sensor networks. In *Proceedings of the 3rd international symposium on Information processing in sensor networks*, pages 20–27. ACM New York, NY, USA, 2004.

[13] Sheng-Yuan Tu and Ali H Sayed. Foraging behavior of fish schools via diffusion adaptation. In *CIP*, pages 63–68. IEEE, 2010.

[14] M. Yukawa and I. Yamada. Set-theoretic adaptive filtering based on data-driven sparsification. *International Journal of Adaptive Control and Signal Processing*, 25:707–722.

[15] M. Yukawa and I. Yamada. A unified view of adaptive variable-metric projection algorithms. *EURASIP Journal on Advances in Signal Processing*, 2009:2–2, 2009.

# Research in parametric optical amplification and injection locking focused on high bit rate optical communication systems applications

Markos Alexandros Fragkos*

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

**Abstract.** In this dissertation, the properties of the all-optical processes of phase sensitive parametric amplification and injection locking, are thoroughly investigated, through numerical simulations and experiments, in order to design integrated optical devices that can improve the performance of coherenct optical communications. Exploiting the phase noise suppression that the phase sensitive parametric amplifier can provide, a novel 40 Gb/s RZ-DPSK regenerator based on the single pump topology is proposed, which adops a realistic solution for the all-optical generation of an idler wave identical to the signal, dealing with the unresolved problem of the practical implementation of the specific devices. Taking into consideration the amplitude noise suppression and the phase replication properties of the injection locked semiconductor laser, we propose the use of the specific device as an alternative solution for (D)PSK and (D)QPSK signal regeneration. From the above properties, the injection locked semiconductor laser is also proposed as an additional unit in a DPSK/ASK receiver in order to provide better discrimination of the two different data streams and improve the performance of the specific modulation format.

**Keywords:** Parametric amplification, injection locking, coherent communication systems, all-optical regeneration, orthogonal modulation formats.

## 1 Dissertation Summary

Coherent optical communication systems that rely on phase modulation formats such as differential phase-shift keying (DPSK), quadrature phase-shift keying (QPSK) etc., currently emerge as an alternative to on–off keying (OOK) for long-haul transmission. The superiority of the coherent systems is based on three fundamental advantages: The first one originates from the balanced detection of the phase modulated signals which offers a 3-dB improvement in receiver sensitivity over direct detection of OOK. The second one is their enhanced tolerance to fiber nonlinearities and basically to interchannel cross-phase modulation (XPM). And the third one is the plethora of the advanced multi-level phase modulation formats providing enhanced bandwidth efficiency. Nonetheless, nonlinear phase noise is the major limitation for

*Dissertation Advisor: Dimitrios Syvridis, Professor

the performance of these systems requiring the periodic regeneration of the signal for long-haul transmissions. The potential for commercial deployment of phase modulation formats has triggered worldwide the research activity on the optical manipulation of signals modulated in phase.

The all-optical regeneration of phase modulated signals is carried out by periodic suppression of the amplitude noise, which limits the build up of the non-linear phase noise, or by direct phase noise suppression.

In this dissertation we thoroughly investigate the processes of phase sensitive parametric amplification, which offers direct phase noise suppression, and injection locking, which offers amplitude noise suppression and phase replication of the injected signal.

In the paragraphs below we present the regenerator architectures that have been emerged from the current investigation and provide performance improvement of coherent optical communication systems.

## 2  Phase sensitive amplification

The FWM is a non-linear process that takes place inside non-linear media such as optical fibers or semiconductor devices (e.g. lasers, optical amplifiers) and depends on the third order susceptability of the medium. In this process, two propagating waves ($\omega_1$, $\omega_2$), inside the medium, create a perturbation to the optical power dependent refractive index which then modulates the phase of each wave with a characteristic frequency that equals to the beating frequency of the two waves, $\omega_{mod} = \omega_2 - \omega_1$. As a result, each wave develops side-bands with spectral distance from the central frequency equal to $\omega_{mod}$, which essentially leads to the generation of two new waves ($\omega_3 = \omega_1 - \omega_{mod}$, $\omega_4 = \omega_2 + \omega_{mod}$).
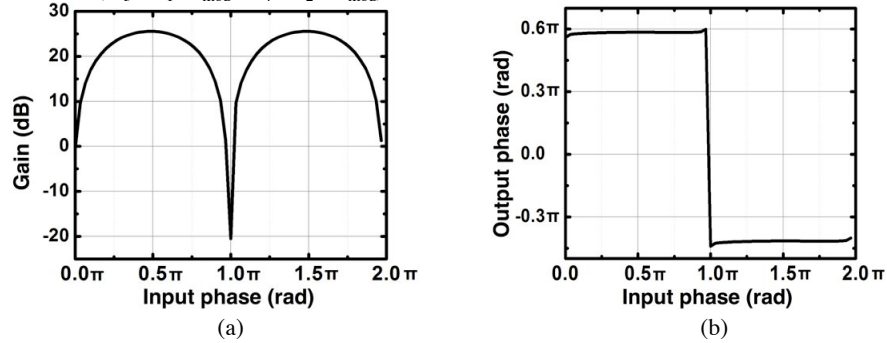


**Fig. 1.** Gain (a) and phase response (b) of a PSA as a function of the input signal's phase.

If four waves with the above spectral distances are being launched at the input of the non-linear medium, then depending on the total phase difference ($\Delta\phi$) between the four waves, energy can be trasnferred from frequencies $\omega_1$ and $\omega_2$ to $\omega_3$ and $\omega_4$, if $\Delta\phi = \pi/2$ (parametric amplification), or vice versa, if $\Delta\phi = -\pi/2$ (parametric absorption).

The parametric amplifiers are based on two main topologies, single and dual pump, depending on the number of strong waves that are used in order to amplify the weak

waves. The weak waves are called signal and idler. The presence or the absence of the idler wave at the input of the parametric amplifier determines the phase sensitive or phase insensitive gain of the amplifier, respectively.

Phase sensitive amplification exhibits a very interesting property. If the phase characteristics of the idler and signal waves are identical, then the gain curve of the amplifier as a function of signal's phase presents periodicity equal to π, while the phase response of the signal at the output of the amplifier is a step function with periodicity and step both equal to π. These properties make PSA the ideal device for all-optical regeneration of phase modulated signals.

## 2.1 Regenerative performance of a single pump PSA

In order to evaluate the noise suppression characteristics of the PSA, we investigate, through numerical analysis, its behavior to sinusoidal perturbations of the intensity and phase of the input signal. Depending on the power of the input signal, the PSA operates in two different regimes. For low input powers, PSA operates in linear regime in which the amplifier exhibits strong phase noise suppression ($\approx 18$ dB) and high phase noise to amplitude conversion ($\approx 26$ dB) . For higher input powers, PSA operates in non-linear regime in which the phase noise suppression is degraded ($\approx 12$ dB) but strong amplitude noise suppression ($< 12$ dB) is observed while the phase noise to amplitude conversion is strongly suppressed ($\approx 20$ dB). The behavior of the PSA in these two regimes is depicted in the diagrams of figure 2.
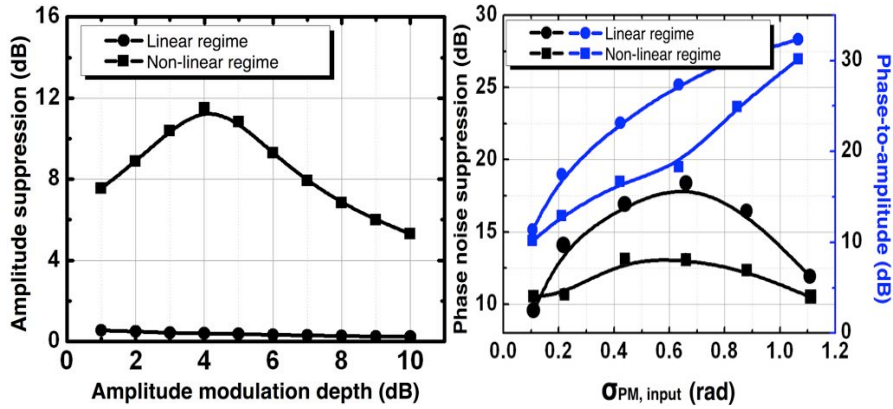


**Fig. 2.** PSA response to sinusoidal perturbation

Subsequently, we investigate the regenerative capabilities of the single pump PSA using RZ-DPSK signals. In order to evaluate its performance, we compare the quality of the RZ-DPSK signal, which is quantified by the Q factor, at the output of three different long-haul transmission links. The first link is a typical transmission link without the use of any regeneration device, consisting of multiple spans of 80 km of dispersion compensated transmission and amplification of the signal at the end of each span. The second link is based on a popular topology which uses recurrent units of saturated PIAs placed every 480 km along the transmission link. Finally, the third

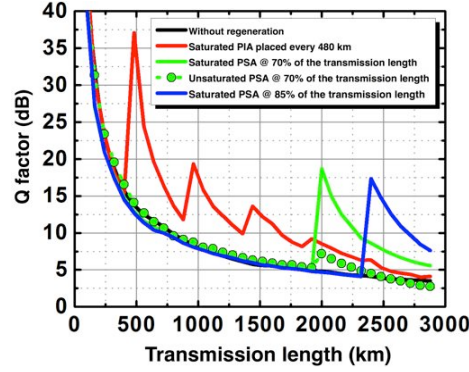link uses a PSA as a regenerative unit which is placed in an intermediate point of the transmission line.



**Fig. 3.** Performance comparison between three different optical transmission links for a 40 Gb/s RZ-DPSK transmitted signal.

As can be seen in the figure above, the saturated (non-linear regime) PSA exhibits remarkable regenerative performance due to its strong amplitude and phase noise suppression, improving the quality of the signal up to 10 dB higher than the unsaturated PSA and the multiple saturated PIAs.

## 2.2 RZ-DPSK regenerator based on single pump PSA

As we mentioned before, the idler wave at the input of the amplifier has to be identical to the signal in order to ensure the regenerative capabilities of the PSA. This constitutes the most significant obstacle for the implementation of this device and till now no realistic solution has been proposed, at least for the topology of the single pump PSA. In this dissertation we propose a novel architecture of a RZ-DPSK regenerator based on single pump PSA, in which the all-optical generation of an identical to the signal idler wave is achieved. The topology of the proposed regenerator is depicted in the block diagram of figure 4 and consists of three main units. The first unit undertakes the all-optical generation of two identical and phase locked intensity modulated RZ waves at the wavelengths $\lambda_s$ and $\lambda_i$, using a PIA implemented in a highly non-linear fiber (HNLF). The RZ pulses have the same repetition rate as the received RZ-DPSK signal and $\lambda_s$ is the same as the received wavelength. The second unit is the most important unit of the regenerator as it undertakes the all-optical and coherent transfer of the phase information of the received signal to the phase of each one of the locally generated RZ waves. Firstly, a delay interferometer (DI) transfers the DPSK information of the received signal to its amplitude.
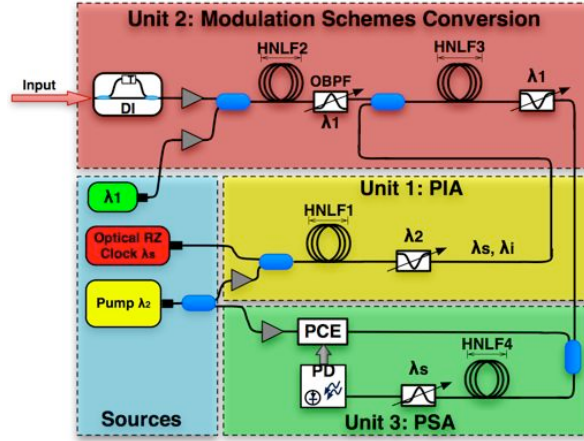
**Fig. 4.** RZ-DPSK regenerator based on single pump PSA. DI: delay interferometer, OBPF: optical band-pass filter, PCE: phase control element and PD: photodiode.

As a result, at the output of the DI we receive a RZ-OOK signal at the amplitude of which is imprinted the phase information and the total noise (amplitude and phase noise) of the received signal. The RZ-OOK signal is then converted into NRZ-OOK using an all-optical technique that is based on cross-phase modulation (XPM), which takes part in a second HNLF, and the final signal is transferred to the phase of $\lambda_s$ and $\lambda_i$ using the XPM process, which takes place in a third HNLF. The two identical RZ-DPSK signals ($\lambda_s$, $\lambda_i$) at the output of this unit are imported into the third unit in which, through phase sensitive amplification, the regeneration of $\lambda_s$ is achieved. The quality of the regenerated RZ-DPSK signal for the two different PSA operating regimes can be observed in the eye diagrams of figure 5.



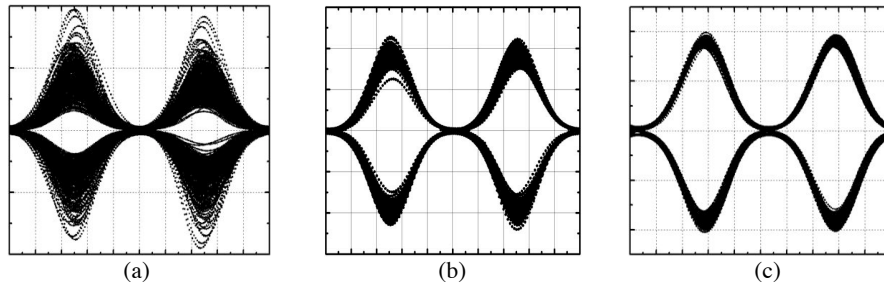(a)                              (b)                              (c)

**Fig. 5.** 40 Gb/s RZ-DPSK signal before (a) and after regeneration for PSA operating in linear (b) and non-linear (c) regime.

As it was expected from the analysis of the previous paragraph, the proposed regenerator exhibits better performance when the PSA operates in the non-linear regime improving the optical SNR of the received signal up to 28 dB.

# 3  Injection locking

Injection locking refers to the frequency effects that can occur when a harmonic oscillator (slave) is disturbed by a second oscillator (master) operating at a nearby frequency. When the coupling is strong enough and the frequencies near enough, the master oscillator can capture the slave oscillator, causing it to have essentially identical frequency as the master. This process can occur in a variety of oscillators in which the semiconductor lasers are included. In semiconductor lasers, the strength of the coupling is proportional to the injection ratio, R, which is given by the ratio of the injected power of master laser to the emitted power of the free running slave laser. As injection ratio increases the locking range of slave laser broadens enabling it to follow higher frequency perturbations of the injected signal. Through injection locking, slave laser acquires a variety of characterics such as the ability of replicating the phase information and suppress the amplitude noise of the injected signal, the enhancement of its modulation frequency, the reduction of its relative intensity noise, the suppression of its linewidth and the suppression of its side modes. Taking into consideration the above capabilities of the injection locked laser we propose its use as a cost effective, energy efficient and low complexity regenerator for phase modulated signals.

## 3.1  Regenerative performance of a single mode semiconductor laser

Using the modified rate equations appeared in [??] we modeled the master-slave system and investigated the phase replication and amplitude noise suppression capabilities of a semiconductor single mode injection locked laser for different phase modulation schemes and bit rates of the injected signal. The injection locked laser exhibits remarkable regenerative performance demonstrating modulation scheme transparency, supporting QPSK signals, high bit rate data replication, supporting up to 25 Gbaud/s, and strong amplitude noise suppression reaching up to 10 dB. In the polar diagrams of figure 6 is depicted the quality improvement that the injection locked laser can provide to a 25 Gb/s PSK and 50 Gb/s QPSK injected signal degraded by amplified spontaneous emission (ASE) noise.



**Fig. 6.** Polar diagrams of PSK (a) and QPSK (b) regenerated signals of 25 Gbaud/s. Black points correspond to master laser complex optical field and red points correspond to slave laser complex field.

To verify the performance of the proposed regenerator we realized the experimental set up of figure 7 and investigated the behavior of slave laser for 10 Gb/s DPSK injected signal degraded by two different noise cases; ASE noise and sinusoidal amplitude perturbations.



**Fig. 7.** ... n lock ... tor, IM: int ... : photo ... nce generat ...

In ... ed sig ... uses and in ... nding ... as a functio ... t of th ...



**Fig. 8.** Eye diagrams of 10 Gb/s DPSK signal before and after regeneration. Figures (a) and (c) correspond to ASE degraded signal and (b) and (d) correspond to amplitude noise degraded signal.

**Fig. 9.** (a) BER measurement as a function of the receiver power for ASE degraded signal of OSNR equal to 30 dB. (b) BER measurement as a function of the receiver power for a signal degraded only by strong amplitude perturbation of 1 GHz.

Due to phase to amplitude noise conversion in the case of ASE degraded injected signal, the power penalty that the regenerator provides is limited to 1.5 dB. On the other hand, if the injected signal is degraded only by amplitude noise then the power penalty increases to the remarkable value of 11 dB demonstrating error free data recovery at the receiver.

### 3.2 Regenerative performance of a Fabry-Perot semiconductor laser

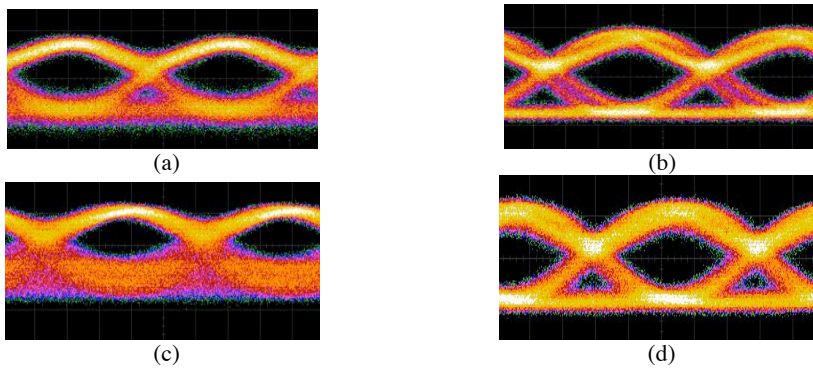Using the same experimental set up as before (fig. 7), we replaced the single mode semiconductor laser with a Fabry-Perot (FP) laser in order to evaluate its capability of single mode emission, that injection locking can provide to it through side mode suppression, and its regenerative performance. As can be seen in the figures below, the injection locked FP laser demonstrates single mode operation in a bandwidth of 16.4 nm arround its central emission wavelength, with side mode suppression ratio (SMSR) equal to 45 dB, for an injection ratio of -11 dB.



**Fig. 10.** Single mode operation of the injection locked Fabry-Perot laser for three representative emission modes.

In figure 11, the eye diagrams of the degraded by amplitude noise input and the regenerated output differentially phase shift keying (DPSK) signal are depicted for the above three different wavelengths (1538.89 nm, 1549.09 nm and 1555.29 nm). The figure shows clearly that FP laser is capable of suppressing the unwanted amplitude noise within a wavelength band of 16.4 nm. Outside this region, FP laser is able to lock but the injection level needed for single-mode operation is extremely high to allow limiting amplification and noise suppression.



(a)

(b)

(c)

(d)

(e)

(f)

**Fig. 11.** 10-Gb/s DPSK eye diagrams of the input signal degraded by ASE noise at 1538.89 nm (a), 1549.09 nm (c) and 1555.29 nm (e) and the corresponding regenerated output (b, d, f).

Figure 12 shows a series of BER measurements conducted for input OSNR equal to 23 dB, which further prove the remarkable amplitude limitation properties of the proposed injection locked laser amplifier. The amplitude noise level was adjusted appropriately so as the input signal to exhibit constant BER for every wavelength inside the investigated spectral range. The FP laser provides a reduction of 12 dB in the required receiving power for the achievement of BER performance equal to $10^{-3}$.

**Fig. 12.** BER performance as a function of the receiveing power for degraded master signal (black trace) and the corresponding regenerated signal at the output of the slave at 1538.89 nm (red trace), 1549.09 nm (green trace) and 1555.29 nm (blue trace).

### 3.3 Alternative application of injection locked semiconductor laser

Apart from its regenerative use, the injection locked laser can find a number of alternative applications providing quality improvement and capacity enhancement of future transmissions. In this dissertation, we propose a modified receiver based on the injection locking technique for the performance improvement of the novel orthogonal modulation scheme DPSK/ASK.

DPSK/ASK transmitter relies on the simultaneous modulation of both phase and amplitude of the signal with two different data streams, which doubles the bandwidth efficiency of the signal. Once received, the DPSK/ASK signal is driven through a 3 dB coupler into a direct detection receiver and a DPSK receiver for the demodulation of each data stream. Due to amplitude fluctuations at the DPSK receiver, caused by the remaining ASK information, the initial ASK extinction ratio must be kept at low values (<3dB) in order to ensure high quality of the demodulated DPSK signal. However, if we insert an injection locking unit at the input of DPSK receiver, the unwanted amplitude fluctuations can be suppressed improving DPSK quality and enabling the usage of higher ER of the ASK signal providing total performance improvement of the DPSK/ASK tranceiver. The experimental set up of the proposed modified receiver is depicted in the block diagram of figure 13 and its performance compared to the typical DPSK/ASK receiver is depicted in figure 14.

**Fig. 13.** Block diagram of the experimental setup. PM:phase modulator, IM: intenesity modulator, PRBS: pseudorandom binary sequence, EDFA1: boost amplifier, EDFA2, pre-amplifier, EDFA3: saturated amplifier, EDFA4: injection ratio controller, VA: variable attenuator, PD: photodiode, OBPF: optical band-pass filter, PC: polarization controller, DI: delay interferometer, OI: optical isolator.



**Fig. 14.** BER measurements as a function of ER for the 10-Gb/s ASK signal after direct demodulation (red squares) and the 10-Gb/s DPSK signal demodulated with (black circles) and without injection locking (blue triangles).

The proposed receiver allows the usage of ER up to 8 dB maintaining the high quality of the ASK and DPSK data streams for longer transmission lengths. Numerical simulations of the above transeiver demonstrated the error free detection of 10 WDM DPSK/ASK channels for transmissions up to 800 km for 20 Gb/s/channel and up to 400 km for 50 Gb/s/channel.

# 4  Conclusions

In this dissertation, three novel device architecture for all-optical processing have been proposed providing performance improvement of transmissions based on high efficiency modulation formats.

A RZ-DPSK signal all-optical regenerator based on single pump PSA is for the first time proposed, employing a realistic solution for the all-optical generation of an identical to the received signal idler wave. The specific regenerator can handle RZ-DPSK signals up to 40 Gb/s providing SNR improvement of up to 28 dB. The proposed regenerator can also be implemented using semiconductor optical amplifiers leading to footprint minimization, lower energy consumption and lower complexity.

The injection locked laser is for the first time proposed as an alternative all-optical regenerator for phase modulated signals. The specific regenerator can be placed along the optical path or at the input of the receiver improving the quality of the signal and reducing the required power for a given BER. It can also be placed at the output of a PSA operating in the linear regime suppressing the amplitude noise generated by the specific amplifier and leading to total noise suppression. The injection locked laser regenerator provides modulation format and wavelength transparency and can handle signals up to 25 Gbaud/s.

Finally, the injection locked laser is proposed as an additional unit at the input of the DPSK part of the DPSK/ASK receiver allowing the usage of higher ER of the ASK data stream, reaching the value of 8 dB, improving the total performance of the DPSK/ASK signal and therefore enabling its use as a cost effective alternative for channel capacity enhancement for access and metro networks.

## References

1. Fragkos, A. Bogris, and D. Syvridis, "All-Optical Regeneration Based on Phase-Sensitive Nondegenerate Four-Wave Mixing in Optical Fibers," IEEE Photon. Technol. Lett., vol. 22, no. 24, pp. 1826-1828, 2010

2. A. Fragkos, A. Bogris, D. Syvridis, and R. Phelan, "Amplitude Noise Limiting Amplifier for Phase Encoded Signals Using Injection Locking in Semiconductor Lasers," J. Lightw. Technol., vol. 30, no. 5, pp. 764-771, 2012

3. A. Fragkos, A. Bogris, D. Syvridis, and R. Phelan, "Colorless Regenerative Amplification of Constant Envelope Phase-Modulated Optical Signals Based on Injection-Locked Fabry–Pérot Lasers,"IEEE Photon. Technol. Lett., vol. 24, no. 1, pp. 28-30, 2012

4. A. Fragkos, A. Bogris, and D. Syvridis, "Efficient Orthogonal Modulation Enabled by Injection Locked Limiting Amplifiers," IEEE Photon. Technol. Lett., vol. 25, no. 7, pp. 667-670, 2013

5. A. Fragkos, A. Bogris, and D. Syvridis, "Black-box Optical Regenerator exploiting Non-Degenerate Phase-Sensitive Amplification," ECOC, P1.03, September (2010)

6. A. Fragkos, A. Bogris, D. Syvridis, R. Phelan, J. O' Carroll, B. Kelly, and J. O' Gorman "Amplitude Regeneration of Phase Encoded Signals Using Injection Locking in Semiconductor Lasers," OSA/OFC/NFOEC, OWG1, March (2011)

7. A. Fragkos, A. Bogris, and D. Syvridis, "Spectrally efficient and High Extinction Ratio DPSK/ASK Orthogonal Modulation Schemes Based on Injection Locking Limiting Amplifiers," ECOC, P3.04, September (2012)

# Changing Representation of Curves and Surfaces: Exact and Approximate Methods.

Tatjana Kalinka[*]

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

**Abstract.** In this thesis we explore recent methods for computing the Newton polytope of the implicit equation and study their applicability to the representation change from the parametric form to implicit. Computing a (super)set of the monomials appearing in the implicit equation allows us to determine the interpolation space. Following this phase we implement interpolation by exact or numeric linear algebra (applying singular value decomposition). We evaluate the monomials at the points, most suitable for the task, thus building a numeric matrix, ideally of corank 1, whose kernel vector contains the coefficients of the implicit equation. We propose techniques for handling the case of higher corank. This yields an efficient, output-sensitive algorithm for computing the implicit equation. The method can be applied to polynomial or rational parameterizations of planar curves or (hyper)surfaces of any dimension including parameterizations with base points. Moreover, this technique can be used for problems such as the computation of the discriminant of a multivariate polynomial or the resultant of a system of multivariate polynomials.

**Keywords:** implicitization, interpolation, Newton polytope, sparse resultant, linear algebra.

## 1 Introduction

The modern CAGD and CAM systems operate with several different representations of geometric objects, where each is more suitable for some applications. For instance, parametric and implicit representations have complementary features: with parametrization it is easy to obtain points on the geometric object while the implicit equation allows to check quickly if a given point is inside or outside a given object. Hence the need for the robust methods to change between the two representations.

In this thesis we focus on implicitization, i.e. process of changing the representation of a geometric object from parametric to implicit (algebraic). The main objective of our work was exploring applicability of the recently developed method for computing Newton polytope of a resultant [1,2] to computing the implicit equation.

**Definition 1.** *Given a polynomial $f = \sum_a c_a t^a \in \mathbb{R}[t_1, \ldots, t_n]$, $t^a = t_1^{a_1} \cdots t_n^{a_n}, a \in \mathbb{N}^n, c_a \in \mathbb{R}$, its support is the set $\{a \in \mathbb{N}^n : c_a \neq 0\}$; its Newton polytope $N(f)$ is the convex hull of its support.*

Let us now define the problem formally. A *parametrization* of a geometric object of *co-dimension one*, in a space of dimension $n + 1$, can be described by parametric map:

$$f : \mathbb{R}^n \to \mathbb{R}^{n+1} : t = (t_1, \ldots, t_n) \mapsto x = (x_0, \ldots, x_n),$$

---

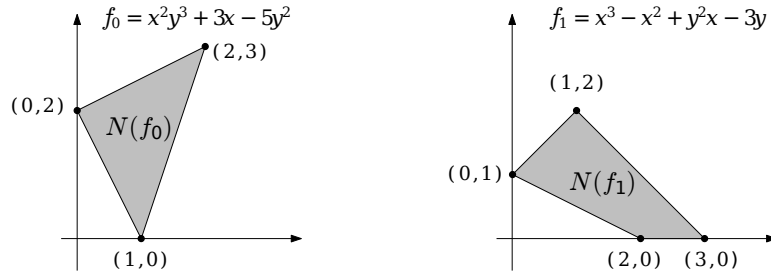[*] Dissertation Advisor: Ioannis Z. Emiris, Professor

**Fig. 1.** Newton polygons $N(f_i)$ of polynomials $f_i \in \mathbb{Z}[x, y]$.

where $t$ is the vector of parameters and $f := (f_0, \ldots, f_n)$ is a vector of continuous functions, including polynomial, rational, and trigonometric functions, also called coordinate functions. These are defined on some product of intervals $\Omega := \Omega_1 \times \cdots \times \Omega_n, \ \Omega_i \subseteq \mathbb{R}$.

**Definition 2.** *The implicitization problem asks for the smallest algebraic variety containing the image of the parametric map $f : t \mapsto f(t)$. This image is contained in the variety defined by the ideal of all polynomials $p(x_0, \ldots, x_n)$ s.t. $p(f_0(t), \ldots, f_n(t)) = 0$, for all $t$ in $\Omega$. We restrict ourselves to the case when this is a principal ideal, and we wish to compute its defining polynomial*

$$p(x_0, \ldots, x_n) = 0,$$

*given its Newton polytope, or a polytope that contains it.*

There have been numerous approaches to implicitization, including those based on Gröbner bases, resultants, residues, moving lines and surfaces, and $\mu$-bases. Our approach, presented in [3,4,5,6], follows the standard method of interpolating the unknown coefficients of the implicit polynomial given a superset of its monomials.

Using the recently developed method to determine potential monomials of the implicit equation we build a numeric matrix, ideally, of codimension 1, whose kernel yields (up to a nonzero scalar multiple) the coefficients corresponding to the predicted implicit support. This is a standard case of *sparse interpolation* of the polynomial from its values.

Since the kernel can be computed numerically, our approach also yields an approximate sparse implicitization method.

The kernel space of the numerical matrix may be of high dimension. We relate it to the geometry of the predicted support, which is a superset of the true implicit support. Another reason for obtaining a high-dimensional kernel is that the numeric evaluation of the support monomials may not be sufficiently generic.

Our implicitization method can be applied to planar curves, surfaces, or hypersurfaces of any dimension, given by a polynomial, rational or trigonometric parametrization, including those with base points. Moreover, our method can be used to compute discriminants of well-constrained systems as well as resultants, since the latter can be viewed as a special case of discriminants.

## 1.1   Prior work

Our implicitization algorithm is closely related to the interpolation-based method presented in [7]. We employ the same, most direct, method to reduce implicitization to linear algebra: construct

a $|S| \times |S|$ matrix $M$, indexed by monomials with exponents in $S$ (columns) and $|S|$ different values (rows) at which all monomials get evaluated, and compute kernel vector $p$ of $M$ containing coefficients of the implicit equation. This idea was used also in [8,9]; the approach was extended to an approximate implicitization as well.

Evaluation of the potential monomials at unitary $\tau \in (\mathbb{C}^*)^n$, one of the evaluation strategies examined in our work, was proposed in [10], Another approach, described in [11], is based on integration of matrix $M = SS^T$, over each parameter $t_1, \ldots, t_n$. Then, $p$ is in the kernel of $M$. This method covers a wide class of parametrizations, including polynomial, rational, and trigonometric representations, but the size of $M$ is large and matrix entries take big values, so it is difficult to control its numeric corank.

In practical applications of CAGD, precise implicitization often can be impossible or very expensive to obtain. Moreover, exact implicit equations usually are of high degree and contain unwanted branches and singularities. Hence the need for approximate implicitization proposed by T. Dokken [12,13]. The idea is to interpolate the coefficients using successively larger bounds on the total degree of the target polynomial, starting with a quite small support and extending it until a satisfying accuracy of the approximation is reached.

In order to determine the space of interpolation we have used the Newton polytope of the implicit equation, or *implicit polytope*.

There are several methods for computing the implicit polytope, such as those based on tropical geometry, or mixed fiber polytopes, for instance [14,15,9]. In this thesis the implicit polytope is computed from the Newton polytope of the sparse (or toric) resultant, or *resultant polytope*, of polynomials defined by the parametric equations. In particular, we use software tool `ResPol` [2] to compute the implicit support in general case and, in the case of curves, the method presented in [1] and implemented in `Maple`. We shall note, however, that our implicitization method does not depend on the approach used to compute the implicit polytope.

## 2   Algorithm and implementation

In this section we present our implicitization algorithm, discuss the importance of suitable evaluation points for building the matrix and give some details on the implementation.

The main steps of our algorithm are the following:

*Input:* Polynomial or rational parametrization $x_i = f_i(t_1, \ldots, t_n)$.

*Output:* Implicit polynomial $p(x_i)$ in the monomial basis in $\mathbb{N}^{n+1}$.

1. Determine (a polytope containing) the implicit polytope.
2. Compute all lattice points $S \subseteq \mathbb{N}^{n+1}$ in the polytope.
3. Repeat $\geq |S|$ times: Select value $\tau$ for $t$, evaluate $x_i(t)$, $i = 0, \ldots, n$, thus evaluating each monomial with exponent in $S$. This yields a matrix $M$.
4. Given matrix $M$, solve $Mp = 0$ for kernel $p$.
5. Let the kernel vector $p_i$ correspond to a polynomial of least total degree.
6. Return the primitive part of polynomial $p_i^\top \cdot m$, where $m$ is the set of monomials with exponent in $S$.

The complexity of our algorithm is $O(\mu |S|^2)$.

Let us describe the construction of matrix $M$ in Step 3.

Consider $S := \{s_1, \ldots, s_{|S|}\}$; each $s_j = (s_{j0}, \ldots, s_{jn})$ is an exponent of a (potential) monomial $m_j := x^{s_j} = x_0^{s_{j0}} \cdots x_n^{s_{jn}}$ of the implicit polynomial, where $x_i = f_i(t)/g_i(t)$. We evaluate $m_j$ at some $\tau_k$, $k = 1, \ldots, \mu$, avoiding values that make the denominators of the parametric expressions

close to 0, and obtain $m_j|_{t=\tau_k} := \prod_i \left(\frac{f_i(\tau_k)}{g_i(\tau_k)}\right)^{s_{ji}}$. Thus, we build an $\mu \times |S|$ matrix $M$ with rows indexed by $\tau$'s and columns indexed by $m_j$'s:

$$M = \begin{bmatrix} m_1|_{t=\tau_1} & \cdots & m_{|S|}|_{t=\tau_1} \\ \vdots & \cdots & \vdots \\ m_1|_{t=\tau_\mu} & \cdots & m_{|S|}|_{t=\tau_\mu} \end{bmatrix}$$

We compute the kernel of the matrix $M$ either symbolically or numerically, by applying singular value decomposition (SVD).

## 2.1  Choosing the evaluation points

Experiments with curves and surfaces in the monomial basis as well as in the Bernstein basis, show that when building the matrix $M$, it is important to choose $\tau$ values that are suitable for the specific instance.

Choosing $\tau$ for implicitization of classical algebraic curves and surfaces, we have experimented with

 – random integers in the range $-\mu^2 \ldots \mu^2$,
 – random rational numbers,
 – complex $\mu$-th roots of unity,
 – random complex numbers modulo 1.

Random integers offer the most numerically stable results, however with large matrices they result in fast growth of matrix entries. Random rational values have proved to be unreliable when implicitizing classical algebraic curves and surfaces, although complex values are numerically stable.

In the case of curves and surfaces in the Bernstein basis we have used evaluation by

 – random rational numbers,
 – uniformly distributed rational numbers,
 – complex roots of unity,
 – Chebyshev nodes in $[0, 1]$:

$$\tau = \frac{1}{2} + \frac{1}{2}\cos\left(\frac{2i-1}{2n}\pi\right), \ i = 1, \ \ldots, n.$$

While for classical algebraic curves and surfaces rational numbers led to a loss of numerical stability, here rational numbers chosen randomly in $[0, 1]$ provide the fastest results.

Our experiments affirm the results of [13], as the evaluation with Chebyshev nodes allows to minimize the approximation error in numerical computations. Complex roots of unity gave the slowest timings and introduced complex coefficients into the resulting approximate implicit equation.

## 2.2  Implementation

Our algorithm is implemented in `Maple`[1] and `SAGE`, based on the software for computing implicit polytopes [2], available as a C++ implementation[2]. The main functions are *imgen* (general

---

[1] http://ergawiki.di.uoa.gr/index.php/Implicitization
[2] http://sourceforge.net/projects/respol/files/

implicitization, applicable for curves, surfaces and hypersurfaces, requires parametric equations and predicted polytope vertices as an input) and *imcurve* (for curves only, support prediction is part of the routine).

For exact computations we prefer `Maple`, while for numerical ones `SAGE`. In our `Maple` implementation the computation of the lattice points in Step 2 is done, for up to four dimensions, by routines that utilize the `Maple` package `convex` [16], whereas our `SAGE` implementation uses its built-in functions for the same task. For higher dimensions we have employed the software package `Normaliz`.

When the kernel computation in Step 4 is done numerically, we build a rectangular overconstrained matrix $M$ in order to increase the numerical stability.

### 2.3 Accuracy of the approximate implicitization

It is important to estimate the numeric accuracy, or quality, of the result when solving numerically. We use the matrix condition number and the ratio between the two smallest singular values to evaluate the error in the coefficient vector computed by SVD.

We employ two measures to quantify the accuracy of approximate implicitization:

(a) Coefficient difference: measured as the Euclidean norm of the difference of the two coefficient vectors $V_{exact}, V_{app}$, obtained from exact and approximate implicitization, after padding with zero the entries of each vector which do not appear in the other.
(b) Evaluation norm: measured by considering the maximum norm of the approximate implicit equation when evaluated at a set of sampled points on the given parametric object.

Both accuracy measures feature in Table 2 where we compare running time and accuracy of our approximate implicitization against another method.

## 3 Results

In this section we discuss some of the key results of our work. First we prove the relation between the size of the predicted implicit polytope and the presence of the extraneous factors in the resulting expression. Next, we also talk about performance of the implementation of our method compared with others and present examples of alternative, non-geometrical application of the method.

### 3.1 The extraneous factors cases

By the construction of matrix $M$ using values $\tau$ that correspond to points on the parametric surface, we have the following:

**Lemma 1.** *Any polynomial in the basis of monomials indexing $M$, with coefficient vector in the kernel of $M$, is a multiple of the implicit polynomial $p$.*

The following theorem establishes the relation between the dimension of the kernel of $M$ and the accuracy of the predicted support. It remains valid even in the presence of base points. In fact, it also accounts for them since then, $P$ is expected to be much smaller of $Q$.

**Theorem 1.** *Let $P = N(p)$ be the implicit polytope and $Q$ the predicted polytope. Then, assuming $M$ has been built using sufficiently generic evaluation points, the dimension of its kernel space equals $\#\{m \in \mathbb{Z}^n : m + P \subseteq Q\} = \#\{m \in \mathbb{Z}^n : N(x^m \cdot p) \subseteq Q\}$.*

We assume genericity of the resultant whose symbolic coefficients are then specialized to the actual coefficients of the parametric equations. If this does not hold, then the actual implicit equation divides the specialized resultant.

In order to produce the exact implicit equation in the instance when the matrix $M$ has corank $> 1$ we propose the following:

– Reduce the predicted Newton polytope.
– Compute gcd of two or more polynomials corresponding to kernel vectors. In case of numeric solving approximate methods for computing the gcd can be applied.
– Apply factoring, then determine which of the factors vanishes when the $x_i$ variables are substituted by the parametric expressions.
– In practice, an actual implicit equation usually is present among the polynomials corresponding to kernel vectors. Hence the solution: sort the polynomials, return the one of the least degree.

*Example 1.* Consider its parametrization:

$$x_0 = \frac{2s}{1 + t^2 + s^2}, \ x_1 = \frac{2st}{1 + t^2 + s^2}, \ x_2 = \frac{-1 - t^2 + s^2}{1 + t^2 + s^2}.$$

Predicted implicit polytope has vertices: $(0,0,0)$, $(0,0,2)$, $(0,0,4)$, $(0,2,0)$, $(0,4,0)$, $(4,0,0)$. Implicit equation of the sphere being quadratic, here implicit polytope $P \subset Q$, where $Q$ is predicted polytope, which contains the actual implicit polytope. It contains 35 lattice points. We build $M$ of size $\mu \times 35$ ($\mu \geq 35$) of corank 10. The polynomials corresponding to the kernel vectors are:

$g_1 = -y^2 + y^2z^2 + y^4 + x^2y^2,$
$g_2 = -z^2 + z^4 + y^2z^2 + x^2z^2,$
$g_3 = -1 + z^2 + x^2 + y^2,$
$g_4 = -x + xz^2 + xy^2 + x^3,$
$g_5 = -yz + yz^3 + y^3z + x^2yz,$
$g_6 = -y + yz^2 + y^3 + x^2y,$
$g_7 = -xz + xz^3 + xy^2z + x^3z,$
$g_8 = -z + z^3 + y^2z + x^2z,$
$g_9 = -xy + xyz^2 + xy^3 + x^3y,$
$g_{10} = -1 + 2z^2 - z^4 + 2y^2 - 2y^2z^2 - y^4 + x^4.$

Computing the gcd of two randomly chosen polynomials yields, either the actual implicit equation $p = -1 + z^2 + x^2 + y^2$, or a multiple of $p$ of degree 3.

Let us have a closer look at the numeric solving in the case of $\dim(\text{kernel}(M)) = 10$. Applying SVD in we obtain approximate results, i.e. polynomials with non-integer coefficients. Computing the kernel of $M$ approximately yields polynomials with real coefficients.

The approximate gcd of the first two is:

$-0.9999998548199414 + 0.9999999857259533x^2 + 1.000000000052092y^2 + 1.000000000000000z^2$, which is accurate to 7 decimal digits.

## 3.2   Comparison to other methods

Here we report on a comparison of our method, implemented in `Maple`, against existing implicitization software. All the experiments mentioned have been performed on an Intel©Core2 Duo CPU, 2.20GHz, 3Gb memory, `Maple 14`.

Table 1 features the running times for implicitization of some examples of algebraic curves by different methods, all implemented in `Maple`. Namely our function *imcurve*, only for curves, that includes support prediction routine, $\mu$-bases method only for curves [17], and *Maple* function *Implicitize*, which employs integration of matrix $M$ over each parameter [11] and can be run in exact and numerical mode.

| Curve | degree | *Implicitize* exact | *Implicitize* numeric | Our software | $\mu$-bases |
|---|---|---|---|---|---|
| Trisectrix of Maclaurin | 3 | 1.92 | 0.064 | 0.02 | 0.016 |
| Folium of Descartes | 3 | 9.3 | 0.08 | 0.012 | 0.024 |
| Tricuspoid | 4 | 1.92 | 0.064 | 0.044 | 0.016 |
| Bean | 4 | 129.7 | 0.12 | 0.036 | 0.028 |
| Talbot's | 6 | 18.98 | 0.252 | 0.324 | 0.072 |
| Fifth heart | 8 | 799.74 | 0.44 | 0.104 | 0.08 |
| Ranunculoid | 12 | >3000 | 1.64 | 1.376 | 0.3 |

**Table 1.** Comparing runtimes (sec) of: Maple function *Implicitize* (exact and numeric), our method, and $\mu$-bases.

Of the three methods *Implicitize*, even in numerical mode, is the slowest, however the method has less restrictions on the parametrization accepting non-rational representations. Our method is faster than *Implicitize* but slower than the $\mu$-bases method.

While in case of curves our method may not be the best choice, experiments show that for geometric objects of higher degree and dimension it is competitive to the popular Gröbner bases method.

Consider Table 2, where we show the results of our experiments with surfaces. Here we compare our `Maple` implementation *imgen* against `Maple`'s native function *Implicitize* in numerical mode and implicitization using Gröbner bases in `Maple`.

The input consists of a family of classical algebraic surfaces, the so called Plücker's conoid. $x_0 = t$, $x_1 = s$, $x_2 = \frac{Re((t+I\cdot s)^a)}{|(t+I\cdot s)^a|}$. By choosing appropriate values of parameter $a = 2b$ we obtain rational parameterizations of the surfaces with desired total degree. While implicitization of the Plücker's conoid is trivial task, we have chosen this example to demonstrate robustness of our method when the properties of the surface family allow us to compute comparatively small implicit polytope. This is the reason our exact implicitization shows here better results that the Gröbner bases method. We should note that, compared with the `Implicitize` function, our approximate method is not only faster but also more precise (we use accuracy measures (a) and (b) as defined in [2.3]).

In general, the results of our experiments show that for low degree curves ($\leq 6$) or surfaces ($\leq 4$), Gröbner bases outperform our software. The situation is reversed for higher degree: for instance, the ranunculoid curve (degree 12) was computed in 1.3 sec. by our method and in 7.3 sec. using Gröbner bases. For the standard benchmark of the bicubic surface (degree 18) the timings are 42 min. and over 4 hours, respectively.

### 3.3 Non-geometrical applications

Our algorithm finds other applications besides the implicitization. Since the support prediction software `ResPol` actually computes a resultant support, its straightforward application is to

**Table 2.** Comparison of our method (exact and numerical) to `Maple`'s function *Implicitize()* and Gröbner bases. Runtimes are given in seconds.

| Surface degree | Our exact | Gröbner | Our numerical | | | Implicitize(numerical) | | |
|---|---|---|---|---|---|---|---|---|
| | runtime | | runtime | accuracy (a) | accuracy (b) | runtime | accuracy (a) | accuracy (b) |
| 3 | 0.016 | 0.031 | 0.031 | $10^{-15}$ | $9.07 \cdot 10^{-10}$ | 46.07 | $10^{-15}$ | $1.98 \cdot 10^{-9}$ |
| 5 | 0.016 | 0.046 | 0.032 | $10^{-10}$ | $3.57 \cdot 10^{-8}$ | 85.43 | $3.67 \cdot 10^{-7}$ | $6.83 \cdot 10^{-6}$ |
| 7 | 0.031 | 0.078 | 0.046 | $10^{-11}$ | $9.97 \cdot 10^{-8}$ | 359.49 | $9.06 \cdot 10^{-7}$ | $2.94 \cdot 10^{-4}$ |
| 9 | 0.046 | 0.078 | 0.063 | $10^{-10}$ | $1.35 \cdot 10^{-7}$ | 695.65 | $2.86 \cdot 10^{-6}$ | $7.55 \cdot 10^{-3}$ |
| 11 | 0.078 | 0.141 | 0.078 | $10^{-11}$ | $1.07 \cdot 10^{-6}$ | > 2000 | - | - |

reduce resultant computation to interpolation; this is also the premise of [18,19]. The main difference with interpolating the implicit equation is the absence of a parametric form, however, the latter can be derived using Horn-Kapranov parametrization [20], as demonstrated below.

*Example 2.* Let $f_0 = a_2 x^2 + a_1 x + a_0$, $f_1 = b_1 x^2 + b_0$, with supports $A_0 = \{2, 1, 0\}$, $A_1 = \{1, 0\}$. Their (Sylvester) resultant is a polynomial in $a_2, a_1, a_0, b_1, b_0$.

The algorithm in [2] computes its Newton polytope with vertices $(0, 2, 0, 1, 1)$, $(0, 0, 2, 2, 0)$, $(2, 0, 0, 0, 2)$; it contains 4 points, corresponding to 4 potential monomials $a_1^2 b_1 b_0$, $a_0^2 b_1^2$, $a_2 a_0 b_1 b_0$, $a_2^2 b_0^2$.

The Horn-Kapranov parametrization of the resultant yields: $a_2 = (2t_1 + t_2) t_3^2 t_4$, $a_1 = (-2t_1 - 2t_2) t_3 t_4$, $a_0 = t_2 t_4$, $b_1 = -t_1 t_3^2 t_5$, $b_0 = t_1 t_5$, where the $t_i$'s are parameters.

We substitute these expressions to the predicted monomials, $-t_1^2 t_3^4 t_5^2 (-2t_1 - 2t_2)^2 t_4^2$, $t_1^2 t_3^4 t_5^2 t_2^2 t_4^2$, $-t_1^2 t_3^4 t_5^2 t_2 t_4^2 (2t_1 + t_2)$, $(2t_1 + t_2)^2 t_3^4 t_4^2 t_1^2 t_5^2$, evaluate at 4 sufficiently random $t_i$'s, and obtain a matrix whose kernel vector $(1, 1, -2, 1)$ yields $\mathcal{R} = a_1^2 b_1 b_0 + a_0^2 b_1^2 - 2 a_2 a_0 b_1 b_0 + a_2^2 b_0^2$.

Another possible application is computing the discriminant of a multivariate polynomial.

Computation of the discriminant is a difficult problem, since explicit formulas only exist for low-degree uni-variate polynomials. We reduce discriminant computation to sparse implicitization.

**Definition 3.** *$A$-discriminant is an irreducible polynomial $D_A = D_A(c)$ with integer coefficients in the vector of coefficients $c = (c_a : a \in A)$, defined up to sign, which vanishes for each choice of $c$ for which $F_A$ and all $\partial F_A / \partial t_i$ have a common root in $(\mathbb{C} \backslash \{0\})^n$.*

Given $A$, we form the $(n + 1) \times m, m > n + 1$ integer matrix (also called $A$ by abuse of notation) whose first row consists of ones, and whose columns are given by the points $(1, a)$ for all $a \in A$.

Let $B = (b_{ij}) \in \mathbb{Z}^{n \times (m-n-1)}$ be a matrix whose column vectors are a basis of the integer kernel of matrix $A$. Then $B$ is of full rank. Since the first row of $A$ equals $(1, \ldots, 1)$, the column vectors of $B$ add up to 0.

We illustrate computation of $A$-discriminant by the following example from [5].

*Example 3.* Consider a generic polynomial of two variables of degree 3,

$$F_A(t_1, t_2) = c_1 t_1 + c_2 t_2 + c_3 t_1 t_2 + c_4 t_1^2 + c_5 t_1^3$$

where $A = \{[1, 0], [0, 1], [1, 1], [2, 0], [3, 0]\} \subset \mathbb{Z}^2$.
We build the matrix A:

$$A = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 0 & 1 & 2 & 3 \\ 0 & 1 & 1 & 0 & 0 \end{pmatrix}$$

then the matrix B is as follows:

$$B = \begin{pmatrix} -1 & -1 \\ 1 & 2 \\ -1 & -2 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}$$

Let $l_1 = -1 - s, l_2 = 1 + 2s, l_3 = -1 - 2s, l_4 = 1, l_5 = s$, then we have the parametrization

$$f_1 = \frac{l_2 l_4}{l_1 l_3} = \frac{1 + 2s}{(-1 - 2s)(-1 - s)}, f_2 = \frac{l_2^2 l_5}{l_1 l_3^2} = \frac{(1 + 2s)^2 s}{(-1 - s)(-1 - 2s)^2}$$

Support prediction yields 4 Newton polygon vertices: $[0,0]$, $[2,0]$, $[3,0]$, $[3,2]$. The Newton polygon has 7 lattice points. Applying our method, we obtain implicit equation $\Delta_B(x, y) = x - y - 1$.

We perform substitution following

$$\Delta_B(f_1, f_2) = D_A\left(1, 1, 1, \frac{l_2 l_4}{l_1 l_3}, \frac{l_2^2 l_5}{l_1 l_3^2}\right),$$

which gives us A-discriminant of $F_A$: $D_A(c) = c_2 c_3 c_4 - c_2^2 c_5 - c_1 c_3^2$.

## 4    Conclusions

We have developed an algorithm for computing implicit equations that combines linear algebra with promising support prediction methods. The method applies to polynomial, rational and trigonometric parameterizations of classical algebraic equations of curves and (hyper)surfaces. Moreover, it can be used for implicitization of geometric objects represented in NURBS form, after converting them to the monomial base. The method works even in the presence of base points, which are known to raise important issues for some other implicitization methods.

Our method has its limits: geometric objects have to be presented using the monomial basis and in the case of trigonometric parameterizations they have to be convertible to rational functions.

In some instances the polynomial computed using our algorithm contains an extraneous factor. We have analyzed such cases the and propose techniques for handling them.

While our implicitization algorithm was intended and applied in this work to implicitize curves and surfaces of codimension 1 only, interpolation can be sufficiently applied to implicitize space curves; this can be a challenge for future work.

Moreover, it is possible that our algorithm can be adapted in a way that, while bypassing actual computing of the equation, effectively provides an answer if the point belongs to the the geometric object. Our method represents an implicit (hyper)surface by a kernel vector. It is challenging to devise suitable CAGD algorithms that exploit this representation, for instance to compute surface-surface intersection, as in [12,21].

Yet another topic for future work is approximate implicitization of piecewise parametrized Bézier and NURBS curves and surfaces. In this thesis we have limited our experiments to implicitization of single patches of the objects represented in Bernstein basis, however the same interpolation principle can be applied to computing implicit equations of curve or surface splines. As

demonstrated in the thesis, the need for conversion from Bernstein basis to power basis presents a major drawback: in the NURBS format curves and surfaces are usually given by floating point coefficients and recalculating the parametrization in power basis furthers the precision loss.

# References

1. Emiris, I. Z., Konaxis, C.,Palios, L.: Computing the Newton polygon of the implicit equation. Mathematics in Computer Science, Special Issue on Computational Geometry and Computer-Aided Design, 4(1), 25–44 (2010)
2. Emiris, I. Z., Fisikopoulos, V., Konaxis, C., Peñaranda, L.: An output-sensitive algorithm for computing projections of resultant polytopes. In: Proceedings of the 2012 symposuim on Computational Geometry SoCG '12 New York, NY, USA: ACM. Final version to appear in IJCGA pp. 179–188 (2012)
3. Emiris, I. Z., Kalinka, T., Konaxis, C.: Implicitization of curves and surfaces using predicted support. In: Electr. Proc. Inter. Works. Symbolic-Numeric Computation San Jose, Calif. (2011)
4. Emiris, I. Z., Kalinka, T., Konaxis, C., Luu Ba, T.: Implicitization of curves and (hyper)surfaces using predicted support. Theoretical Computer Science (2012)
5. Emiris, I. Z., Kalinka, T., Konaxis, C., Luu Ba, T.: Sparse implicitization by interpolation: Characterizing non-exactness and an application to computing discriminants. Computer-Aided Design. 45(2), 252–261 Special Issue Conference on SPM (2013)
6. Emiris, I. Z., Kalinka, T., Konaxis, C.: Sparse implicitization via interpolation. To appear in SAGA Volume (Springer) (2013).
7. Emiris, I. Z., Kotsireas, I. S.: Implicit polynomial support optimized for sparseness. In: Proc. Intern. Conf. Computational science appl.: Part III Berlin: Springer. pp. 397–406 (2003)
8. Marco, A., Martínez, J.-J.: Implicitization of rational surfaces by means of polynomial interpolation. CAGD. 19, 327–344 (2002)
9. Sturmfels, B., Yu, J.: Tropical implicitization and mixed fiber polytopes. In: Software for Algebraic Geometry, volume 148, of IMA Volumes in Math. & its Applic. pp. 111–131 Springer New York (2008)
10. Sturmfels, B., Tevelev, J., Yu, J.: The Newton polytope of the implicit equation. Moscow Math. J. 7(2) (2007)
11. Corless, R. M., Giesbrecht, M., Kotsireas, I. S., Watt, S. M.: Numerical implicitization of parametric hypersurfaces with linear algebra. In Proc. AISC. pp. 174–183 (2000)
12. Dokken, T., Thomassen, J. B.: Overview of approximate implicitization. Topics in algebraic geometry and geometric modeling. 334, 169–184 (2003)
13. Barrowclough, O. J. D., Dokken, T.: Approximate implicitization of triangular Bézier surfaces. In Proceedings of the 26th Spring Conference on Computer Graphics SCCG New York, NY, USA. pp. 133–140 (2010)
14. D'Andrea, C., Sombra, M.: The Newton polygon of a rational plane curve. Math. in Computer Science. 4(1), 3–24 (2010)
15. Esterov, A., Khovanskiĭ, A.: Elimination theory and newton polytopes. arXiv:0611107[math] (2006)
16. Franz, M.: Convex: a maple package for convex geometry, version 1.1.3. Available at: http://www-math.uwo.ca (2009)
17. Busé, L.. Luu Ba, T.: Matrix-based implicit representations of algebraic curves and applications. Computer Aided Geometric Design. 27(9), 681–699 (2010)
18. Cueto, M. A., Dickenstein, A.: Some results on inhomogeneous discriminants. In Proc. XVI Latin Amer. Algebra Colloq., Bibl. Rev. Mat. Iberoamericana. arXiv:math/0610031v2 [math.AG] pp. 41–62 (2007)
19. Tanabé, S.: On Horn-Kapranov uniformisation of the discriminantal loci. Adv. Studies Pure Math. 46, 223–249 (2007)
20. Kapranov, M.: A characterization of A-discriminantal hypersurfaces in terms of the logarithmic Gauss map. Mathematische Annalen. 290, 277–285 (1991)
21. Dokken, T., Thomassen, J. B.: Weak approximate implicitization. In Proc. IEEE Intern. Conf. Shape Modeling Appl. p. 31 (2006)

# Automated Negotiations between Intelligent Entities Participating in Electronic Markets

Kostas Kolomvatsos[*]

Department of Informatics and Telecommunications
National and Kapodistrian University of Athens
15784, Ilissia, Athens, Greece
kostasks@di.uoa.gr

**Abstract.** Automated negotiations comprise an interesting research domain for many years. A scenario, mostly depicting real life negotiations, defines that entities act under no knowledge on the characteristics of the rest of them. This means that their behavior should incorporate mechanisms for handling uncertainty imposed by the lack of knowledge as well as intelligent methods for modelling every aspect of the discussed scenario. In this thesis, we adopt computational intelligence techniques in order to propose efficient mechanisms for the definition of the behavior of entities participating in Electronic Markets. We cover the entire framework defined in a marketplace by proposing methodologies for the definition of basic parameters together with decision making models. We take into consideration the uncertainty in such scenarios in combination with profit maximization. The proposed models are based on Fuzzy Logic, Swarm Intelligence, Optimal Stopping Theory and Machine Learning techniques. We describe methods for the selection of middle entities and products. We utilize Quality of Service parameters in order to increase the efficiency of the proposed models. We study negotiations between one buyer and one seller as well as concurrent negotiations between a buyer and multiple sellers.

**Keywords.** Negotiations, Fuzzy Logic, Prediction, Neural Networks, Swarm Intelligence, Optimal Stopping Theory

## 1 Dissertation Summary

Nowadays, users are confronted with a huge amount of information resources. Apart from information, users are able to find and purchase a very large number of products. Providers can find new ways to reach customers in an open and dynamic environment such as Web. However, a number of difficulties are present in purchasing products. The first is the number of product resources. It is out of human capabilities to find and navigate in a huge amount of Web stores. Moreover, it is very difficult for users: a) to collect the necessary information for various products, and b) to identi-

---

[*] Dissertation Advisor: E. Hadjiefthymiades, Associate Professor

fy providers' intentions. For example, customers cannot be aware of the providers' pricing strategies. Therefore, customers should spend a lot of time and effort in order to conclude successful transactions.

The solution to the above mentioned problems can be the combination of the intelligent agent technology with Electronic Markets (EMs). An Intelligent Agent (IA) is a software or hardware component capable of acting in order to accomplish the tasks delegated by its owner. EMs are virtual environments where a set of entities try to agree upon the exchange of goods. Usually, there are groups of market members such as: the buyers, the sellers and members that are in the middle between them helping in their tasks. Buyers aim to buy products while sellers offer a number of specific products. Middle entities deal with administration or mediation tasks.

In this dissertation, we focus on modelling the behavior of entities as well as on the negotiation between them. The interaction between autonomous entities for exchanging offers with the final objective of dealing for purchases can be defined as negotiation. Usually, in negotiations, entities are selfish and try to maximize their profit. Negotiations could be either single (bilateral) or concurrent (multi-lateral). Bilateral negotiations are related to the negotiation between a buyer and a seller. The negotiation could involve either a single issue (product characteristic) or multiple issues (e.g., price, delivery time, etc). Finally, a negotiation could involve a complete or an incomplete information setting related to the knowledge of the opponent characteristics (e.g., deadline, pricing strategy, etc).

A number of research efforts focus on the discussed domain. The majority of them are based on Game Theory while others adopt Fuzzy Logic. Both of them are used for defining various parameters of the negotiation. However, the proposed approaches have a number of disadvantages. For example, game theoretic models require the definition of players' strategies and, in many cases, assume common knowledge of some of the characteristics of players (e.g., deadlines distribution). Additionally, fuzzy logic schemes are based on a specific rule base that describes the actions followed at every round. The following list summarizes the use of fuzzy logic in negotiations:

- for evaluating the difference of issues values or the attributes of each offer or specific constraints in successive rounds of the negotiation.
- for predicting the reservation prices (e.g., upper – lower acceptable price) of the opponent.
- for deciding the appropriate action at every round of the negotiation.
- for evaluating the satisfaction level that an offer produces in the players' decision mechanism.

In this thesis, we deal with a number of open issues in negotiations. These issues are:

- The efficient definition of specific behavior parameters for each player.
- The definition of the equilibrium path under no knowledge on the players' characteristics.
- The definition of an efficient decision making mechanism to be used by the buyer and the seller.

- The definition of an efficient mechanism for selecting the appropriate products and middle entities.
- The definition of an efficient trust framework to be used in dynamic environments like EMs.

In the above listed issues, we propose specific methodologies and models. Computational Intelligence techniques can provide the basis for solving problems arising when entities interact with each other. A real life scenario defines that entities have no knowledge on the characteristics of other entities also involved in the EM setup. We adopt computational intelligence techniques in order to propose efficient mechanisms for the definition of the behavior of entities participating in EMs. We cover the entire framework defined in a marketplace by proposing methodologies for the definition of basic parameters together with decision making models at every step of the negotiation. We take into consideration the uncertainty in such scenarios in combination with profit maximization. We propose decision making models that are based on different aspects of the discussed scenario in order to reveal the optimal one. We cover the research gap by proposing an efficient decision making mechanism, for the buyer [3] and the seller side, based on fuzzy logic [9] and utilizing a number of parameters (instead of using a limited number as research efforts found in the literature).

We describe methods for the selection of middle entities [4] and products [5]. The proposed methods result the appropriate middle entity or product that best matches the buyers' needs. We utilize Quality of Service parameters in order to increase the efficiency of the proposed model. We study negotiations between one buyer and one seller as well as concurrent negotiations between a buyer and multiple sellers. In the first case, we rely on the game theory principles with the objective to provide a model that maximizes the expected profit. For the second case, we rely on Swarm Intelligence theory in order to have a framework where threads, used by the buyer, converge to the best solution (the best agreement) through a team work. Additionally, Optimal Stopping Theory gives us the tool to change the view of the problem. Based on Optimal Stopping Theory, we propose models trying to find the best time to take a decision instead of finding the best action as the response to the opponent move.

Finally, we propose a technique for defining the trust level of entities [11]. The reason is that our scenario involves an open and very dynamic environment like EMs. The trust level of an entity affects the decision taken by the rest of them for their involvement in negotiations with her. If the entity is very trusted, then the risk of negotiations with an unknown entity (we are not sure that the entity is going to offer what is promoting) is eliminated. From the above, we see that we try to cover all the aspects of negotiations starting from the selection of entities to negotiate with, to decision making mechanisms. We reveal the problems in the specific research area and propose specific solutions.

## 2    Results and Discussion

### 2.1    Negotiation Setup

At every round of the negotiation, the involved entities propose a specific price to the opponent. Should this price be accepted, the negotiation ends with an agreement and specific profit for both parties. The seller starts first and the buyer follows if the proposed offer is rejected. If a player is not satisfied with the offer then she has the right to reject it and issue a counter-proposal. If an agreement is reached then the negotiation ends with profit for both parties. A conflict leads to zero profit for both.

In this interaction, there are two factors that affect the decision making of the entities. The first factor is the seller's cost and the second one is the buyer's valuation about the product. The proposed offers have a lower limit defined by the cost (seller side). The buyer has a specific valuation about the product and is not willing to pay more than this value. Evidently, only in the case where the seller's cost is smaller than the buyer's valuation an agreement can be reached. However, the two players do not know if this pre-condition holds true. Finally, there is a specific time horizon for the negotiation [10]. The buyer has a specific deadline posed by her owner while the seller calculates her deadline as discussed in [7, 8]. If one of the deadlines expires and no agreement is reached till then the negotiation ends with a conflict.

The characteristics of the buyer are: the *valuation* about the product ($V$), the *discount factor* ($\delta_b$), the *utility function* ($U_b$), the *deadline* ($T_b$) and the *pricing strategy* ($p_b$). On every round, she proposes a price according to the following pricing strategy:

$$p_b(i) = p_0 + \left(V - p_0\right) \cdot \left(i \cdot T_b^{-1}\right)^k \qquad (1)$$

where $p_0$ is the first proposed price (usually it is a very small price) and $i$ is the current round. The parameter $k$ defines the strategy and could be: patient, aggressive, neutral.

The seller negotiates for a number of products with a number of buyers. It is of high importance to note that the buyer is not aware of any of the seller's characteristics. The characteristics of the seller are: the product *cost* ($c$), the *discount factor* ($\delta_s$), the *utility function* ($U_s$), the *deadline* ($T_s$), the intended *profit* ($\varepsilon$) and the *pricing strategy* ($p_s$). Usually, the seller starts by proposing a large price which equals to $c + \varepsilon$ and according to her strategy she can reduce the offers as the negotiation progresses. Furthermore, she can change the strategy at every negotiation round. There can be four types of sellers: a) *neutral*, b) *patient*, c) *impatient* and d) sellers that change the strategy at every round (*mixed behavior*). These strategies are reflected by:

$$p_s(i) = c + \varepsilon \cdot (1 - i \cdot T_s^{-1})^k \qquad (2)$$

where $i$ is the round index. The parameter $k$ denotes the policy of the seller. An aggressive seller wants to conclude the negotiation process as soon as possible. Following this strategy, her intention is to quickly reduce her prices in order to challenge the buyer to accept her offers. For this, we propose the strategy defined in [8]. The proposed pricing function fully adapts the resulting values to the product characteristics. The goods, available at the seller, are ranked according to their popularity. We can defined the product popularity based on Zipf's Law. The proposed pricing strategy is:

$$p_S(i) = \frac{\varepsilon}{i^{q+1}} + c \tag{3}$$

where q is the popularity measure. The described pricing function indicates a very aggressive seller that tries to conclude as many transactions as she can in a certain period of time. Based on the above described strategy, we adopt the deadline calculation process proposed in [6, 7, 8]. The seller deadline is calculated as follows:

$$T_S \approx (\alpha \cdot \varepsilon \cdot (q+1))^{\frac{1}{q+2}} \tag{4}$$

where $\alpha$ is a scaling factor which depends on the seller's strategy. If the seller follows a patient policy the $\alpha$ factor assumes a relatively high value.

In [7], we present a fuzzy logic system for the derivation of the $\alpha$ value. The proposed system is based on parameters q, $\varepsilon$ and the final result is the value of $\alpha$. For each parameter, specific linguistic values are defined $A_1 = A_2 = B_1$ {low, medium, high} as well as the corresponding trapezoidal fuzzy sets. Concerning the scaling factor $\alpha$, a fuzzy value *low a* indicates that, the seller is an impatient player which stays for a few rounds in the BG. A *medium* and a *high* value of *a* indicates a medium and high value of patience respectively. For a more fine-grained resolution of the fuzzy linguistic values of *a*, we use the linguistic modifier *very*; $very(\mu(a)) = \mu(a)^2$. Specifically, *very low a* denotes that the seller wants to sell the product as soon as possible, thus, participating in only a few negotiation rounds and *very high a* denotes that, the seller is a very patient player. The strategy of the seller is mapped into a set of fuzzy rules in order for the seller to estimate / calculate the time limit *a* for the specific negotiation with a specific buyer. Finally, in [6], we propose the automatic fuzzy rules generation from data provided by experts. The proposed methodology is simple, yet, efficient as experiments show.

### 2.2    Sequential Equilibrium Definition

The outcome of the negotiation mainly depends on two issues: a) the players' deadlines, and, b) the players' strategies. However, as no knowledge is present, players should predict both issues based on the offers made in order to take the most appropriate decision. Let us examine first the prediction of deadlines. We describe the buyer side. A similar approach stands for the seller side. We consider that the buyer could adopt a Uniform distribution for the seller deadline estimation:

$$P(t \le T_S \le t + \Delta_t, \Delta_t \to 0) = \frac{1}{t_{max}}, \text{if } 0 < t < t_{max} \tag{5}$$

where $P(t \le T_S \le t + \Delta_t, \Delta_t \to 0)$ represents the probability that the seller deadline is equal to t. If $t \ge t_{max}$, we consider that the probability of the seller deadline expiration is equal to 1. For the pricing strategy distribution estimation, we adopt the known Kernel Density Estimation (KDE) methodology. KDE is a methodology for estimating the PDF of an unknown distribution. The Kernel estimator of this distribution is defined by the following equation:

$$\hat{f}_h(x) = \frac{1}{N \cdot h} \sum_{i=1}^{N} K\left(\frac{x - x_i}{h}\right) \tag{6}$$

where x is the examined variable, N is the sample's size, h is the bandwidth of the kernel and $K(\cdot)$ is the Kernel function. In our model, we use the Gaussian function and, thus, the probability distribution of the seller pricing strategy can be given by:

$$P(x) = \frac{1}{N \cdot h} \cdot \sum_{i=1}^{N} K\left(\frac{x - p_{si}}{h}\right) \tag{7}$$

In our scenario, without loss of generality and for simplicity in our calculations, we take the bandwidth equal to 1. Thus, Equation (7) is transformed to the following:

$$P(x) = \frac{1}{N} \cdot \sum_{i=1}^{N} \frac{1}{\sqrt{2 \cdot \pi}} \cdot e^{-\frac{(x - p_{si})^2}{2}} \tag{8}$$

In the above equations, $p_{si}$ depicts the seller price proposed at every round of the negotiation. Based on the above analysis, we take the cumulative distribution function of the seller pricing strategy:

$$CDF(x) = \frac{1}{N} \cdot \sum_{i=1}^{N} \frac{1}{2} \cdot \left(1 + erf\left(\frac{x - p_{si}}{\sqrt{2}}\right)\right) \tag{9}$$

**Definition:** *After an offer made by the seller, the buyer decides on the next action to be taken based on:*

$$D_b^t(p_s^t, V, t) = \begin{cases} \text{Accept,} & \text{if } [[(T_s = t) \text{ or } (T_s = t+1)] \text{ and } (p_s^t \leq V)] \text{ or } [U_b(p_s^t) \geq U_b(p_b^t)] \\ \text{Defect} & \text{if } [(T_b = t) \text{ or } (T_b = t+1)] \text{ and } (p_s^t > V) \\ \text{Reject and Propose} & \text{otherwise} \end{cases} \tag{10}$$

The utility that the buyer gains at round t of the negotiation process is given by the following equation:

$$\delta_b^{t-1} \cdot U_b^t \cdot P(\text{accept}) \tag{11}$$

Where P(accept) represents the probability that the buyer accepts the seller offer. Based on the KDE methodology, we can make the following proposition:

**Proposition:** *For the negotiation model described in (10) there is a strategy combination which satisfies the sequential equilibrium. If the buyer adopts a Uniform distribution for estimating the seller deadline and the KDE for estimating the seller pricing strategy the buyer should reject the seller's offers at every round of the negotiation and accept only when:*

$$p_s^t \leq \begin{cases} V - \dfrac{t_{max} \cdot (V - p_b^t) \cdot z'}{2 \cdot z + z' \cdot t_{max} - 2 \cdot z \cdot z'}, & \text{if } t < t_{max} \\[4mm] V - \dfrac{(V - p_b^t) \cdot z'}{2 \cdot z + z' - 2 \cdot z \cdot z'}, & \text{if } t \geq t_{max} \end{cases} \tag{12}$$

*with*

$$z = \sum_{i=1}^{N} \frac{1}{2} \cdot \left(1 + \sqrt{1 - e^{-Q \cdot \frac{\frac{4}{\pi} + a \cdot Q}{1 + a \cdot Q}}}\right) \quad Q = \frac{(V - p_{si})^2}{2} \tag{13}$$

$$z' = \sum_{i=1}^{N} \frac{1}{2} \cdot \left(1 + \sqrt{1 - e^{-W \cdot \frac{\frac{4}{\pi} + a \cdot W}{1 + a \cdot W}}}\right) \quad W = \frac{(p_b^t - p_{si})^2}{2} \tag{14}$$

*and N is the number of seller proposals till the current round.*

### 2.3 Seller Decision Making Mechanism

In the seller side, we also propose a decision making mechanism based on fuzzy logic. The proposed mechanism is used at every round in which the buyer proposes a price. The decision refers to whether the seller should accept or reject the proposed offer. The decision is based on the *Acceptance Degree* (AD) which shows when the seller should accept the buyer's offer and depends on the following parameters: a) the *time difference* between the current time of the negotiation and the seller's deadline (t), b) the *belief* about the expiration of the buyer's deadline (b), c) the absolute value of the *price difference* between the buyer's proposal and the upcoming seller's offer (d), and, d) the *number of buyers* waiting/interacting for/with the seller (N). For the reasoning process, an input to the fuzzy system might be described as: a round which has increased time difference with the deadline $T_s$ is represented by the value high ($l(u_i)$ = High), medium difference ($l(u_i)$ = Medium) or low difference ($l(u_i)$ = Low) where $u_i$ = the time difference between the current round and the deadline of the seller. The same rationale stands for the remaining parameters. The form of the rules is:

$R_j$ : *If t is $A_{1(j)}$ AND b is $A_{2(j)}$ AND d is $A_{3(j)}$ AND N is $A_{4(j)}$ Then AD is $B_{(j)}$.*

where $A_{i(j)}$ and $B_{(j)}$ is the fuzzy set representing the j$^{th}$ linguistic value for the input parameter i and for the output parameter AD, respectively. The linguistic expressions of the values for the parameters t, b, d, N and AD are defined in the sets $A_1 = A_2 = A_3 = A_4 = B_1$ {Low, Medium, High} and we use for them trapezoidal fuzzy sets. We consider three sigmoid functions for parameters t, d, and N in order to produce values in the range [0,1]. Concerning the acceptance degree AD, *Low AD* indicates that the seller should not accept the buyer's proposal and make a counter offer, a *Medium* and *High AD* indicates a neutral and positive attitude to the buyer's offer. Specially, a high AD value means that the seller should accept the buyer's proposal and conclude the negotiation before the expiration of her deadline. Finally, the strategy of the seller can be mapped into a set of fuzzy rules in order for the seller to decide if she will accept or reject the buyer's offer. Rules are defined by experts.

## 2.4    Buyer Decision Making Mechanism

In this section, we focus on the buyer side and shortly describe the reasoning mechanism adopted by her in order to decide on the acceptance or rejection of a seller's offer [3]. We developed a fuzzy logic system, which determines the buyer's reaction to the seller's proposals. We define as *Acceptance Degree* (AD) the capability of the buyer to accept the seller's offer. The AD parameter reflects the willingness of the buyer to accept the price for a product offered by the seller hoping to maximize her utility. High AD degree indicates that the buyer accepts the seller's offer and concludes the negotiation. Specifically, the *AD* degree depends on the following parameters: a) the *relevance factor* (r) which shows to which extend the product corresponds to the buyer's needs, b) the absolute value of the *price difference* (d) between the seller's proposal and the upcoming buyer's offer, c) the *belief* (b) about the expiration of the seller's deadline, d) the *time difference* (t) between the current time of the negotiation and the buyer's deadline, and, e) the buyer's *valuation* (V) about the product. The system relies on a rule base for inference. We adopt the multi-input single-output (MISO) form of the linguistic rule $R_j$, with $r = u_1$, $d = u_2$, $b = u_3$, $t = u_4$, $V = u_5$ and $y = AD$, that is,

**$R_j$ : If r is $A_{1(j)}$ AND d is $A_{2(j)}$ AND b is $A_{3(j)}$ AND t is $A_{4(j)}$ AND V is $A_{5(j)}$ Then AD is $B_{(j)}$.**

where $A_{i(j)}$ and $B_{(j)}$ are the fuzzy sets representing the $j^{th}$ linguistic value for the input parameter *i* and for the output parameter *AD*, respectively. The system involves a three-step process: a) the fuzzification step transforms the input parameter-values into fuzzy subsets, b) using the fuzzy rule base an inference takes place for the output value (fuzzified AD), and c) the defuzzification process converts the output of the fuzzy inference into the crisp outputs for the parameter AD. For the defuzzification process, we use the *Center-of-Gravity* (COG) approach. The linguistic expressions of the values of the parameters r, d, b, t, V and *AD* are defined in the sets $A_1 = A_2 = A_3 = A_4 = A_5 = B_1 = \{low, medium, high\}$ while we utilize trapezoidal fuzzy for each of them. Specifically, a linguistic value of *low r* indicates that, the relevance of the product with the buyer's needs is low. A linguistic value of *medium r* denotes that, the relevance of the product is medium and a linguistic value of *high r* indicates that the product has increased relevance with the buyer's needs. For a more fine-grained resolution of the linguistic values of *r*, we use the linguistic modifier very: $very(\mu(r)) = \mu(r)^2$. We adopt three sigmoid functions in the range [0, 1] for parameters *d*, *t*, and *V*. Concerning the AD, a fuzzy value of *low AD* indicates that, the buyer should not accept the seller's proposal while a *medium* and a *high* value of *AD* indicate a neutral and positive attitude to the seller's offer respectively. Through the fuzzy rule-base, we imitate the human behavior when acting in a trading environment with no information on the characteristics of the other party (i.e., the seller). Our system contains ten fuzzy rules, which are defined by experts on the e-commerce domain. Our results (Table 1) show an increased utility value better than those reported in the literature (maximum value is to 0.9 with the vast majority to be equal to 0.6).

We extend the proposed fuzzy system and propose and adaptive mechanism for the buyer side [2]. The adaptive mechanism of the buyer consists of two parts: a) the

*seller price predictor*, and, b) the *fuzzy controller*. The price predictor is responsible for estimating the upcoming seller proposal. The fuzzy controller receives the estimation *error* and the *error change* at every round and produces the appropriate values for basic parameters of the buyer strategy such as the *belief (b)* and the *pricing policy (k)*. Belief shows how much the buyer beliefs that the negotiation ends at the upcoming round while the pricing policy influences the upcoming buyer price.

The seller price predictor should be objective and efficient. For this reason, we use three sub-predictors: A *linear*, a *polynomial* and a *neural network predictor*. The values provided by these predictors can be linearly combined in order to produce the final predicted price. The final predicted price is used in a fuzzy controller in order to obtain two important parameters: the buyer *pricing policy factor* and the buyer *belief* about the intentions and the deadline of the seller. In Fig. 1, we compare the performance of the 'simple' fuzzy system with the adaptive system. Both of them are compared with an optimal stopping model. We see that the agreement percentage is at the same level as well as the steps required for an agreement. However, the adaptive model achieves better agreement price compared to a theoretical optimal model.

Finally, we propose a scheme for the automatic generation of the fuzzy rule base [5]. Based on the proposed scheme, the fuzzy rule base is extracted by a number of crisp values defining the behavior of the buyer. The discussed process is more efficient as we do not need experts to define specific rules that are very difficult to cover all the aspects of a negotiation.

**Table 1.** Average intrinsic utility for the proposed fuzzy system.

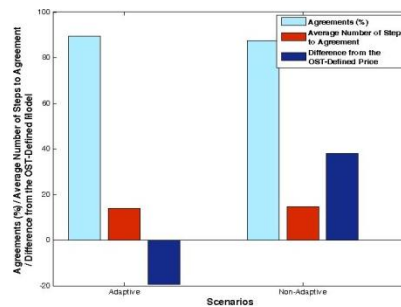| V | Average Intrinsic Utility |
|---|---|
| 5 | 0.70 |
| 20 | 0.84 |
| 50 | 0.95 |
| 80 | 0.97 |
| 100 | 0.95 |
| 150 | 0.98 |
| 200 | 0.97 |



**Fig. 1.** Comparison between the 'simple' fuzzy system and the adaptive case.

### 2.5 Selection of Products

We extend the work presented in [2, 3] and provide a methodology for defining the product *relevance factor (r)* [1]. The proposed methodology is based on the buyer request description and the product description. The buyer request can be described by: a) the context, b) the description of the desired product by using a set of keywords or simple sentences and c) a set of constraints. At the seller side, we adopt a similar approach for product description. Thus, a product description could be defined by: a) the context, b) the description by using simple sentences and c) a set of attributes. Every attribute has a name and a value.

The proposed methodology is based on the use of similarity algorithms. We choose to utilize linguistic similarity as semantic techniques require more time and resources. In order to have efficiency, we utilize a large number (16) of similarity algorithms. By using so many algorithms, we aim to avoid extreme results (e.g., very pessimistic or very optimistic). Every time we obtain the algorithms results, we calculate the values variance. If the variance is over a pre-defined threshold, we reject the minimum and the maximum value from those 16 similarity values. The final similarity value is the average of the results. It should be noted that in case where the variance is over the threshold we can reject the first and the last two values from the ranked list of algorithms' values. The developer can choose the scenario that best matches to her needs. We apply the similarity measure on the product context (request / demand) and the product description. The request context is matched against the seller product context and the request keywords are matched against the seller product description. The same process is applied for matching constraints with attributes. If the request context matches to the product context, we obtain two results: (a) keywords similarity (kFactor), and, (b) constraints similarity (cFactor). Furthermore, we combine those results with results concerning the QoS characteristics of the product. The final relevance factor value could be calculated by following a '*hard*' or a '*soft*' approach. Based on '*hard*' approach the relevance factor is calculated by the following equation:

$$r = kFactor \cdot cFactor \cdot QoSFactor \qquad (16)$$

where:

$$kFactor = \frac{k}{|\,keywords\,|} \qquad\qquad cFactor = \frac{c}{|\,constraints\,|} \qquad (17)$$

with the number of successful matches for keywords and c is the number of successful matches for constraints. Symbols $|\ |$ depict the number of keywords/constraints. QoSFactor calculation is based on price, delivery time and seller trust. Following the '*hard*' approach, the buyer is very pessimistic in characterizing a product as relevant to her goals. Following the '*soft*' approach in the calculation process, the relevance factor could be calculated through the following equation:

$$r = w_1 \cdot kFactor + w_2 \cdot cFactor + w_3 \cdot QoSFactor \qquad (18)$$

## 2.6    Concurrent Negotiations

In concurrent negotiations, the buyer could negotiate with a number of sellers trying to achieve the best agreement. For this, she is based on a number of threads. We propose a model that utilizes the Particle Swarm Optimization (PSO) algorithm in order to reach to the best solution (best agreement price). Each buyer thread can be considered as particle in the PSO algorithm. They should converge to the optimal solution which is the best price for the specific group of sellers.

The buyer will accept offers that are below her valuation. Each particle initially defines its own pricing strategy. The buyer threads follow the equilibrium path. All the buyer threads have the same deadline. Each particle negotiates autonomously with a specific seller. If an agreement is "reached" then the specific thread sends the agreement message to the rest of them. We consider that the communication time is negligible. The personal best position is the smallest price for which each particle negotiates with the seller. The global best is the smallest agreement price defined in a negotiation. The global best is defined when an agreement takes place. If no agreement is present then all the particles have velocity equal to 0 and continue to propose prices according to the pricing strategy. If a particle does not have an agreement and receives an agreement message from another particle then it switches its state and: a) if her current negotiated price is smallest than the global best, she remains at the current status or b) if the global best is smaller than current negotiated price, she changes the pricing strategy in order to reach the global best position. Particles velocity is defined when particles want to change their position (an agreement was "announced" in a better price). The velocity is initially calculated when an agreement is announced and for every round after that. If the global best is not smaller than the current price the velocity is set equal to 0 or else the velocity is calculated by the PSO algorithm. The velocity affects the pricing strategy and the proposed by the buyer price respectively.

## 3    Conclusions

The interaction between autonomous entities in dynamic environments (such as EMs) is a very interesting research issue. In this thesis, we present decision making mechanisms for the buyer and the seller side. The mechanism utilizes fuzzy logic that is appropriate for handling uncertainty. We also propose models for choosing the most appropriate product in the buyer side while we analyze the equilibrium path for the negotiation process with a seller. Moreover, we propose a prediction mechanism for the seller pricing strategy. The prediction engine in combination with the negotiation parameters provides the necessary information for the buyer to adapt her behavior. Additionally, we study concurrent negotiations and propose the use of the PSO algorithm. The advantage is that buyer threads through a team work find the optimal solution. The difference of our work from the efforts found in the literature is that we do not any coordination between threads. Experimental results show increased number of agreements in combination with the increased utility for both parties. The fuzzy logic system is proved to be very efficient for both the buyer and the seller.

# References

1. Kostas Kolomvatsos, and Stathes Hadjiefthymiades, 'An Extended Q-Gram Algorithm for Calculating the Relevance Factor of Products in Electronic Marketplaces', accepted for publication in *Elsevier Electronic Commerce Research and Applications (ECRA)*, 2013.

2. Kostas Kolomvatsos and Stathes Hadjiefthymiades, 'Buyer Behavior Adaptation Based on a Fuzzy Logic Controller and Prediction Techniques', *Elsevier Fuzzy Sets and Systems (FSS)*, February 2012, pp. 30-52.

3. Kostas Kolomvatsos, Christos Anagnostopoulos, and Stathes Hadjiefthymiades, 'A Fuzzy Logic for Bargaining in Information Markets', *ACM Transactions on Intelligent Systems and Technology (ACM TIST)*, February 2012, vol. 3(2), pp. art. No 32.

4. George Boulougaris, Kostas Kolomvatsos, and Stathes Hadjiefthymiades, 'Building the Knowledge Base of a Buyer Agent Using Reinforcement Learning Techniques', In Proc. of the *2010 IEEE World Congress on Computational Intelligence (WCCI 2010), IJCNN*, July 18th - 23rd, Barcelona, Spain, pp. 1166-1173.

5. Roi Arapoglou, Kostas Kolomvatsos, and Stathes Hadjiefthymiades, 'Buyer Agent Decision Process Based on Automatic Fuzzy Rules Generation Methods', In Proc. of the *2010 IEEE World Congress on Computational Intelligence (WCCI 2010), FUZ-IEEE*, July 18th - 23rd, Barcelona, Spain, pp. 856-863.

6. Kostas Kolomvatsos and Stathes Hadjiefthymiades, 'Automatic Fuzzy Rules Generation for the Deadline Calculation of a Seller Agent', In Proc. of the *9th International Symposium on Autonomous Decentralized Systems (ISADS 2009)*, Athens, Greece, March 23-25, 2009, pp. 429-434.

7. Kostas Kolomvatsos, Christos Anagnostopoulos, and Stathes Hadjiefthymiades, 'On The Use of Fuzzy Logic in a Seller Bargaining Game', In Proc. of the *32nd Annual IEEE International Computer Software and Applications Conference (COMPSAC 2008)*, July 28th - August 1st, Turku, Finland, 2008, pp. 184-191.

8. Kostas Kolomvatsos and Stathes Hadjiefthymiades, 'Implicit Deadline Calculation for Seller Agent Bargaining in Information Marketplaces', In Proc. of the *2nd International Conference on Complex, Intelligent and Software Intensive Systems (CISIS 2008)*, March 4th - 7th, Polytechnic University of Catalonia, Barcelona, Spain, 2008, pp. 184-190.

9. Kostas Kolomvatsos and Stathes Hadjiefthymiades, 'On the Use of Fuzzy Logic in Electronic Marketplaces', in the book 'Cross Disciplinary Applications of Artificial Intelligence and Pattern Recognition: Advancing Technologies', ed. Vijay Mago, IGI Global, 2011.

10. Kostas Kolomvatsos and Stathes Hadjiefthymiades, 'Defining Time Constraints for Sellers in Electronic Markets', in the 'Encyclopedia of E-Business Development and Management in the Global Economy', ed. In Lee, IGI Global, 2010.

11. Kostas Kolomvatsos and Stathes Hadjiefthymiades, 'How Can We Trust Agents in Multi-Agent Environments?', Chapter in 'Intelligence Integration in Distributed Knowledge Management', eds D. Krol and N. T. Nguyen, IDEA Group Inc., 2008.

# Wordspotting on Historical Document Images

Thomas Konidaris [†]

National and Kapodistrian University of Athens
Institute of Informatics and Telecommunications
`tkonid@iit.demokritos.gr`

**Abstract.** In this dissertation innovative methods of wordspotting on historical printed documents are presented. In particular, two methods based on document segmentation on word level have been developed. The first method uses a hybrid feature scheme for word matching based on zones and projections. It also uses a process of creating query keyword images for any word using synthetic data. The synthetic words are created using images of individual characters taken from the processed documents. The method also presents a process allowing user feedback in order to improve the final results. The second method uses the Dynamic Time Warping (DTW) algorithm for comparing word images. It assist the transition between the synthetic data and real data comparison. Synthetic data and real data differ and DTW allows a better alignment between the features of the two images. Again, feedback can be applied to improve the results. Furthermore, a method that uses no segmentation on the document images has been also developed. The method overcomes the problem of incorrect segmentation that affect the final results since it detects query keyword images directly on entire document page images. It also allows for partial matching such as detecting word that are included in larger ones. The evaluation of the aforementioned methods showed satisfactory results presenting better performance against competitive methods of wordspotting.

## 1 Introduction

World wide libraries hold a vast amount of historical documents in terms of books, papers, drawings, journals, etc. These documents are highly valuable items due to the information they contain as well as the historical importance and rarity that characterizes them. The digitization of such archival and historical collections is an ongoing process that results to digital content which allow access to the information without distorting the original material. It is clear that efficient indexing and retrieval are important prerequisites of any system that manipulates such digital content. Optical Character Recognition (OCR) is a standard technology that is widely used in indexing documents with noticeable results in contemporary documents. However, historical documents are prone to a number of difficulties such as typesetting imperfections, document degradations and low print quality which decrease the performance level of OCR systems [14], [23], [10], [9].

---

[†] Dissertation Advisor: Sergios Theodoridis, Professor

## 2 Dissertation Summary

Word spotting is an alternative methodology for document indexing based on spotting words directly on images without the use of any OCR procedures. Thus, a word spotting system attempts to detect words as a whole rather than to exactly recognize the characters as in OCR. In a typical scenario, the query image is selected from a set of predefined keywords of interest or is interactively defined by the user by cropping a rectangular image area that serves as query example. The word spotting system uses the query and detects similar words in document images based on image matching techniques without any conversion of the images into readable text. Extensive studies have shown that indexing terms in documents automatically using a word spotting system makes it possible to use costly human labor more sparingly than a full transcription would require [22]. Several word spotting methods rely on a pre-processing step where the document image is segmented into words. The segmented words are then compared to the query image in order to detect potential matches. Recently, segmentation-free methods have been also proposed that do not require any segmentation and the document image is treated as one entity by-passing any errors that may occur due to poor segmentation results. In this line, we propose a segmentation-free word spotting method for historical printed documents. The method is based on local keypoint correspondences and consists of two distinct steps that determine candidate image areas in order to accurately extract the final bounding boxes which indicate the word instances in the document page. The method is evaluated using two different datasets of different languages and the experimental results show that the proposed method outperformed significantly the competitive approaches.

The word spotting literature can be divided into two main categories depending upon whether segmentation of the document image is applied or not. Indeed, there are several methods that are based on page segmentation as a pre-processing step, while others are applied directly to the document image. Additionally, a variety of features are used in order to describe the query word as well as the document image. These features strive to efficiently express the geometric and local information of the visual content and include projection profiles, Gabor features, zones and gradient-based features, to name a few. Keypoint-based local features have been also successfully used in order to describe document images as a set of local feature vectors that are invariant to scale changes, illumination and distortions. The Scale Invariant Feature Transform (SIFT) [19] is a well known technique in this category that produces an adequate number of distinctive features even for small visual objects. In the following, we summarize some word spotting techniques that rely at least in part on segmentation as well as approaches where no segmentation is required.

### 2.1 Segmentation based methods

There are three levels of page segmentation that are typically used for detecting words in documents, namely segmentation into lines [11], [20], words [12],

[13], [21] [24], [25], or even characters [2], [8]. Profile features, such as upper or lower word profiles, projection, density or transition profiles have been reported to successfully represent words in a document image that has undergone word level segmentation [22], [24]. Fusion of multiple features is also adopted in several studies in order to improve the word image description. For example, in [26] a multiple feature scheme is used consisting of projection profiles, upper/lower word profiles and background-to-ink transitions. Similarly, Jawahar et al [6] involve word profiles to describe the outline shape of the word, structural features to extract statistical information like moments or variation and finally Fourier coefficients as a compact representation of the features in the frequency domain. In [12] a hybrid feature scheme based on a combination of projection profiles and upper/lower word profiles is used for matching words segmented from document images. In [8], a word spotting method is proposed based on mesh features and in [29] and [33] the feature scheme used is gradient-based binary features. Another feature used for word spotting is based on skeletons and is used in the works of [18] and [7]. Gabor features can also be applied for word spotting as proposed in [2]. In Li et al [17] a word image is decomposed into vertical strokes and a stroke-based coding scheme is built for all the word in the document database. Considering features basedon local keypoints, Ataer et. al. [1] use SIFT features in order to match segmented words from Ottoman documents. Similarly, in [32] a word image matching method is presented using SIFT descriptors on keypoints that are extracted using the Fast-Corner-Detection algorithm [27]. These features are quantized into visual terms (visterms) using hierarchical K-Means algorithm and indexed using an inverted file. In [30] a word spotting method based on line segmentation is presented. The method uses a sliding window over each line. The matching is performed using dynamic programming and slit style HOG features. In the previous methods, the segmented words are presented as feature vectors and Dynamic Time Warping is an algorithm that has been extensively used to match words based on these vectors [25], [26], [6], [13], [11]. Other matching techniques are based on morphological variants [21], voting schemes [18], [1], [32], similarity distances [12], [7], character or string matching [8], [17] and correlation measures [29], [20], [33]. Overall, document segmentation results to higher level structures that are semantically important and can be further explored. On the other hand, detection methods based on segmentation results are intrinsically prone to errors like over- or under- segmentation as well as well as partial occlusion and mis-segmentation.

## 2.2 Segmentation-free approaches

Although, there is a very large collection of published work concerning the segmentation approach, in the recent years there is a growing research interest concerning segmentation-free methods. There are cases where documents cannot be segmented correctly leading to insufficient results. The segmentation-free approaches overcome the problems associated to bad segmentation results by treating the document image as a whole. In [4] a template matching method

based on pixel densities is used for locating words in documents without segmenting them. Although the method provides rotation and scale invariance, this is applied on limited extent. In [16] an alphabet is used that is manually selected from each document collection processed. The alphabet is used to create word instances that serve as queries. The features extracted are based on gradient values. In [15] gradient information is also used as features for the word images. Word interest points are matched against document images and try to locate zones of interest presenting similar features. Local image features have been also used in segmentation-free methods trying to benefit from the scale and rotation invariance they offer as well as their robustness to noise. Such methods usually involve a voting scheme in order to detect and localize potential word matches in the document image. In this line, Rusinol et al [28] opt SIFT features in a bag-of-visual-words approach. The method is applied on both handwritten and printed documents. The SIFT features are extracted using small predefined squared areas that are assumed to cover most of the font sizes. The search space does not correspond to the entire document image but rather to overlapping local patches of fixed geometry. However, several assumptions concerning the size of the patch and the expected font sizes seriously affect the generalization and the applicability of the method. Furthermore, as the authors mention, the performance of the system is highly related to the length of the queried words.

## 2.3 Contribution of Research

Our research concentrated into both approaches. In particular we have developed two methods that are based on the segmentation approach and one method that follows the segmentation-free approach. The methods that are based on the segmentation approach segment the document on word level. The first method uses synthetic data to create query keyword images. Each query keyword is created in a synthetic manner using individual character images taken from the processed documents. This way we are able to construct any query word image we like. The feature scheme used combines two different features. The first one is zones and the second is projection profiles. The features are matched using a simple distance metric. The advantage of the method lies on the fact that is very fast in producing an initial set of results. These are further improved through a user's feedback process. The second segmentation-based method uses synthetic data to create query keyword images and a combination of four different features. However, unlike the above method, the features are compared using the Dynamic Time Warping (DTW) algorithm. The DTW algorithms manages to overcome local distortions between the compared feature vectors. This method performs better than the former segmentation method but needs more time to complete as DTW algorithm poses bigger complexity.

On the other hand, the segmentation-free method that was developed comes to solve the problem of incorrect segmentation. There are cases where the documents are not segmented correctly. This error percentage affects the overall performance of the methods. The segmentation-free method does not require the documents to be segmented at any level. Rather, the query keyword images

are compared with the entire document page images. The method uses the SIFT algorithm for the extaction of keypoints and their descriptors. The performance of the method is very satidfactory.

# 3 Results and Discussion

In this section we will discuss the segmentation-free method that was developed during our research. In the proposed method we adopt a segmentation-free word spotting approach in order to overcome the poor segmentation results that usually characterize historical documents. We are based on SIFT features that have been proved to provide robustness concerning low image quality and image degradation. However, detection of word instances based on direct matching between query and image SIFT keypoints leads to unsatisfactory results. This is due to the spatial scattering of matching correspondences since a query keypoint may be similar to a large number of document page keypoints. These document keypoints do not a priori belong to correct word instances. Furthermore, the existence of multiple word instances in the same document page does not allow the query image to be matched by a sufficient number of correspondences.

The proposed method does not adopt the original matching process described in the SIFT algorithm, but instead a two step approach is followed. In the first step, for every keypoint in the query keyword image, the nearest K points are found in the document page image. These document keypoints are used as indicators in order to create candidate image areas. In the second step, each candidate image area is matched against the query keyword image. The keypoint correspondences are used by the RANSAC algorithm in order to estimate the final bounding boxes indicating the detected word instances. Furthermore, we use the strength of SIFT descriptors in a way that multiple instances of the desired word can be found on the document page.

## 3.1 Detection of Candidate Image Areas

The first step of the proposed method involves the matching of the query keyword keypoints to the document keypoints. The purpose is to find point correspondences on the document image that will serve as indicators of candidate image areas. For each keypoint of the query keyword we locate the K most similar keypoints on the document image. The value of $K$ is experimentally defined as discussed in section 4. Let $f_q$ and $f_d$ be the SIFT feature vectors of the $i^{th}$ keypoint in the query keyword image and the $j^{th}$ keypoint in the document image, respectively. The distance between these two keypoints is calculated as follows:

$$d(i,j) = cos^{-1}(\langle f_q^i, f_d^j \rangle) \tag{1}$$

where $\langle f_q, f_d \rangle$ denotes the dot product between the two normalized vectors.

Each pair of corresponding keypoints defines a candidate image areas on the document page. Since we know the relative position of the query keyword

keypoint in respect to the edges of the query keyword image we define a bounding box around the keypoint on the document image taking into account the position of the corresponding keypoint of the query keyword image. Let $p_q(x_q, y_q)$ be a point on the query keyword image and $p_d(x_d, y_d)$ be its corresponding point on the document page image. Let dx, dy be the distance of the query keypoint from the left and the top edge of the query keyword image respectively. The bounding box surrounding the candidate image area is defined by its top-left $(x_{min}, y_{min})$ and bottom-right $(x_{max}, y_{max})$ corners is given by the following equations:

$$x_{min} = x_d - (\frac{sc_d}{sc_q} \cdot dx \cdot t_s) \tag{2}$$

$$y_{min} = y_d - (\frac{sc_d}{sc_q} \cdot dy \cdot t_s) \tag{3}$$

$$x_{max} = x_d + [(w_q - x_q) \cdot \frac{sc_d}{sc_q} \cdot t_s] \tag{4}$$

$$y_{max} = y_d + [(h_q - y_q) \cdot \frac{sc_d}{sc_q} \cdot t_s] \tag{5}$$

where $w_q$ and $h_q$ are the width and height of the query keyword image, respectively. Parameter $t_s$ is the boundary size factor, which gives extra space to the boundaries of the candidate image areas and has been experimentally set to 1.1. Variables $sc_q$ and $sc_d$ are the scales of the query keypoints $p_q$ and $p_d$, respectively, as provided by the SIFT algorithm.

### 3.2 Detection of Word Instances

In the previous section we matched the query keyword image with the entire document page image in order to use the matching keypoints as indicators for creating candidate image areas. These areas cannot guarantee that they contain the query word under consideration. For this reason the keypoints of the query keyword image are matched against the keypoints of each candidate image area. For each keypoint on the query keyword image we find the most similar keypoint on the candidate image area using Eq. 1. In order to estimate a model that describes the efficiency of these keypoint correspondences the RANSAC algorithm [3] is involved. RANSAC is an iterative method that can efficiently estimate the parameters of a model even when the measurements contain outliers. Using RANSAC the number of inliers is calculated, that is, the number of corresponding pairs that are conveniently described by the model. Moreover, the keypoint correspondences are used in order to calculated a homography that serves as a transformation matrix from the query keyword image to the candidate image area plane [5]. There must be at least four point correspondences to calculate the homography matrix. Let $P_q = (x_q, y_q)$ be a point in the query keyword image and $P_c = (x_c, yc)$ be the corresponding point in the candidate image area. The transformation between these two points can be given by the following equation:

$$P_c = H \cdot P_q \tag{6}$$

where $H$ is the homography matrix. The above equation can take the form:

$$\begin{bmatrix} x_q \\ y_q \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} x_c \\ y_c \\ 1 \end{bmatrix}$$

This process is applied for all candidate image areas aiming to produce a set of bounding boxes that are afterwards ranked according to their matching efficiency. The inliers percent provides an indicator of the goodness of fit regarding the RANSAC model. However, there may be a number of detected bounding boxes having equal inliers percent value. In order to distinguish them, we propose to divide the inliers percentage value L by a quantity D which corresponds to the sum of the distances between the query and the candidate image area keypoints. Thus, for a candidate image area the ranking value V is calculated as follows:

$$V = \frac{L}{D} \tag{7}$$

### 3.3 Removing Overlapping Results

We have seen that the candidate image areas are created using the point correspondences between the query keyword image keypoints and the keypoints of the document page image. There are cases where more than one candidate image areas correspond to the same word in the document image. Therefore, we end up with overlapping bounding boxes, each of them having different ranking values $V$ as calculated by Eq. 7. On the document image, two bounding boxes $B_i$ and $B_j$ are considered overlapping if the following equation holds:

$$IoU = \frac{B_i \cap B_j}{B_i \cup B_j} \geq t_v \tag{8}$$

where $t_v$ has been experimentally defined equal to 0.3. The bounding box that has the larger ranking value $V$ among the overlapping bounding boxes is the one kept while the others are discarded from the list. Figure 6 illustrates an example of resulting bounding boxes concerning the same candidate image area. The two bounding boxes are considered overlapping since their intersection over union ration exceeds the threshold $t_v$ , as shown in Figure 1(a). However, the bounding box in Figure 1(b) has a smaller ranking value $V$ than the bounding box in 6(a) and it is discarded from the list of bounding boxes. The remaining bounding boxes are further filtered out using the word length of the query keyword image. The bounding boxes which the following equation holds are excluded from the list of the results.
where $w_b$ is the length of a bounding box from the results list, $w_q$ is the length of the query keyword image and $t_w$ is the threshold which has been experimentally set to 0.4.

$$|\frac{w_b - w_q}{w_q}| \geq t_w \tag{9}$$

(a) $IoU = 0.3212$      (b) $V_i = 0.0207$      (c) $V_j = 0.0183$
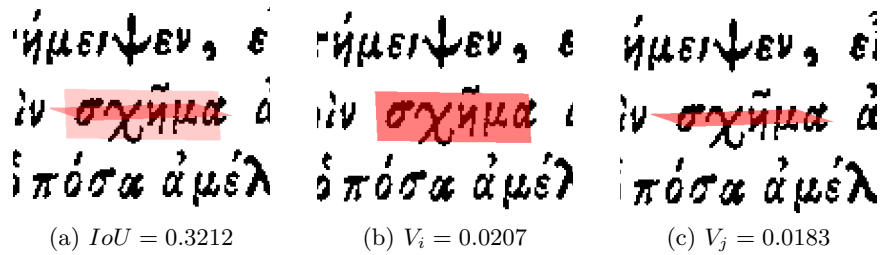
Fig. 1: Resulting bounding boxes $B_i$ and $B_j$ for the same word on a document page image. (a) Their intersection over union ratio, (b) The bounding box $B_i$ with ranking value of 0.0207, (c) The bounding box $B_j$ with ranking value of 0.0183. The bounding box $B_j$ is discarded since it has lower ranking value $V$.

### 3.4 Experimental Results

The experiments that were conducted in order to evaluate the proposed method included two different datasets. The first dataset consists of 100 pages from a Greek historical typewritten book from the period of Renaissance and Enlightenment. The second dataset consists of 100 pages of a German historical typewritten book of Eckartshausen which was published in 1788 and is owned by the Bavarian State Library [31]. For the Greek dataset we have used seven (7) query keyword images as queries and for the German dataset we have used ten (10) keyword images as queries.

The proposed method was compared against the method presented in [4] and original SIFT. Figure 2 shows the performance of the proposed method against the competitive ones concerning the Greek dataset. Likewise, the results for the German dataset are shown in Figure 3.
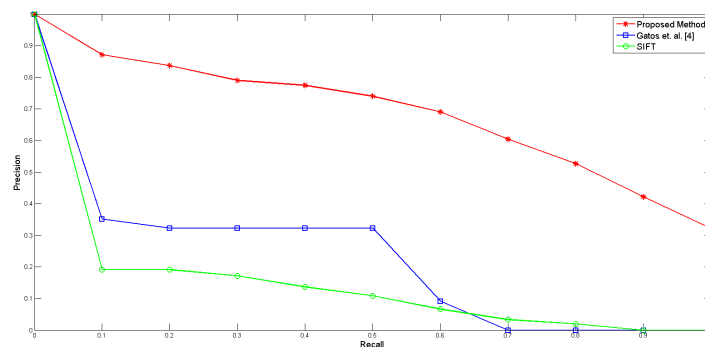


Fig. 2: Precision-Recall curves for the different methods concerning the Greek dataset.
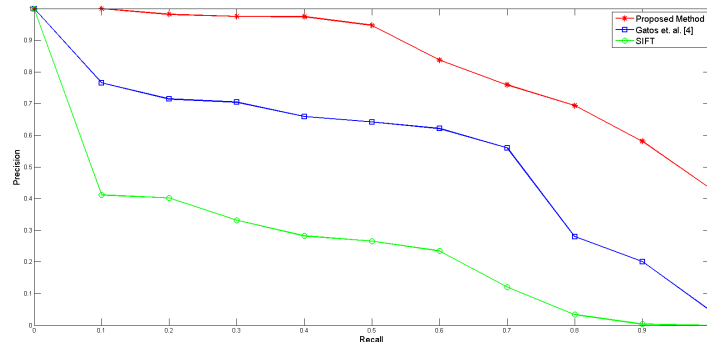
Fig. 3: Precision-Recall curves for the different methods concerning the German dataset.

The proposed method outperforms the competitive methods significantly. Furthermore, there is clear evidence that if we apply original SIFT matching between a query keyword image and a document page image in order to get keypoint correspondences that will serve as indicators to the creation of candidate image areas, the results are very poor.

## 4  Conclusions

The research aimed at creating methods for wordspotting. Our methods touched both segmentation and segmentation-free approaches. In particular, we have developed two methods that are based on the segmentation approach. These are a fast method that uses synthetic data, user feedback and a feature scheme that combines two different features. The method performed very well and in small amount of time. The second method used the DTW algorithm in order to solve the problem of variations and distortions that is found between words. This method, outperforms the first segmentation based method giving much better results. As far as the segmentation-free method is concerned, we have developed a method that does not require any prior segmentation of the document images. The query keywords are matched against the entire document page image. Such segmentation-free methods starting to gain great attention since they overcome the problems of bad segmentation and can even be used at documents that segmentation fails dramatically. The results of the mehtod are very encouraging as well as satisfactory.

## References

1. Esra Ataer and Pinar Duygulu. Matching ottoman words: an image retrieval approach to historical document indexing. In *Proceedings of the 6th ACM international conference on Image and video retrieval*, pages 341-347, 2007.

2. H. Cao and V. Govindaraju. Template-free word spotting in low-quality manuscripts. In *6th International Conference on Advances in Pattern Recognition (ICAPR'07)*, pages 45-53, 2007.

3. M. A. Fishler and R. C. Bolles. A paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the Association of Computer Machinery*, 24(6):381-395, 1981.

4. B. Gatos and I. Pratikakis. Segmentation-free word spotting in historical printed documents. In *10th International Conference on Document Analysis and Recognition (ICDAR'09)*, pages 271-275, Barcelona, Spain, 2009.

5. R. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2nd edition, 2003.

6. C. V. Jawahar, A. Balasubramanian, and M. Meshesha. Word-level access to document image datasets. In *Proceedings of the Workshop on Computer Vision, Graphics and Image Processing (WCVGIP)*, pages 73-76, 2004.

7. P. Keaton, H. Greenspan, and R. Goodman. Keyword spotting for cursive document retrieval. In *Workshop on Document Image Analysis*, pages 74-81, San Juan, Puerto Rico, 1997.

8. S. Kim, S. Park, C. Jeong, J. Kim, H. Park, and G. Lee. Keyword spotting on korean document images by matching the keyword image. In *Digital Libraries: Implementing Strategies and Sharing Experiences*, volume 3815, pages 158-166, 2005.

9. V. Kluzner, A. Tzadok, Y. Shimony, E. Walach, and A. Antonacopoulos. A complete optical character recognition methodology for historical documents. In *Tenth International Conference on Document Analysis and Recognition (ICDAR)*, pages 501-505, Barcelona, 2009.

10. A. Koerich, R. Sabourin, and C. Y. Suen. Devanagari ocr using a recognition driven segmentation framework and stochastic language models. *International Journal on Document Analysis and Recognition (IJDAR)*, 6(2):126-144, 2003.

11. A. Kolcz, J. Alspector, M. Augusteijn, R. Carlson, and G. Viorel Popescu. A line-oriented approach to word spotting in handwritten documents. *Journal of Pattern Analysis and Applications*, 3(2):153-168, 2000.

12. T. Konidaris, B. Gatos, K. Ntzios, I. Pratikakis, S. Theodoridis, and S. J. Perantonis. Keyword-guided word spotting in historical printed documents using synthetic data and user feedback. *Internation Journal of Document Analysis and Recognition (IJDAR), special issue on historical documents*, 9(24):167-177, 2007.

13. T. Konidaris, B. Gatos, S. J. Perantonis, and A. Kesidis. Keyword matching in historical machine-printed documents using synthetic data, word portions and dynamic time warping. In *The eighth IAPR Workshop on Document Analysis Systems*, pages 539-545, 2008.

14. F. Lebourgeois, J.-L. Henry, and H. Emptoz. An ocr system for printed documents. In *Proceedings of IAPR Workshop on Machine Vision Applications*, pages 83-86, Tokyo, Japan, 1992.

15. Y. Leydier, F. LeBourgeois, and H. Emptoz. Text search for medieval manuscript images. *Pattern Recognition*, 40:3552-3567, 2007.

16. Y. Leydier, A. Ouji, F. LeBourgeois, and H. Emptoz. Towards an omnilingual word retrieval system for ancient manuscripts. *Pattern Recognition*, 42(9):2089-2105, 2009.

17. L. Li, S. J. Lu, and C. L. Tan. A fast keyword-spotting technique. In *Ninth International Conference on Document Analysis and Recognition (ICDAR)*, pages 68-72, 2007.

18. J. Lladós and G. Sánchez. Indexing historical documents by word shape signatures. In *Ninth International Conference on Document Analysis and Recognition (ICDAR)*, pages 362-366, 2007.

19. D. G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 60(2):91-110, 2004.

20. A. Marcolino, V. Ramos, M. Ramalho, and J.R. Caldas Pinto. Line and word matching in old documents. In *Fifth IberoAmerican Symposium on Pattern Recognition (SIAPR)*, pages 123-135, 2000.

21. M. Meshesha and C. V. Jawahar. Matching word images for content-based retrieval from printed document images. *International Journal on Document Analysis and Recognition (IJDAR)*, 11(1):29-38, 2008.

22. A. Murugappan, B. Ramachandran, and P. Dhavachelvan. A survey of keyword spotting techniques for printed document images. *Artificial Intelligence Review*, 35(2):119-136, 2011.

23. P. Natarajan, I. Bazzi, Z. Lu, J. Makhoul, and R. M. Schwartz. Robust ocr of degraded documents. In *International Conference on Document Analysis and Recognition (ICDAR)*, pages 357-361, 1999.

24. T. M. Rath and M. Manmatha. Word spotting for historical documents. *International Journal on Document Analysis and Recognition (IJDAR)*, 9(24):139-152, 2007.

25. T. M. Rath and R. Manmatha. Features for word spotting in historical manuscripts. In *International Conference of Document Analysis and Recognition*, pages 218-222, 2003.

26. T. M. Rath and R. Manmatha. Word image matching using dynamic time warping. *Computer Vision and Pattern Recognition*, (2):521-527, 2003.

27. E. Rosten and T. Drummond. Machine learning for high-speed corner detection. In *European Conference on Computer Vision*, volume 1, pages 430-443, 2006.

28. M. Rusinol, D. Aldavert, R. Toledo, and J. Lladós. Browsing heterogeneous document collections by a segmentation-free word spotting method. In *11th International Conference on Document Analysis and Recognition (ICDAR'11)*, pages 63-67, China, 2011.

29. S. N. Srihari, Srinivasan H, C. Huang, and S. Shetty. Spotting words in latin, devanagari and arabic scripts. *Indian Journal of Artificial Intelligence*, 16(3):2-9, 2006.

30. Kengo Terasawa and Yuzuru Tanaka. Slit style hog feature for document image word spotting. In *10th International Conference on Document Analysis and Recognition (ICDAR'09)*, pages 116-120, 2009.

31. Carl von Eckartshausen. *Aufschlüsse zur Magie aus geprüften Erfahrungen über verborgene philosophische Wissenschaften und verdeckte Geheimnisse der Natur*. Bavarian State Library, 1778.

32. Ismet Zeki Yalniz and R. Manmatha. An efficient framework for searching text in noisy document images. In *Proceedings of Document Analysis Systems (DAS)*, pages 48-52, 2012.

33. B. Zhang, S. N. Srihari, and C. Huang. Word image retrieval using binary features. In *Document Recognition and Retrieval XI (SPIE)*, pages 45-53, 2004.

# Management of Network and Energy Resources in Cognitive and Self-Organizing Wireless Networks

Apostolos D. Kousaridas[*]

Department of Informatics and Telecommunications
National and Kapodistrian University of Athens

akousar@di.uoa.gr

**Abstract.** The reduction of the consumed energy in modern self-organizing communication systems in a dense urban environment is a challenging task that requires coordination in management operations for the most effective use of network resources. Configuration and performance optimization tasks affect energy consumption of specific components and energy-related metrics of different devices. We propose a novel approach for energy saving and resource management in a wireless urban environment. Central to our approach is the organization of WLAN access points into clusters to facilitate local management and coordination. In each cluster, a cluster head access point monitors the energy consumption changes during the transmission and reception, at both the access point and user equipment sides, and decides on the appropriate adaptation action. The energy consumption reduction and performance improvement attained under the proposed solutions, at both the network and the user equipment sides, is evaluated via simulation.

**Keywords:** network management, wireless networks, resource management, energy saving, self-organization, cognition

## 1    Dissertation Summary

Contrary to common belief, information and communication technologies contribute a significant portion both to world energy consumption (2-4%) and environmental pollution (2-2.5% of greenhouse gas) [1]. Wireless network energy efficiency plays a primary role in reducing the impact of communication systems on energy consumption and environmental pollution [2]. Apart from

---

[*] Dissertation Advisor: Lazaros Merakos, Professor

an environmental responsibility, energy saving is important for the reduction of communication networks operational costs.

In the last decade, there has been a continuous increase in the number of wireless access points (AP) installations (e.g., IEEE 802.11) in private and public places to cope with user mobility and capacity requirements for emerging and future Internet services. Such APs are often not part of the same administrative entity, and the configuration of their locations and operational features are not necessarily planned for the "network welfare". This unstructured network environment results in dense AP topologies, especially in urban areas, with high coverage and frequency overlapping. The above in conjunction with users' varying traffic volume and service requirements, create optimization opportunities in energy saving and wireless resources.

The need to cope with complexity that derives from the interaction of hundreds or even thousands of network devices for the identification and realization of optimization opportunities calls for a distributed and localized solution. Self-organizing networks (SON) is considered as one of the most promising approaches for the management of networks that operate in highly dynamic and dense environments [3], [4].

For the deployment of a SON, each AP incorporates a Cognitive Network Manager (CNM), where energy saving as well as Coverage and Capacity Optimization (CCO) algorithms are placed [5]. In the context of this dissertation we present two types of a CNM a) simple CNM that is referred to as Network Element Cognitive Manager (NECM) and b) domain CNM entitled Network Domain Cognitive Manager (NDCM) [6]. NECM implements the cognitive cycle at the network element level, providing an intelligent adaptation layer to the conventional control plane. Management problems that cannot be addressed directly at the network element level, due to computational or communicational constraints, are escalated to the respective NDCM level. The NDCM incorporates the required cognitive capabilities to identify optimization opportunities and solve problems that require a greater view of network status. The distributed software architecture of cognitive managers is described [7], [8]. We implemented the distributed cognitive framework (software agents and artificial intelligence algorithms) that is deployed in access points and base station of a real heterogeneous access network composed of a Broadband Worldwide Interoperability for Microwave Access (WiMAX) BS and WiFi AP. Interference management and load balancing through channel reselection and vertical assisted handover algorithms respectively are the management tasks of the NECM and NDCM in this experimentation phase. Useful findings and the recommendations from the deployment of the cognitive network management architecture in a real life implementation are provided (av-

erage utilization of processing resources, memory usage, and delay of cognitive cycle phases) [9].

Central to our approach is the organization of WLAN APs into clusters to facilitate local management and coordination. Clusters are organization structures used for the collaborative tackling of network management problems, and they are formed following a common known scheme. Clusters facilitate the cooperation and the coordination of a group of network nodes for identifying and solving network management problems. In the literature, there are several clustering algorithms, which are mainly targeting wireless sensor networks or mobile ad-hoc networks, but the majority of them are application-specific (e.g., energy-efficient, mobility-aware). We propose SYSTAS algorithm, for the distributed discovery and establishment of clusters among network nodes, based on the features of the physical network topology. The density of the network graph and the preferential attachment model are used in order to form logical topologies [10]. The application of the proposed algorithm leads to the election of the head and the specification of the borders of the clusters through the allocation of the member nodes to the elected heads. The number of elected heads defines the number of the formed clusters. Clusters are non-overlapping and consist of two types of nodes:

- Simple member node
- Head node.

The head node of each cluster has the role of a NDCM, while simple member nodes instantiate the NECM. Both types of nodes implement CCO and energy saving management tasks. The simulation results, using various network graphs, show the effective cluster formation and the resulted high modularity.

In each cluster, the elected cluster head monitors usage of resources as well as the energy consumption changes and decides on the appropriate adaptation action (Fig. 1). In this dissertation, we propose a novel approach for energy saving and wireless resources management in a WLAN urban environment, where dependencies among different types of nodes and components are taken into account. CCO adapts network connectivity at all desired locations and provides bandwidth according to the communication needs of the users, avoiding the overutilization and underutilization of network resources.

The Capacity Usage Ratio (CUR) of a cluster network area with n APs is defined as the fraction of the available capacity that is actually being used:

$$CUR = \frac{\sum_{i=1}^{n} C_i}{\sum_{i=1}^{n} C_i^{\max}} \tag{1}$$

where $C_i$ and $C_i^{\max}$ is the used (uplink and downlink) capacity and the maximum available (uplink and downlink) capacity, respectively, of AP *i*.

For the calculation of the degree of coverage overlap in a cluster we introduce the overlapping factor (OF), which is based on the clustering coefficient (CC) [11]. The correlation of the CUR with the OF of the APs in a cluster area allows for more effective interpretation of the information that CUR provides, by taking into account the overlap level of the offered bandwidth. For this reason we use the composite metric of Coverage Optimization Opportunity (COOP) introduced in [12]:

$$COOP = CUR^{OF}$$

(2)

The COOP metric is useful for the identification of optimization opportunities for low load situations, where less capacity needed, as well for high load situations, where more capacity is required. A low COOP value means that too much capacity is provided in a very dense area, while a too high COOP value indicates an overloaded network area, where more resources are needed.

Energy consumption is measured at both the AP and the User Equipment (UE) side, focusing on the communication component (transmission, reception). CCO is applied via a novel scheme for the dynamic deactivation or reactivation of APs. This scheme aims at the rational usage of the radio resources according to traffic intensity and network density. The mechanism for UE load balancing after the de(re)-activation of an AP is also provided. The effect of an AP deactivation on UE and other APs energy consumption is assessed, triggering an additional adaptation action in the case that an energy efficiency problem has been detected. A scheme for multi-hop relay communication mode is proposed for the energy saving of UE transmission phase, exploiting local networking opportunities [13], [14], [15]. In addition, a novel channel reallocation scheme is introduced for the reduction of energy consumption during data reception phase. CCO and energy saving algorithms in a SON WLAN have been evaluated using OPNET simulation environment.

A novel scheme for the management of the interactions of various optimization or configuration problems in a SON is proposed. We identify and resolve conflicts on metrics and parameters (i.e., configuration actions), which arise from the deduction phase of the cognitive managers that are placed in a SON, in the context of the same or neighboring devices. The proposed scheme consists of three steps. Firstly, a time series-based mechanisms is used in order to avoid checking a configuration action (adaptations) that is triggered by a performance metric, which value appears continuous variations due to

temporary changes of the network area. This phase helps a SON to act proactively on conflicts and dependencies resolving, avoiding the trigger of adaptations that appear high uncertainty. In the next step, we check for conflicts on the triggered configurations actions. Only one "direction" per configuration action is prioritized, according to the severity and the priority of the performance metrics that trigger the corresponding configuration action. Finally, the impact of non-conflicting configuration actions on the other performance metrics is analyzed using a cost-benefit analysis scheme. The goal is to select the highest priority configuration action that creates fewer conflicts among available configuration actions and has the minimum possibility to deteriorate other (high priority) performance metrics. The proposed scheme for the coordination of various SON problems has been tested using OPNET simulation environment addressing CCO, energy saving, and interference tasks.
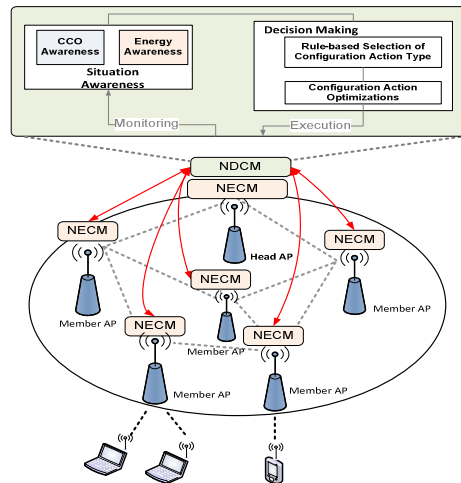


**Fig. 1.** Cluster head management tasks

In addition, an algorithmic framework for the extension of cognitive capabilities in network management has been described, facilitating performance management tasks. This framework is used for the improvement of Voice over IP (VoIP) QoS in a congested WiMAX network [16]. Despite the WiMAX related introduction, the proposed algorithmic framework solution is not access technology specific, but is equally feasible to other wireless network technologies as well, such as WLAN. The proposed algorithmic framework consists of the decision making, the execution and the learning phase. The decision making part includes the scheme for the identification of the most appropriate action for the packet loss reduction of VoIP service; selecting between a) the change of VoIP flows priority at the WiMAX base station, exploiting Medium Access

Control (MAC) features, and b) the change of VoIP flows selected codec, exploiting service level features. The solution and the quantification of the derived action (e.g., number of VoIP flows, the type of codec transition) is achieved using either a history-based scheme that takes advantage of previous events, or a heuristic approach for un-classified (i.e., unknown) situations. A k-Means learning algorithm is introduced to process the accumulated knowledge from all applied actions and evolve the decision making scheme [17], [18]. The performance and feasibility evaluation of the proposed solution has been tested using FIRE Panlab WiMAX experimental facility.

## 2 Results and Discussion

The proposed solutions for the effective utilization of network resources and energy saving have been deployed evaluated using both a simulation environment as well as real WiFi/WiMAX infrastructure.

### 2.1 SYSTAS: Algorithm for Cluster-based Structure of a Self-Organizing Wireless Network

The distributed algorithm for the organization of APs control loops in clusters has been evaluated in different topologies (sparse, dense, and real). In the literature, there are several clustering algorithms, which are mainly targeting wireless sensor networks or mobile ad-hoc networks, but the majority of them is application-specific (e.g., energy-efficient, mobility-aware). The results show efficient discovery of clusters and resulted modularity [19], comparing with algorithms from the area of data mining and graph clustering algorithms.
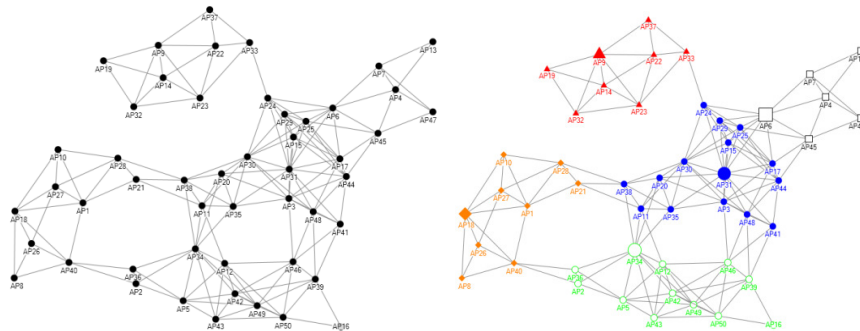


**Fig. 2.** Sample Topology of 50 Nodes (a) Graph visualization, (b) Formed clusters visualization

## 2.2 Energy Saving and Efficient Utilization of Wireless Resources in a cluster-based Self-Organizing Network

After the formation of the cluster structures and the activation of NECM and NDCM, the head AP retrieves topological, performance, and energy consumption information from the member APs. Moreover, the head receives monitoring data that the associated UEs provide to the All these data allow the head to build its situation awareness, which includes energy and CCO awareness of the cluster. They both are the input for the decision-making phase, which consists of two steps. Firstly, the head uses a rule-based scheme in order to evaluate the existing the situation awareness and select the appropriate adaptation. Then, the deduced configuration is resolved.
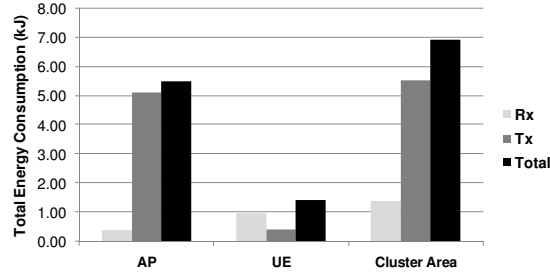


**Fig. 3.** Total energy consumption – Disabled energy saving and CCO

In this dissertation, we have focused on the communication component (Rx, Tx) of both APs and UEs in a dense WLAN network (Fig. 3) [20]. The decomposition of energy consumption into device and component levels facilitates the analysis and the identification of their dependencies.

From the obtained simulation results it is evident that coverage and capacity optimization is a tool for achieving energy efficiency. However, the extend of energy reduction as well as the impact on other system components (Rx, Tx) or nodes depends on the specific network topology, the network configuration features and traffic conditions. The simulation results for the selected topology configuration, show that an AP deactivation action, reduces cluster level $EC_{AP}$, mainly due to the reduced energy spent for packets reception (13% decrease) and processing. However, it increases the energy consumption of UEs for the reception and the transmission phase ($EC_{UE}^{Tx}$, $EC_{UE}^{Rx}$) under specific topology and traffic conditions (e.g., high overlap of frequency channels, UL/DL ratio). On the other hand, the consumed energy of UEs for the reception and the transmission phases increases after APs deactivation. The change of the selected channels in the cluster leads to the reduction of the

energy that UEs consume for the sensing of data packets that come from neighboring cells. APs gain also benefit from this adaptation. $EC_{UE}^{Rx}$ and $EC_{AP}^{Rx}$ measured in nJ/bit are improved by 35.7% and 40.5%, correspondingly. Furthermore, the handover of a UE to a more distant AP leads to an increase of the energy used for transmissions ($EC_{UE}^{Tx}$), especially for a UE with high UL traffic. In this case, the formation of UEs multi-hop relays improves the energy consumption in the data packet transmission phase by 20% (Fig. 4).
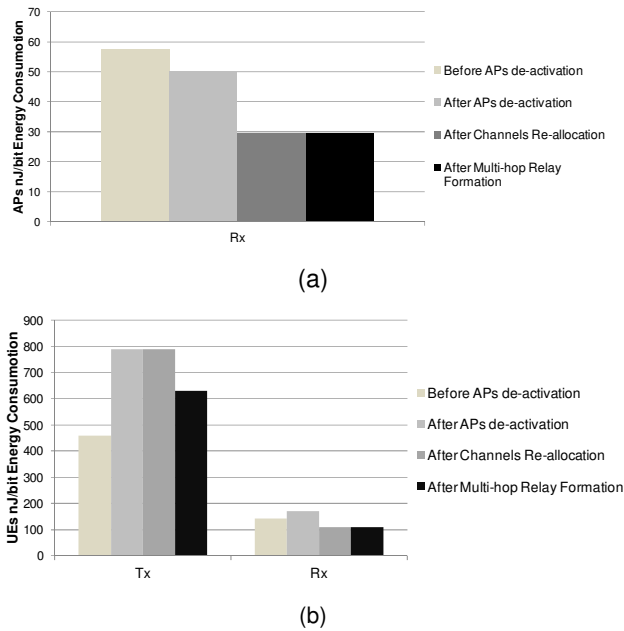


(a)



(b)

**Fig. 4.** Energy consumption in nJoules/bit (a) APs Rx, (b) UEs Tx and Rx phases

## 2.3  Coordination of Conflicts and Dependencies in Self-Organizing Networks

In the conducted experiments that described above, the degradation of a performance metric, after the enforcement of a re-configuration is addressed by using an additional optimization action. However, in the case that two or more configurations actions are triggered concurrently it is important for the SON to identify and resolve conflicts on metrics or parameters, so as to assure the stability of the communication network. The goal of the introduced algorithmic

scheme is to select, solve and apply the most appropriate configuration action in a cluster area taking into account a) problems severity, b) the total performance improvement of the network nodes of the cluster, c) and the number of reconfigurations or system's oscillations. Simulation results show that the proposed scheme for the coordination of the self-optimization functions of thye different cognitive cycles improve the performance of the system (throughput BER, network side energy consumption), according to the defined priorities.
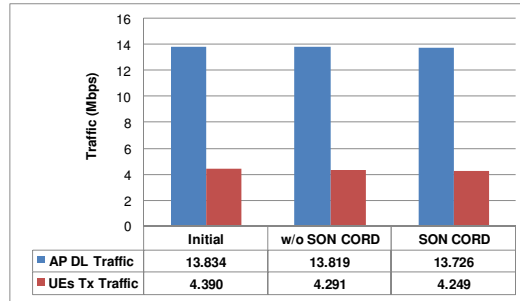
| | Initial | w/o SON CORD | SON CORD |
|---|---|---|---|
| AP DL Traffic | 13.834 | 13.819 | 13.726 |
| UEs Tx Traffic | 4.390 | 4.291 | 4.249 |

**Fig. 5.** Cluster-level throughput (UL/DL)

| | Initial | w/o SON CORD | SON CORD |
|---|---|---|---|
| BER | 0.0000910 | 0.0000250 | 0.0000097 |

**Fig. 6.** Cluster-level BER

| | Initial | w/o SON CORD | SON CORD |
|---|---|---|---|
| AP Tx | 534.000 | 534.000 | 534.000 |
| AP Rx | 29.033 | 23.091 | 12.060 |

**Fig. 7.** Cluster-level AP energy consumption (nJ/bit)

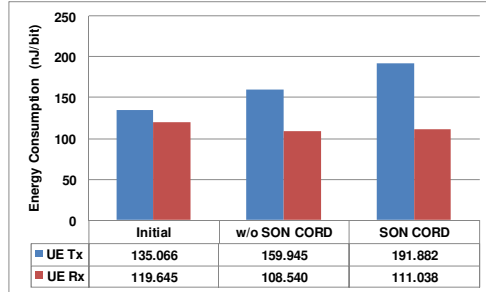| | Initial | w/o SON CORD | SON CORD |
|---|---|---|---|
| ■ UE Tx | 135.066 | 159.945 | 191.882 |
| ■ UE Rx | 119.645 | 108.540 | 111.038 |

**Fig. 8.** Cluster-level UE energy consumption (nJ/bit)

## 2.4 Cognitive Capabilities for a Service-aware Self-Managed Network

The proposed architecture and the introduced algorithmic framework for the extension of cognitive capabilities in a self-managed network system have been deployed and evaluated for the improvement of VoIP QoS (packet loss, delay, jitter) in a congested WiMAX network. The conducted experiments, using FIRE Panlab and CORE facilities, prove the feasibility and the strengths of our work. Both network and service side adaptation actions improve the detected packet loss level. However, the change priority of VoIP flows at WIMAX BS cannot reach the target PL threshold (<1%) for high PL events; although service continuity is always satisfied. On the other hand, the VoIP codec modification is more drastic (Fig 9.).
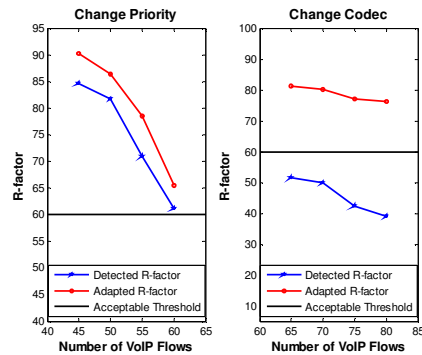


**Fig. 9.** R-Factor after Change Priority and Change Codec adaptation

The utilization of the history of previous adaptations reduces the transitory period and the iterations that a heuristic scheme requires for QoS improvement. The results of the learning phase show that the accuracy of the decision making scheme is improved, avoiding adaptations that are not effective. Final-

ly, we have showed that other QoS metrics such as delay, jitter and R-factor are also improved by the proposed solution (Fig. 10).
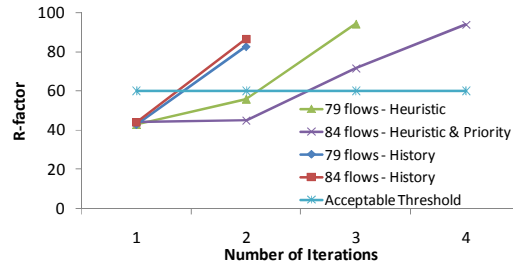


**Fig. 10.** R-factor vs. Number of Iterations

## 3    Conclusions

The efficient usage of network resources and the reduction of the consumed energy in modern communication systems that operate in a dense urban environment are challenging tasks, due to the complexity and the spatio-temporal dynamics of wireless networks. Self-organization is considered as one of the most promising paradigms for the management of networks that operate in highly dynamic and dense environments. In this thesis a novel approach has been proposed to dynamically control the size and configuration of a wireless network for the effective utilization of network resources and energy saving. The energy consumption reduction and performance improvement (BER, throughput, QoS) attained under the proposed solutions, at both the AP and the user equipment sides, is evaluated via simulation.

## 4    References

1. G. Fettweis and E. Zimmermann, "ICT Energy consumption – Trends and Challenges", in Proc. IEEE WPMC, Lapland, Finland, Sept.2008.

2. T. Chen, H. Zhang, Z. Zhao, X. Chen, "Towards green wireless access networks", in Proc. CHINACOM, Beijin, China, Aug. 2010, pp.1-6.

3. 3GPP TS 32.500, "Telecommunication management; Self-Organizing Networks (SON); Concepts and requirements", 2011.

4. C. Prehofer, and C. Bettstetter, "Self-Organization in Communication Networks: Principles and Design Paradigms", IEEE Commun. Mag, vol. 43, no. 7, pp. 78-85, Jul. 2005.

5. A. Kousaridas, C. Polychronopoulos, N. Alonistioti, A. Marikar, J. Mödeker, A. Mihailovic, G. Agapiou, I. Chochliouros, G. Heliotis, "Future Internet elements: cognition and self-management design issues", In Proc. ICST/ACM 2nd International Conference on Autonomic Computing and Communication Systems, 2008.

6. A. Kousaridas, N. Alonistioti, "On a Synergetic Architecture for Cognitive Adaptive Behavior of Future Communication Systems", In Proc. IEEE WoWMoM, 2008, pp. 1-7.

7. A. Kousaridas, G. Nguengang, J. Boite, et al., "An experimental path towards Self-Management for Future Internet Environments", in Towards the Future Internet - Emerging Trends from European Research, G. Tselentis, A. Galis, A. Gavras, et al., Eds, Netherlands: IOS Press, 2010, pp. 95 – 104.

8. A. Mihailovic, G. Nguengang, A. Kousaridas, M. Israel, V. Conan, et al., "An approach for designing cognitive self-managed Future Internet," In Proc. Future Network and Mobile Summit, Jun 2010, pp.1-9.

9. M. Bouet, G. Nguengang, V. Conan, A. Kousaridas, P. Spapis, N. Alonistioti, "Embedding Cognition in the Wireless Network Management: An experimental perspective" IEEE Communications Magazine, vol. 50, no. 12, pp. 150-160, 2012.

10. A. Kousaridas, N. Alonistioti, A. Mihailovic, "Dynamic compartment formation for coverage optimization of cognitive wireless networks", in Proc. IEEE PIMRC, Istanbul, 2010, pp. 2255-2260.

11. E. Schaeffer, "Graph clustering", J. Computer Science Review, vol. 1, no. 1, pp. 27-64, 2007.

12. A. Mihailovic, A. Kousaridas, A. Jaron, P. Pangalos, N. Alonistioti, H. A. Aghvami, "Self-Management of Future Access Networks for dynamic Configuration and Optimisation" Springer Wireless Personal Communications, 2013.

13. A. Kousaridas and N. Alonistioti, "Self-Organizing Cognitive Network Elements for Next Generation Communication System", In Proc. NAEC, 2009.

14. A. Kousaridas and N. Alonistioti, "Topology Control in self-Managed Wireless Networks", In Proc. MOBILIGHT, 2010.

15. A. Kousaridas, G. Katsikas, N. Alonistioti, E. Piri, M. Palola and J. Makinen (2011), "Testing End-to-End Self-Management in a Wireless Future Internet Environment", In. John Domingue, Alex Galis, Anastasius Gavras, Theodore Zahariadis, Dave Lambert. The Future Internet (pp. 259-270), Berlin, Heidelberg: Springer-Verlag.

16. P. Magdalinos, A. Kousaridas, P. Spapis, G. Katsikas, N. Alonistioti, "Enhancing a Fuzzy Logic Inference Engine through Machine Learning for a Self-Managed Network", Springer MONET 16(4), pp. 475-489, 2011.

17. P. Magdalinos, A. Kousaridas, P. Spapis, G.Katsikas, N. Alonistioti, "Feedback-based Learning for Self-Managed Network Elements", IEEE/IFIP International Symposium on Integrated Networks Management, pp. 666-669, 2011.

18. A. Kousaridas, A. Kaloxylos, et al, "Integrating the Self-Growing Concept in Self-Organizing Networks: Topology Optimization using Motion Sensors", Accepted for publication to Wiley International Journal of Network Management, 2013.

19. M.E.J. Newman, M. Girvan, "Finding and evaluating community structure in networks", J. Phys. Rev. E 69, 026113, 2004.

20. W. Heinzelman, A. Chandrakasan, and H. Balakrishnan, "Energy-efficient communication protocol for wireless microsensor networks", in Proc. Annual Hawaii International Conference, 2000.

# Open Source Software: Management, Diffusion and Competition

Spyridoula Lakka*

National and Kapodistrian University of Athens,
Department of Informatics and Telecommunication
`lakka@di.uoa.gr`

**Abstract.** The aim of this thesis is to contribute to the Open Source Software (OSS) research by a comprehending study of the factors that determine OSS diffusion, as well as the economic and social impact of this diffusion. The research focuses on the process of diffusion over time and assesses cause and effect relationships of the phenomenon. Firstly, it identifies and assesses the factors determining the diffusion and sustainability of OSS (cause). Secondly, it examines the effects of this diffusion at an economic and socio-political level (effect). At the economic level, the changes in markets structure and dynamics as a result of the OSS diffusion are analyzed. At the socio-political level, the study focuses on the effects in eGovernment and education.

The study grounds its results on a number of conceptual models that are based on a theoretical background with elements from the theories of technology acceptance and diffusion of innovations (DOI), as well as from social and economic theories. The models' evaluation is performed with the aid of rigorous methodological frameworks of mathematics and econometrics. Results can be a valuable input for both research and practice. For research, they provide with more accurate, a-priori estimations of the diffusion rate and the market competition. They also provide with the assessment of technological, social, economic and institutional factors that determine the OSS technology diffusion. As a result, they can become useful tools for strategic planning and policy making, in a continuously evolving and competitive environment such as the ICT market.

## 1    Introduction

OSS is an alternative model of software production and use, where source code is open for inspection, modification and distribution. OSS technology has introduced an innovative model of software development, based on self-organized communities that are open for participation to both users and developers. OSS innovation is twofold. First, the innovative method of organization and management of human and technology resources of OSS communities. Second, the OSS philosophy of open participation and the values of collaboration and sharing. According to Von Hipel, OSS is an innovation with a different value creation model, in which value is an outcome of collective intellect achieved through the OSS community [1].

---

* Dissertation Advisor: Draculis Martakos, Associate Professor

Over the last years OSS has moved to mainstream, creating a rapidly evolving eco-system that provides with thousands of software solutions. From a managerial perspective, OSS offers critical advantages that have turned it to a competitive option to most organizations. As a result, most software production companies have reshaped their business models and strategies so as to include OSS development procedures. It can be deduced that OSS plays a critical role in the ICT markets. This, in turn, places research interest in OSS at a high level.

Although its technological aspects are the object of extensive research [2-6], many researchers have also stressed out the socio-economic changes caused by the emergence and rapid diffusion of OSS [1, 7, 8]. The multi-dimensional nature of OSS has attracted academics from different research fields, like software engineering, economics, sociology and even political economy [3]. Studies that have performed a thorough review on OSS [2-6], detected a severe gap in the literature concerning the diffusion process of OSS as well as the factors that underlie this diffusion. Also, very few studies have been found to examine the factors that determine existence and long-term sustainability of OSS products [9], which is closely related to its diffusion. This, in turn, calls for a comprehending analysis and study of the parameters that determine OSS diffusion and sustainability. Moreover, though the socio-economic aspects of OSS have been extensively discussed, important research questions like the impact of OSS on market structure and competition and the social effects of OSS, in relation to other open initiatives, like open government and education remain unresolved.

The aim of this thesis is to fill this gap in the literature and perform a holistic analysis on the OSS phenomenon in two aspects. First, to evaluate the factors that are critical to its diffusion and sustainability. Second, to evaluate the economic and socio-political effects of this diffusion. The methodologies and results of the studies implemented into the context of this thesis, are briefly presented in the following sections.

## 2    Sustainability of OSS

The study investigated factors that affect the long-term sustainability of OSS projects. For that purpose, the study thoroughly examined the development methodologies and processes of the SourceForge portal [10] and identified projects' characteristics that could affect sustainability. The empirical data were queried from a database that is provided by the University of Notre Dame (UND) for research purposes [11] and contains full SourceForge projects' activities. The proposed conceptual model is based on the Unified Theory of Acceptance and Use of Technology (UTAUT) [12] and the IS success model [13]. The model defines the metrics of the impacting factors based on project's characteristics and considers the performance of these metrics in three distinct time periods. The aim of using different time segments is the evaluation of the impact of the factors of one time period as a cause for the users' behavior on subsequent time periods. The model consists of a number of structural equations and evaluates the weight of impact of the different time-dependent factors on long-term sustainability. The methodology for the model's evaluation is the Structural Equation Modeling (SEM).

Results indicate that the ability to attract the users' interest initially and active users and developers in the next period are the two critical factors for a project's sustainability. The user's choice is influenced by the performance and productivity of the community. The final decision depends on the social influence exerted through dis-

cussions of community fora. The research has been submitted to the scientific journal 'European Journal of Information Systems (IS)' and is under review process.

## 3 OSS diffusion.

The research estimates the cross-country OSS diffusion and the factors that shape this diffusion at a country level. Taking into account the multidimensional nature of OSS, a new theoretical framework is proposed as a lens for the identification of the possible impacting factors. The framework consists of the theories of exogenous and endogenous growth and institutionalism. A country is conceptualized as a socio-economic system within which OSS growth occurs. The model is based on the idea that the forces of growth to an economic system comprise of institutional, endogenous and exogenous factors and is specified as:

$$OSS_{it} = F(X^{endog}, X^{exog}, X^{inst}) \tag{1}$$

Where $OSS_{it}$ is the OSS growth rate determined by the three vectors of factors relevant to endogenous growth ($X^{endog}$), to exogenous growth ($X^{exog}$) and institutional theories ($X^{inst}$), for each country $i$, at time $t$. In this sense, growth is not restricted to economic development, but includes social, institutional and technological aspects. Into this context, two studies were implemented.

The first study examines the OSS diffusion process at a country level, under the prism of DOI theory and a parameterized diffusion model proposed by Dekimpe et al. [14]. The factors shaping the diffusion were drawn from the theoretical framework described by (1). The parameterized diffusion model allows for the comparison of diffusion parameters across countries and the evaluation of variables that may affect diffusion at different stages of the diffusion process. The model parameters are estimated and the significance of each of the factor is evaluated by a means of non-linear regression methods.

Research results indicated that a country's innovation level and human capital have a significant effect in all stages of diffusion. An important finding is also the different impact of the factors, depending on the diffusion stage. Technological infrastructure, economic status and institutions are more significant at the initial stage of OSS diffusion, while a nation's human capital and innovation activity are more important for subsequent stages. The research has been submitted to the scientific journal 'Information Technology and Management' and is under review process.

### 3.1 Diffusion of OSS: the case of the Apache web server

Previous research in OSS diffusion is very limited [15, 16]. Taking a case study approach, this study focuses on the diffusion process of a well established OSS, that is, the Apache web server. The study proposes a theoretical framework that consists of the DOI theory and the socio-economic theories as presented in (1) and aims to: (i) Estimate and forecast the market saturation of the Apache web server and provide with critical information for the interpretation of the diffusion process, namely the current stage of the market with respect to its saturation level and its inflection point. (ii) Evaluate the impact of socio-economic, country-level factors that affect the Apache's market saturation in different economic environments.

The theoretical framework is illustrated in Fig. 1. As shown, the research consists of two distinct steps. Firstly, the estimation of Apache's diffusion and market satura-

tion is based on the mathematical modelling drawn from the DOI theory, by means of a dynamic diffusion model [17]. The model assumes a time-variant market saturation, able to reflect the rapid diffusion of the web servers market. Secondly, a number of factors are evaluated for their impact on the Apache's market saturation. The factors are drawn out of the theories of institutionalism and exogenous and endogenous growth, as described in (1). The study's results are published in [18].
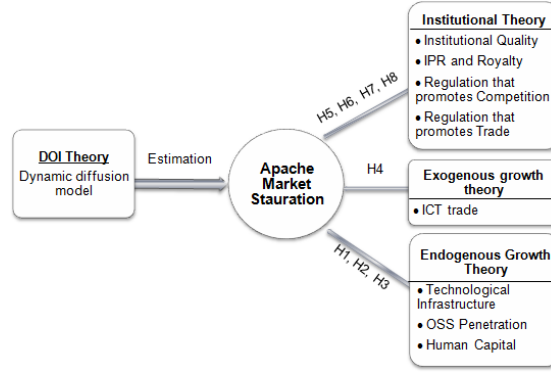


**Fig. 1.** Theoretical Framework

**Estimation and forecasting of the Apache web server diffusion .**

The diffusion process is estimated by the dynamic diffusion model proposed by Mahajan and Peterson [17]. The model assumes that the saturation level $\overline{N}(t)$, is not constant, but a function of time and can be expressed as $\overline{N}(t) = f(S(t))$, where $S(t)$ represents the vector of all relevant exogenous and endogenous factors affecting $\overline{N}(t)$. The proposed model takes into account the influence of one only factor, namely, the total market population growth (denoted by $P(t)=S(t)$) and assumes that the rate of increase in the market saturation with respect to the total market population, at any time t, is a constant. That is:

$$\frac{d\overline{N}(t)}{dP(t)} = k_2 \tag{2}$$

Integration of equation (2) yields (3), where $k_1$ is the integration constant and $k_2$ is the growth rate of market saturation with respect to total market population.

$$\overline{N}(t) = f\left(S(t)\right) = f\left(P(t)\right) = k_1 + k_2 P(t) \tag{3}$$

The final formulation of the dynamic model is given by:

$$\frac{dN(t)}{dt} = \left(a + bN(t)\right)\left(\overline{N}(t) - N(t)\right) = \left(a + bN(t)\right)\left(k_1 + k_2 P(t) - N(t)\right)$$

$$N(t = t_0) = N_0 = 0, \overline{N}(t_0) = f_0 \tag{4}$$

where N(t) refers to the cumulative number of adopters at time *t*. Also, *a* and *b* are the parameters of innovation and imitation of the diffusion process, respectively. $N_0$ represents the number of adopters at time $t_0$ and $f_0$ is the initially estimated saturation level at time $t_0$ .The solution of the differential (4) gives the number of adopters of Apache, at each point of time t. Moreover, *P(t)* can be estimated by the logistic diffusion model and thus is formed as follows:

$$\frac{dP(t)}{dt} = \left(m_1 + m_2 P(t)\right)\left(\overline{P} - P(t)\right),$$

$$P(t = t_0) = P_0 = 0 \tag{5}$$

Where $\bar{P}$ is the population saturation level and the parameters $m_1, m_2$ are the parameters of innovation and imitation. The parameters $\bar{P}$, $m_1, m_2$ are estimated using the Nonlinear Least Squares (NLS) regression method. The estimation of $a, b, k_1, k_2$ is based on the discrete regression analogue of (4), as shown in (6) and NLS. The regression coefficients $x_1$, $x_2$, $x_3$, $x_4$, $x_5$ give the estimates for the model parameters.

$$N(t+1) = x_1 + x_2 P(t) + x_3 N(t) + x_4 N(t)P(t) + x_5 N^2(t), \tag{6}$$

$$x_1 = ak_1, x_2 = ak_2, x_3 = k_1 b - a + 1, x_4 = k_2 b, x_5 = -b$$

The cumulative number of Apache web servers, N(t), are extracted from Netcraft's Web Server Survey [19]. The data span from the year 1996 to 2010 and are on a six months basis. The total population size, P(t), of the market corresponds to the total population of all possible web server adopters and its measure can be approximated by the number of Internet users. This can be justified by the fact that a potential web server adopter should firstly establish an Internet connection. The data for the Internet users penetration are derived from the United Nations (UN) database. The estimation of the parameters $m_1, m_2, a, b, k_1, k_2$ was derived by 24 observations, while the next 5 observations were used as a feedback sample to evaluate forecasting. To further evaluate the dynamic model's performance, the logistic model was used as a benchmark model. The evaluation of the estimation and forecasting is given in Table 1.

**Table 1.** Evaluation of the estimation and forecasting results

|  | Estimation | | Forecasting | |
|---|---|---|---|---|
|  | *Observations:* 24 | | *Observations:* 24 | |
|  | *Dynamic Model* | *Logistic Model* | *Dynamic Model* | *Logistic Model:* |
| *R²* | 0.98 | 0.98 | 0.988 | 0.81 |
| *MSE* | 9.4 | 9.21 | 15.4 | 78.3 |
| *MAPE* | 1.3 | 2.145 | 1.068 | 2.808 |

The above table shows that the statistical measures ($R^2$, MSE and MAPE) confirm the dynamic model's effectiveness in fitting, for both estimation and forecasting. It can be elicited that both the dynamic and logistic models effectively estimate diffusion, yet the dynamic model has a superior forecasting ability as compared to the logistic. This is mainly due to the time variant market potential, which captures the growth of the market and shifts the diffusion curve up to higher values. Graphically, the diffusion curves are illustrated in Fig. 2.

In addition, the inflection points of both the logistic and dynamic models were calculated. As shown in Fig. 2, there is big difference in the estimations with the logistic model showing that Apache has reached the inflection point at t*= 23 (year 2007), while the dynamic model is predicted to reach the inflection point at t*=38 (year 2014) and which is a much more realistic value. The shading parts of the figure correspond to the forecasted values of the diffusion. It can be deduced that the low value of the constant saturation level $\bar{N}$ (estimated at 112.81 million) of the logistic model shapes a downward slope that increases the gap between the last observations of the data. On the contrary, the time variant saturation level (estimated at 167.52 at t=38) of the dynamic model shifts the curves up resulting in much better fitting, especially for the forecasted values (time segments 25-29).This confirms the necessity for a non-constant market saturation for the rapidly diffused Apache web server. Findings suggest that Internet penetration has a positive impact on the market potential for Apache and that the growth of Apache has not yet reached its maximum rate (inflection

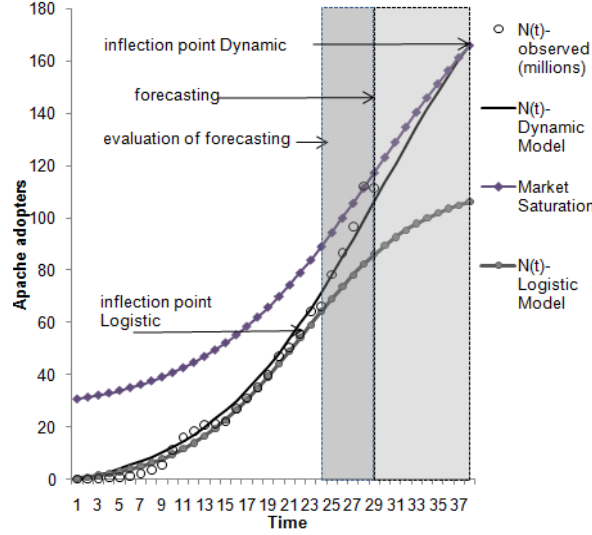point). As a result, the diffusion curve and the market potential are still at a growing stage.



**Fig. 2.** Diffusion curves and forecasting for Apache

**Socio-economic factors that determine the market potential for Apache.**

Taking into account the theoretical framework in (1), the model assumes that market saturation depends on institutional, endogenous and exogenous factors and is specified as:

$$\bar{N}^i(t) = F(X^{it}, Y^{it}, Z^{it}, \vartheta^i) \qquad (7)$$

where $X^{it}$ is a vector of all factors relevant to endogenous growth theory, $Y^{it}$ a vector of all factors relevant to exogenous growth theory, $Z^{it}$ a vector of all factors relevant to institutional theory, for each country i, at time t. In addition, $\vartheta^i$ a country specific variable, which determines developed and developing countries. Market saturation $\bar{N}(t)$ of the Apache web server is the dependent variable, explored in terms of possible influencing factors. The expected market saturation $\bar{N}^i(t)$ for each country i, can be estimated by equation (3), that is,

$$\bar{N}^i(t) = k_1 + k_2 P_i(t) \qquad (8)$$

where $P_i(t)$ is the number of Internet users for each country i and the parameters $k_1$, $k_2$ have been estimated with the dynamic model. By substituting the values of $P_i(t)$, $k_1$, $k_2$ in (8), the expected market saturation for each country is obtained. The factors are evaluated for their impact on the Apache market saturation, by means of a panel data analysis of 25 countries selected so as to represent different regions and economic status. The choice for the possible influencing factors was grounded on certain hypotheses. The statistical tests performed for the econometric model that derived out of (7) showed evidence of endogeneity, thus the Two Stage Least Squares (2SLS) regression with Generalized Method of Moments (GMM) and Heteroscedasticity and Autocorrelation (HAC) errors was the most effective method.

Regression results are presented in Table 2. It can be deduced that the diffusion of Apache depends on both endogenous and exogenous to a country factors, namely technological infrastructure, level of skills and education and ICT trade. The institutions and rules and laws of an organized society were also found to positively affect the growth of the Apache technology. On the contrary, regulation that promotes com-

petition does not appear to have any impact. Finally, Apache saturation levels are higher in developed versus developing countries, an outcome which is in accordance with other technological innovations and the problem of digital divide.

**Table 2.** 2SLS GMM regression results for $\bar{N}(t)$

| **Second regression: ln($\bar{N}(t)$)** | | | | **First regression: IQ** | | |
|---|---|---|---|---|---|---|
| **No of Observations:**112 | | | | **No of Observations:**112 | | |
| **F( 9, 102)**= 86.76*** | | | | **F( 11, 100)**= 66.87*** | | |
| **Vars** | **Coef.** | **S. Err.** | **z** | **Coef.** | **S. Err..** | **t** |
| *lnphone* | 0.570 | 0.049 | 11.72*** | | | |
| *ICTexp* | 0.085 | 0.035 | 2.43** | | | |
| *ICTtrade* | 0.015 | 0.003 | 4.32*** | | | |
| *HCI* | 0.011 | 0.007 | 1.91* | | | |
| *lnOSS* | 0.190 | 0.040 | 4.73*** | | | |
| *educ* | 0.088 | 0.034 | 2.58*** | | | |
| *B_R* | 0.021 | 0.023 | 0.91 | | | |
| *OECD* | -0.411 | 0.158 | -2.60*** | | | |
| *cons* | -5.901 | 0.785 | -7.52*** | | | |
| *IQ* | 0.280 | 0.131 | 2.14** | | | |
| *lnroyalty* | | | | 0.179 | 0.039 | 4.63*** |
| *TradeBar* | | | | 0.137 | 0.043 | 3.16** |
| *IPR* | | | | 0.094 | 0.023 | 4.08*** |

**Notes**: Significance levels: $* = p < .10$, $** = p < .05$ and $*** = p < .01$.

# 4    Economic impact of OSS diffusion.

The research focused on the impact of OSS on ICT markets structure and competition and consists of three studies. Initially, ICT market characteristics and dynamics were explored and analyzed in relation to the OSS special economic attributes that affect competition [20, 21]. The first study identifies and discusses the new dynamics formed in software markets due to the emergence of OSS. The impact of OSS in competition is further evaluated by applying the Herfindahl-Hirshman concentration index (HHI) on market shares data in three widely used software market segments, that is, web servers, web browsers and operating systems. HHI results indicated that though markets exhibit concentration, there are clear upward trends in competition. Market analysis showed that OSS has changed the strategies of most dominant software companies towards the creation of new OSS business models. This, in turn, has created new market entries and raised competition and market dynamics [22].

Further investigation and analysis of the OSS business models (OSS BM) was conducted in the second study. The objective of this study is to provide with a comprehensive and generic OSS BM framework that explicitly defines its structural elements, describing the deeper structure of what firms, adopting an OSS strategy, actually do. The study following the structured-case methodological approach [23] conducted two research cycles. In the first cycle, a sample of 100 popular OSS related firms instances is considered as 'pilots' organizations, in order to explore the different

possible business models cases. The instances were chosen so that to represent all three aspects of ICT markets, i.e. software, hardware and services sectors.

The second research cycle aims to validate, evaluate and further improve the initial findings by means of data collected from questionnaires and a workshop with eighty two participants, experts from the Greek OSS market and Academia. The research cycles revealed new concepts, dimensions and building blocks of the ontological OSS BM. The elements that differentiate OSS BM from the classical business models are the way of organizing production, the different OSS licenses, the innovative models of profit and the OSS community. The taxonomy of OSS BM was derived by a vertical analysis of the structure "*offered value*". Research results were published in [24].

### 4.1 Analysis of the operating systems market dynamics and competition

A deeper research in the effects of OSS in markets structure and equilibria was conducted in the field of the operating systems market. The study considers a highly concentrated market, the desktop (DT) and laptop (LP) operating systems sector, in order to provide some insights of the potential of OSS even in the case of high market concentration.

Based on concepts of population dynamics and organizational ecology, the study analyzes the evolutionary and competitive dynamics of the three leading players of the market, namely the OSS Linux, the partly-OSS Mac OSX and the proprietary Windows operating systems. Market evolution is estimated and forecasted by applying the Lotka-Volterra competition (LVC) model, which describes the competitive interaction of species for a common supply [25, 26]. The model's parameters were estimated by applying genetic algorithms, which are adaptive heuristic search algorithms based on the mechanisms of natural systems and genetics.

The main assumption of the methodology was to consider the three software products as interacting species competing for a common source, the market itself, expressed in terms of market shares. The empirical analysis showed that, at the equilibrium, all operating systems will coexist, while the highly concentrated market will tend to become less oligopolistic. Regarding the dynamics of the market, it was shown that Mac OSX has the highest growth rate and is less affected by the competition with the other systems, while is more affected by its own growth dynamics. An indirect mutualism effect, where the partly OSS Mac OSX is ultimately benefited by the existence of Linux, could also be deduced. Windows on the other hand experiences a decrease in its share, as it faces intense competitive pressures by both Linux and Mac OSX, with Linux being its main opponent. Results show that Linux shares are raised mainly due to Windows users that churn to Linux. However, the low growth rate of Linux is not expected to increase substantially, at least under the current market conditions and Windows will retain its leading position.

The above results add to the issue of the impact of OSS on competition. Firstly, OSS has a direct impact by the emergence of quality OSS such as Linux, which can offset the monopolistic behavior of the software market. Even in the highly concentrated DT/LP OS market, Linux not only survives but also raises its shares. Secondly, OSS allows for the creation of new business models, like the partly OSS that enable successful entrance in the market. The partly OSS Mac OSX paradigm shows that Mac OSX, though a late follower, has successfully entered a highly concentrated market.

As one step further, the study performs a sensitivity analysis of the possible effects on market behavior, induced by a rise in Linux adoption. Such a rise could be attributed to an organizational change of policy towards Linux adoption, as for instance in the public sector, following a governmental initiative. In this case, the Lotka-Volterra model is reformed to accommodate different adoption levels of the Linux operating system. Results demonstrate the effects of such policy on market concentration, according to different levels of Linux adoption. Findings also reveal useful implications for practice, in terms of the role of OSS and its derivative partly-OSS products in markets with high concentration. The main outcomes, which also define the importance of contribution of the proposed methodology, are the estimation of the modeled system dynamics, the provision of forecasts regarding market equilibrium and the estimation of the "churn effect", which reflects the level of users' switching among the operating systems. The model also provides information on the survival or extinction of each species due to the competition effects and the market structure at the equilibrium. he research has been published in [27].

## 5      Socio-political impact of OSS diffusion.

OSS ideology carries the notions and values of freedom, transparency and openness, active participation, cooperativeness and sharing that has created a new philosophical stream. The principles of OSS have extended beyond the software and inspired other forms of "open" initiatives, such as open standards, open access, open content, open science, open education, open government, open innovation and more. This thesis examines the impact of OSS in two sectors that have been highly affected by this openness: eGovernment and education.

### 5.1     OSS and eGovernment.

The relation of OSS diffusion and eGovernment maturity is examined at a national level, under the prism of the theoretical framework of the socio-economic theories of institutionalism, exogenous and endogenous growth, as presented in (1). The theoretical framework was deemed appropriate for two reasons. Firstly, because it was successfully applied for explaining OSS diffusion. Secondly, because the endogenous growth theories and institutionalism have been widely used in the case of eGov. Into this context, three distinct conceptual models were created and evaluated by means of econometric methods on secondary, cross-national data.

The two of the models focus on the impact of OSS on eGov maturity and use FGLS regression to statistically evaluate the corresponding models. Both models agree that OSS diffusion has a significant impact on eGov maturity. It can be concluded that OSS exerts a positive impact in eGov policies, like transparency, democracy and citizens' active participation, which are the characteristics of the eGov's higher maturity levels. The first model's results were published in [28], while the second's in [29].

Taking one step further, the third model investigated the simultaneity in the relation of OSS diffusion and eGov maturity. The model consists of two structural equations that express OSS and eGov mutual effects and is evaluated by means of a Simultaneous Equations Model and 2SLS regression. Results validated the positive impact of OSS diffusion on eGov maturity, yet rejected the assumption of a simultaneous relation. This research has been submitted to the scientific journal 'Technological Forecasting and Social Change' and is under review. Findings also provide with interest-

ing information on the impact of innovation on OSS diffusion and other country level factors that affect eGov maturity. These include institutional (governance effectiveness, freedom of the press, regulatory quality and the quality of institutions) and economic factors (ICT trade). Social development, however, exhibits the highest impact suggesting that higher levels of living and education are essential conditions for effective eGovernment.

### *5.2*    **OSS and education.**

In the education field, OSS has gained wide acceptance. A number of well known and established OSS communities and organizations have created software properly configured for education purposes. However, the success story for OSS, is the online education (E-learning) and the electronic Learning Management Systems (LMS) and Content Management Systems (CMS), with highly diffused software (e.g. Moodle, Sakai). A combination of open content and E-learning is the Open Educational Resources (OERs).

The study, reviewing the literature identifies factors that impact the diffusion of OSS in the education field, that is, the usual OSS advantages, like quality, cost-effectiveness and compatibility. However, there are some additional factors, like the ability for course customization by teachers and students, the low system requirements and the different levels of learning. Especially, important are the fast diffusion of knowledge through the OSS communities, the encouragement and promotion of collaborative and open way of learning, the promotion of open content and open education. The successful implementation of open content initiatives has lead to the development of a new education model, the open education. It can be concluded, that the role of OSS in the digitalization of education is twofold. Firstly, it provides with cost-effective, yet qualitative education software. Secondly, its philosophy sets the stage for new education models based on open standards, collaboration and active participation of both students and teachers.

## 6    Conclusions

The study contributes to the OSS research by the creation and development of new methodologies for the estimation and evaluation of the OSS diffusion process and the underlying critical factors, as well as the social and economic implications of this diffusion. Results, as explicitly described in the previous sections, revealed that OSS diffusion mainly depends on a country's technological infrastructure, innovation and education levels and social development. Also, the continuous attraction of users and developers in OSS communities, ensure its sustainability. For the economic impact, findings suggest that OSS plays a critical role on markets and competition, creating new business models and structures. All prediction models applied into the context of this thesis, were quite favorable to OSS, with steadily increasing trends of adoption. Finally, the positive relation of OSS and eGov verify its socio-political implications.

The research conducted for this thesis provides with useful input for both research and practice. For research, it brings in a new theoretical framework for the study of OSS diffusion that consists of three socio-economic theories: endogenous and exogenous growth theories and institutionalism. Into this context, it has developed methodologies for (i) the estimation of the international Apache diffusion and the factors that

impact market saturation, (ii) estimation of the cross-national OSS diffusion and the impacting factors at different stages of the diffusion process. It also implemented a conceptual model that is able to evaluate the critical factors for OSS sustainability.

Finally, it developed methodologies for the analysis and evaluation of the economic and socio-political impact of OSS. These include (i) the assessment of the OSS impact on market structure and competition, (ii) the creation of a competition model that estimates and forecasts the concentration, dynamics and the future market equilibrium of the operating systems market, (iii) the development of a holistic conceptual framework that provides with insights into the critical ontological elements of OSS BM and a taxonomy of the various OSS BM with an assessment of their risks and opportunities, (iv) the creation of models for the evaluation of the relation of OSS with eGovernment.

For practice, research results can become valuable input for strategic planning and policy making, as they provide with more accurate, a-priori estimations of the diffusion rate, the market competition and equilibrium, in a continuously evolving and competitive environment such as the ICT market. They also provide with the evaluation of the factors that constitute to the diffusion and sustainability of OSS, which are important information for organizations in designing their strategies. For the enterprises, which are planning to encompass OSS into their business models, the proposed OSS BM framework can become a useful tool. The conclusions for the relation of OSS diffusion and eGov maturity could also be taken into account at a political level, as they involve the assessment of technological, social, economic and institutional country level factors.

# References

1. Von Hippel, E. and Von Krogh, G.: Open Source Software and the "Private-Collective" Innovation Model: Issues for organization science. Organization Science, 14, 209-223 (2003)
2. Crowston, K., Wei, K., Howison, J. and Wiggins, A.: Free/Libre OSS Development: What We Know and What We Do Not Know. ACM Computing Surveys, 44, 1-33 (2012)
3. Aksulu, A. and Wade, M.: A Comprehensive Review and Synthesis of Open Source Research. Journal of the Association for Information Systems, 11, 576-656, (2010)
4. von Krogh, G. and von Hippel, E.: The Promise of Research on Open Source Software. Management Science, 52, 975-983 (2006)
5. Scacchi, W. Free/Open Source Software Development: Recent Research Results and Emerging Opportunities. In: ESEC/FSE '07, pp. 459-468. ACM, New York (2007)
6. Yang, J. and Wang, J. Review on Free and Open Source Software. In: IEEE International Conference on Service Operations and Logistics, and Informatics (IEEE/SOLI), 1, pp. 1044 - 1049. IEEE, Beijing (2008)
7. Lerner, J. and Tirole, J.: The simple economics of the Open Source. Journal of Industrial Economics, 52, 197-234 (2000)
8. Ågerfalk, P. J., Deverell, A., Fitzgerald, B. and Morgan, L. Assessing the Role of OSS in the European Secondary Software Sector: A Voice from Industry. In: 1[st] International Conference on Open Source Systems, pp. 82-87, Genova (2005)

9. Chengalur, I., Sidorova, A. and Daniel, S.: Sustainability of Free/Libre Open Source Projects: A Longitudinal Study. Journal of the Association for Information Systems, 11, 657-683 (2010)
10. SourceForge, `http://sourceforge.net/`
11. OSS research portal. University of Notre Dame, `http://zerlot.cse.nd.edu`
12. Venkatesh, V., Morris, M., Davis, G. and Davis, F.: User acceptance of information technology: toward a unified view. MIS Quarterly, 27, 425-478 (2003)
13. DeLone W. H. and R., M. E.: The DeLone and McLean model of information systems success: a ten-year update. Management Information Systems, 19, 9-30 (2003)
14. Dekimpe, M. G., Parker, P. M. and Sarvary, M.: Staged Estimation of International Diffusion Models. Tech. Forecasting and Social Change, 57, 105-132 (1998)
15. Whitmore, A., Choi, N. and Arzrumtsyan, A.: OSS: The Role of Marketing in the Diffusion of Innovation. Information Technology and Control, 38, 91-101 (2009)
16. Gallego, M. D., Luna, P. and Bueno, S.: Designing a Forecasting Analysis to Understand the Diffusion of Open Source Software in the Year 2010. Technological Forecasting and Social Change, 75, 672-686 (2008)
17. Mahajan, V. and Peterson, R. A.: Innovation Diffusion in a Dynamic Potential Adopter Population. Management Science, 24, 1589-1597 (1978)
18. **Lakka, S.**, Michalakelis, C., Varoutas, D. and Martakos, D.: Exploring the determinants of the OSS market potential: The case of the Apache web server. Telecomunications Policy, 36, 51-68 (2012)
19. Netcraft, The Apache Web Server Survey, www.netcraft.com
20. **Lakka, S**., Lionis, N. and Varoutas, D. Social Aspects of Open Source Software: Motivation, Organization, and Economics. In: Lee, C.P. (eds.), Electronic Business: Concepts, Methodologies, Tools, and Applications, pp. 1709-1722, IGI Global, Hershey, 2009.
21. **Lakka, S.**, Michalakelis, C. and Martakos, D. Impact of OSS on Social and Economic Welfare. In: Brooke, H. S., Scott, E. C. (eds.), handbook on Social Change, Nova Science Publishers 2009.
22. **Lakka, S**., Varoutas, D. and Martakos, D. Impact of OSS on Software Markets-an Evaluation. In: 4th MCIS, paper 56, pp. 588-599. AUEB, Athens (2009)
23. Carroll, J. and Swatman, P.: Structured-case: a methodological framework for building theory in information systems research. European Journal of Information Systems, 9 235-242 (2000)
24. **Lakka, S.**, Stamati, T., Michalakelis, C. and Martakos, D.: The ontology of OSS Business Model: an exploratory study. Int. Journal of Open Source Software and Processes, 3, 39-59 (2011)
25. Murrey, J. D. Mathematical Biology. Springer, New York, (2002).
26. Neal, D. Introduction to Population Biology. Cambridge University Press, New York, (2004).
27. **Lakka, S**., Michalakelis, C., Varoutas, D. and Martakos, D.: Competitive dynamics in the operating systems market: Modeling and policy implications. Technological Forecasting and Social Change, 80, 88-105 (2013)
28. **Lakka, S**., Stamati, T. and Martakos, D. Does OSS Affect E-Government Growth? An Econometric Analysis on the Impacting Factors. In:8th IFIP WG 2.13 International Conference OSS, Hammamet, Tunisia, 378, pp. 292–297 (2012).
29. **Lakka, S**., Stamati, T., Michalakelis, C. and Martakos, D.: What drives eGovernment growth? An econometric analysis on the impacting factors. International Journal of Electronic Governance 6, 20-36 (2013)

# Design, Fabrication and Characterization of a vibrational-driven piezoelectric microgenerator

Georgios P. Niarchos[1,2]

[1]Department of Informatics and Telecommunications

National and Kapodistrian University of Athens

[2]Institute of Microelectronics

National Center for Scientific Research "Demokritos"

gniarchos@gmail.com

**Abstract**. The research conducted during this thesis involves the "*Design, Fabrication and Characterization of a Vibrational Piezoelectric Microgenerator*", with use in wireless sensor networks, to make them autonomous. At the beginning, there is a review on some the latest cantileverd-based MEMS piezoelectric microgenerators and their characteristics. Also, there is a small review on the field of Piezotronics and its applications. The prototype microgenerator developed and presented in this thesis is a combination of MEMS technology and nanostructures of the piezoelectric material ZnO, able to convert low vibrations (~100Hz) into electricity. After, there is an extensive review on the basic equations that have been taken into account for the optimization of the design and fabrication processes. Arrays of vertical nanowires as well as uniform nanotextured films of ZnO were fabricated, simulated and characterized, taking into account the influence of each parameter. The MEMS microgenerators were successfully fabricated, packaged and characterized to achieve optimum results. Finally, an alternative approach on the fabrication of flexible nanogenerators is presented. Nanogenerators with either Au or Al electrodes were fabricated on flexible substrates, providing power outputs up to 30nWatts on an external load of 2MΩ.

**Keywords**: Energy Harvesting, MEMS, ZnO nanorods, Hydrothermal Method, Microgenerators, Flexible Nanogenerators

## 1 Introduction

The recent developments in the field of wireless sensor networks has attracted much interest, enabling them for further applications such as temperature or pressure monitoring, detection of toxic chemicals or gases or positioning of people in commercial buildings. Developments also in the field of VLSI have resulted in low power sensor

nodes (~1mW). Continuous powering these nodes is critical because it directly affects their lifetime and proper function, making the common way of powering (batteries) quite undesirable. Also, the total volume of the sensor node is affected by the size of the battery. Therefore, there is a need to develop new and exciting technologies that provide the node with the ability to harvest energy from the environment, making it self-sustainable. In this thesis, we explored the energy harvested from mechanical vibrations in order to fabricate microgenerators, utilizing MEMS processes and piezoelectric nanostructures.

## 2 Related Work

In the past few years there have been a number of publications regarding the fabrication of MEMS piezoelectric microgenerators reporting interesting results. MEMS microgenerators consist of a cantilever beam, usually with an end proof mass, metal electrodes and piezoelectric films. A large number of the known MEMS microgenerators, however, use PZT [1-4], which even though it produces great voltages due to mechanical deformations, it is toxic and therefore, they can be no biological or environmental applications.

ZnO was chosen as the piezoelectric material in this thesis, mainly due to its ability to form in various nanostructures and the fact that it has no toxicity and poses no threat to either the environment or organisms. The originality of this work rests on the combination of MEMS processes and ZnO nanostructures to develop microgenerators to power low-consumption electronics.

## 3 ZnO Nanostructures

3.1 ZnO nanorods

There are a number of known techniques used to produce plain and complex ZnO nanostructures, however our own results were based on the hydrothermal method. ZnO nanorods have beed fabricated on various substrates, taking into account all the necessary parameters involved (temperature, pH, precursor's concentration, time of growth) and the effects each one them has when they change. Patterned growth of vertical ZnO nanorods has been achieved and a thorough statistical analysis has been performed to study the morphological characteristics of the resulting structures (Fig. 1a-b). The inset in Figure 1a shows the verticality, while the one in 1b the different ZnO nanorods we fabricated [5].

(a)

(b)

**Figure 1**: (a) Patterned growth of vertically aligned ZnO nanorods. The inset shows the verticality of the nanorods, (b) Statistical analysis of the morphology of the resulting structures. The inset shows the different types of the resulting hexagons.

Based on the theory of chapter 2, simulations of our ZnO nanorods have been performed, in order to calculate the piezoelectric potential drop when there were subjected on known mechanical deformations. An extensive and novel work was carried away, taking into account not only the morphological results above but also several electrodes' configurations, in order to achieve the maximum theoretical power output for a given area of growth. The results were well within the range of the already published values [6].
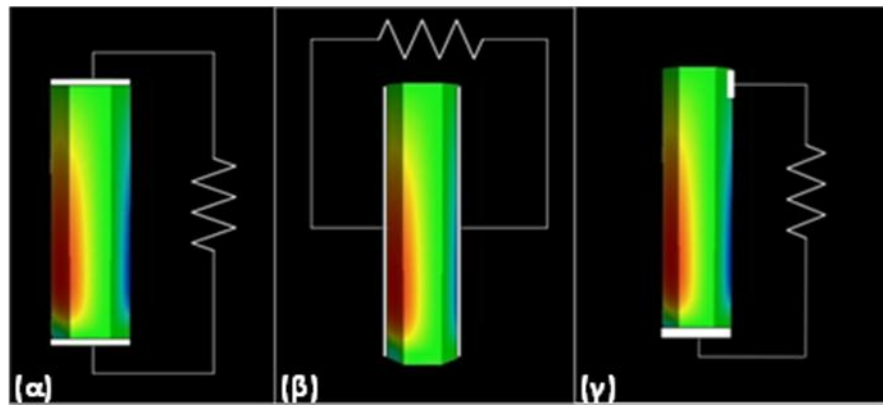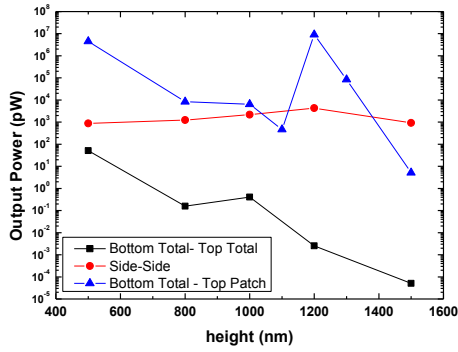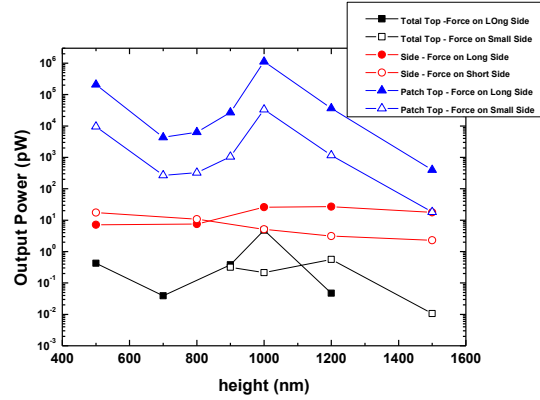


**Figure 2:** Three different topologies for the electrodes on the ZnO nanowire

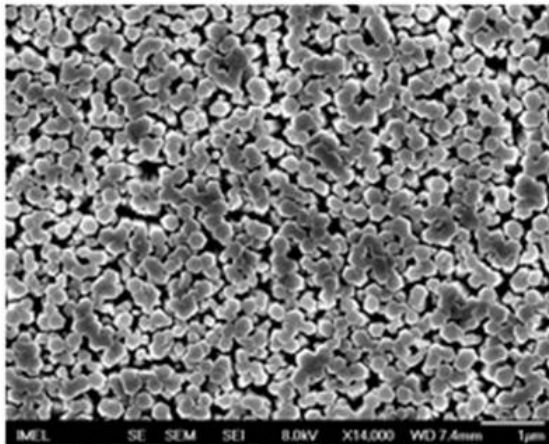**Graph 1**: Simulated power output for regular hexagons with a diameter of 100nm.

**Graph 2**: Simulated power output for elongated hexagons with a diameter of 300nm and the force at both the long and the small side

**Table 1:** Summary of the simulated power outputs based on the morphology of the nanorods and compared to known published values
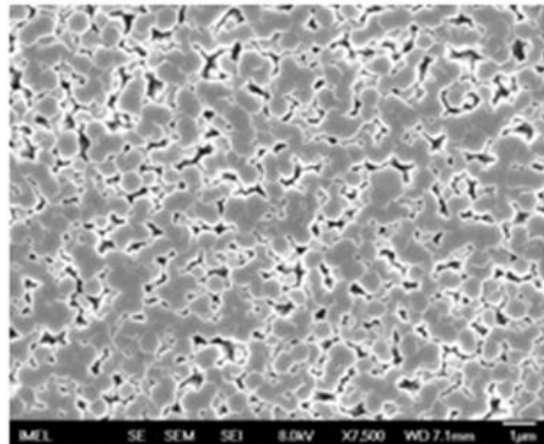
| Electrode Topology | Power (low aspect ratio NRs) | Power (high aspect ratio NRs) | Power measured from regular hexagonal NWs |
|---|---|---|---|
| Total Bottom – Total Top | 33.44 µW/cm$^2$ | 0.47 mW/cm$^2$ | 1 mW/cm$^2$ (Wang et al. 2006) |
| Total Bottom – Top Patch | 0.13 pW/cm$^2$ | **2.47 mW/cm$^2$** | |

## 3.2 ZnO nanotextured film

Further exploiting the hydrothermal method we successfully fabricated a uniform nanotextured ZnO film consisting of vertically aligned ZnO nanorods fused together. The purpose was to use this easy, low-cost method to fabricated uniform columnar films for use as active material in our microgenerators [7].
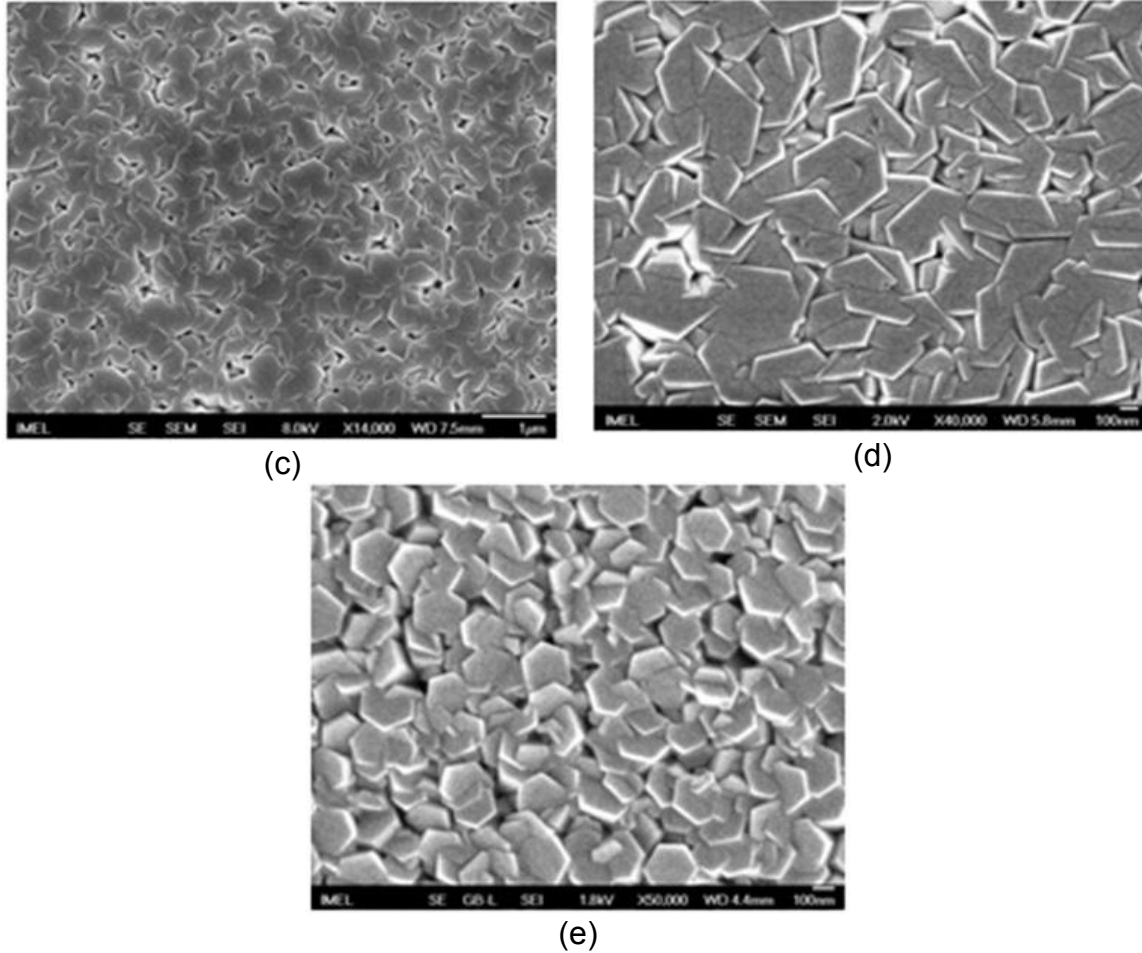




(a)

(b)

(c)



(d)



(e)

**Figure 3**: Fabrication of ZnO nanotextured films for concentrations of (a) 100mM, (b) 200mM, (c) 400mM, (d) 500mM, (e) 1000mM.


3.3 Characterization of ZnO nanotextured films

Further investigation in the fabrication of ZnO nanotextured films was carried away, in order to identify the optimum parameters and therefore produce the film with the best morphological and electrical characteristics. To that end, the nanotextured film was fabricated on two sets of electrodes: (a) Interdigitated and (b) bottom-top. I-V characteristics were performed in order to investigate the contact between the metal and the ZnO.
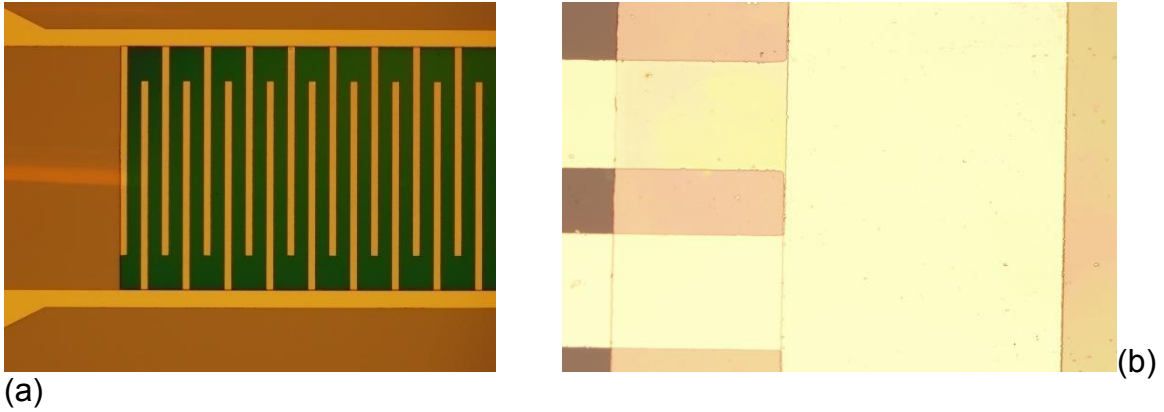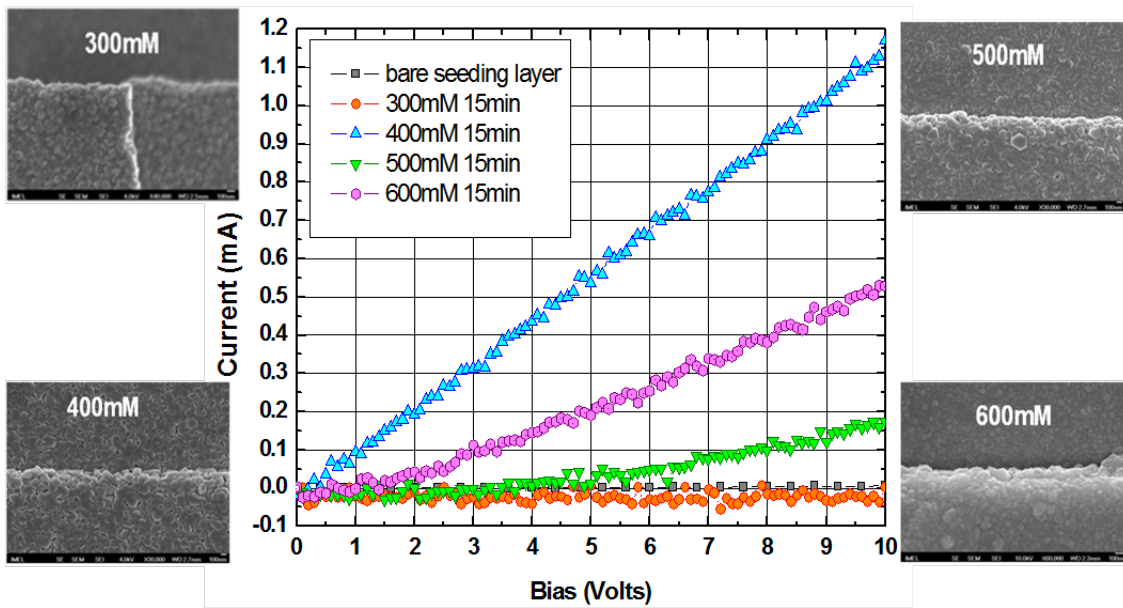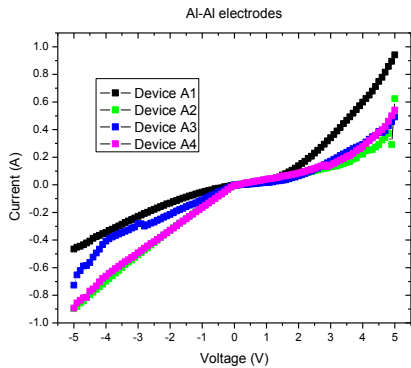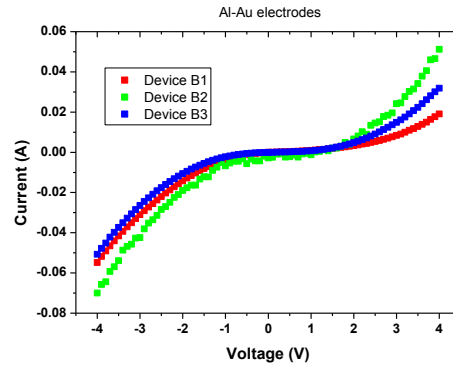
**Figure 4**: (a) Intedigitated electrodes, (b) Top-Bottom electrodes.



**Graph 3**: I-V characteristics on the IDE electrodes, based on the morphology of the resulting nanotextured film (inset SEM pictures).





**Graph 4**: I-V characteristics using bottom and top Al electrodes.

**Graph 5**: I-V characteristics using bottom Al electrode and top Au electrode.

Results indicated that the metal-ZnO contact, when using Al as electrode, is not Ohmic as expected but shows signs of Schottky diode. The microgenerator's principle of operating is based on the Shottky diode between the ZnO and the metal, so the preliminary results indicated that Aluminum can be used for making the electrodes in the microgenerator, making the full process CMOS compatible [8].

## 4 Fabrication – Characterization of MEMS microgenerators

We have successfully fabricated uniform columnar ZnO nanotextured films on metal substrates, as they will provide the electrodes for the generator. The proposed MEMS microgenerator is illustrated on Figure 5. It consists of a cantilever beam made of Si, with top-bottom metal electrodes, a piezoelectric film in between and an end proof mass. The dimensions of the cantilever and proof mass were specifically designed so the whole system resonates at low vibrations ~100Hz.



**Figure 5:** Illustrated schematic of the cantilever-based MEMS microgenerator

The novelty of our microgenerator lies in both the use of SOI technology and the ZnO nanotextured film. SOI wafers was preferred than standard Si ones, in order to simplify the process and make use of the buried oxide as an etch-stop during the etching for the release of the cantilever. The masks for the lithography were designed in AutoCAD and fabricated on Aluminum-coated glasses. Figure 6 shows the front and back side of an SOI half-etched wafer, where approximately 80 devices were fabricated.

(a)

(b)

**Figure 6:** (a) Front and (b) back side of SOI wafer with cantilever-based piezoelectric microgenerators.

A custom-made experimental setup was used to excite externally the fabricated microgenerators, using mechanical vibrations induced by 2 speakers. Using a frequency generator it was possible to control the vibrations and also measure the corresponding acceleration. Each die from the wafer included 4 separate microgenerators (for future parallel or series connection) and was wire-bonded on custom-made PCBs. The output signal was recorded in real-time using a digital oscilloscope. Figure 6 shows the experimental setup and a close look of the die.



(a)

(b)

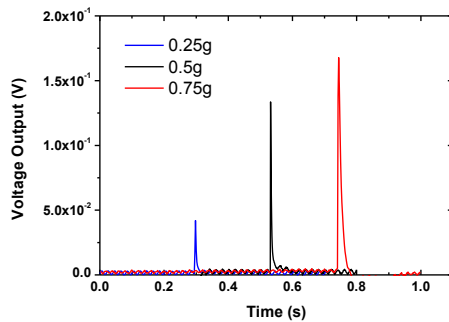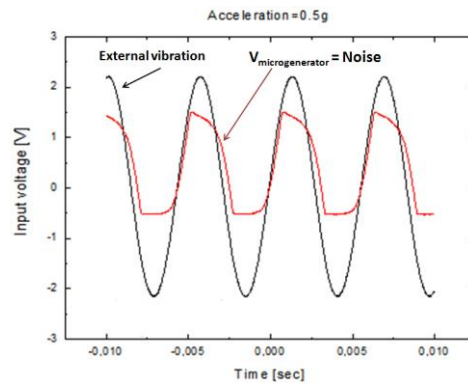**Figure 7:** (a) Wire-bonded die on PCB. Inset shows the size of the die. Inset shows the actual size of the die consisting of 4 microgenerators, (b) Experimental setup for piezoelectric characterization.



**Graph 6**: Voltage spikes of the MEMS microgenerators in different accelerations.

**Graph 7**: Voltage ouput in sinusoid excitation signal. The signal of the microgenerator is embedded in background noise.

Operation of the microgenerators was visible when the excitation was done with short strikes (spikes on Graph 6), but in the use of sinusoid signal there was only noise (Graph 7). Still, the proof-of-concept was proven. Further investigation of the properties of the ZnO nanotextured film was required, so we proceeded with the fabrication of flexible nanogenerators to study these parameters more thoroughly.

## 5 Flexible Nanogenerators

Through the hydrothermal method we were able to successfully fabricate nanogenerators on kapton and PET substrates with various architectures. This work was done in collaboration with Dr. Zhong Lin Wang from Georgia Tech [9]. The ZnO nanotextured film was used as the active material, while Au, Al and Pt were used as metal electrodes.


(a)

(b)

**Figure 8**: Flexible nanogenerators with (a) Au electrodes and (b) Al electrodes.

The characterization of these devices was performed using a techno-pneumatic piston to bend the nanogenerator, while the output signal was measured in real-time with a digital oscilloscope. The output voltage reached in several cases $V_{p-p}$=4 Volts with $I_{p-p}$= 800nA for controlled vibrations. Characterization of the devices occurred both in the labs of Georgia Tech and Institute of Microelectronics, providing interesting results.




**Figure 9:** Experimental setup in Georgia Tech.    **Graph 8:** Current output measured .

**Graph 9**: Voltage ouput vs acceleration for nanogenerators with Al electrodes.

**Graph 10**: Voltage output vs solution concentration for given acceleration and film thickness.

As indicated in Graphs 9 and 10, the voltage output is directly proportional to the acceleration and the morphological characteristics of the nanotextured film. Optimum concentration was 400mM with a growth time of 180min, resulting in film thickness of ~2μm. In order to have a direct comparison between the two different experimental setups, the nanogenerators with the Au electrodes were also measured in our lab giving the results presented in Graph 11.



**Graph 11**: (Top) Voltage output of the nanogenerators with Al electrodes and Au electrodes measured in our lab, (Bottow) Voltage output of the Au nanogenerators measured in Georgia Tech.

Graph 11 shows the comparison between the 2 experimental setups, where it is quite obvious that the voltage output is also depending on the manner in which the flexible substrate is bent. Higher voltages were measured in GaTech which gives room for further investigation and optimization of our own technique.

The next step was to simultaneously measure the voltage and current, using different resistive loads, in order to calculate the power output of such nanogenerators. Graph 12 shows a sample of these measurements, while in Graphs 13 and 14 the power is presented.



**Graph 12**: Voltage and Current ouput of an Al nanogenerator over a $R_L$=2MΩ for given acceleration.



**Graph 13**: Power vs Resistive load for different accelerations for an Au nanogenerator.

**Graph 14**: Power vs Resistive load for different accelerations for an Al nanogenerator.

The power ouput is directly proportional to the acceleration, as expected and we also observed an increase when using Au electrodes. However, the power generated from the Al electrodes is in the same range, which suggests that use of Aluminum as electrode is a safe and very promising choice [10].

## 6 Conclusions

We have presented an easy, low-cost, large-scale method for fabricating high-aspect ratio ZnO nanorods, as well as nanotextured films with high uniformity and promising electrical characteristics. For the first time, there has been fabrication of MEMS

microgenerators using the ZnO nanotextured film and the proof-of-concept was proved. Optimization of the properties of the nanotextured film occurred through the fabrication and characterization of flexible nanogenerators, generating power up to 30nWatts. Further investigation is required to optimize the electrical contacts and therefore the electric properties of the nanotextured films. The flexible approach has lead into interesting applications, such as wearable electronics or smart fibers/clothes.

## References

[1] Beeby SP, Tudor MJ, White NM (2006) Energy harvesting vibration sources for microsystem applications. Meas Sci Technol 17(12):175-195

[2] Choi WJ, Jeon Y, Jeong JH, Sood R, Kim SG (2006) Energy harvesting MEMS device based on thim film piezoelectric cantilevers, J Electroceram 17(2-4):543-548

[3] Fang HB, Liu JQ, Xu ZY, Dong L, Wang L, Chen D, Cai BC, Liu Y (2006) Fabrication and performance of MEMS-based piezoelectric power generator for vibration energy harvesting. Microelectron J 27(11):1280-1284

[4] Jeong S-J, Kim M-S, Songa J-S, Lee H-K (2008) Two layered piezoelectric bender device for micro-power generator Sens Actuators A Phys 148(1):158-167

[5] Niarchos G., Makarona E., Tsamis C., Growth of ZnO nanorods on patterned templates for efficient, large-area energy scavengers. Microsyst. Technol. 16, 669-675

[6] Niarchos G., Makarona E., Tsamis C. Modelling and optimization of ZnO nanostructure arrays for improved energy conversion efficiency. Oral Presentation at ICON 2009, Atlanta, Georgia, USA

[7] Makarona E., Fritz C., Niarchos G., Speliotis Th., Arapoyanni A., Tsamis C. Growth and characterization of uniform ZnO films as piezoelectric materials using a hydrothermal growth technique. Proceedings of SPIE, Volume 8066

[8] Niarchos G., Makarona E., Tsamis C., Low-cost ZnO nanorod arrays for nanogenerators of improved conversion efficiency, Poster Presentation at MNE 2010, Genoa, Italy

[9] Wang Z. L. Nanogenerators for Self-powered Devices and Systems, Georgia Institute of Technology, SMARTech digital repository, Atlanta, USA, 2011

[10] Niarchos G., Makarona E., Kyrasta Th., Voulazeris G., Speliotis Th., Tsamis C., Lin L., Hu Y., Wang Z. L. Comparison of ZnO-based Piezoelectric Nanogenerators on Flexible Substrates, "Hot-Poster" Presentation CIMTEC 2012, Montecantini Terme, Italy

# Document Image Binarization

Konstantinos Ntirogiannis[*]

[1] National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
`kntir@di.uoa.gr`

[2] National Centre for Scientific Research "Demokritos"
Institute of Informatics and Telecommunications
Computational Intelligence Laboratory
`kntir@iit.demokritos.gr`

**Abstract.** Principal stage of the document image analysis procedure is the binarization, according to which the pixels are classified into text and background. It is a crucial stage that can affect further stages including the final character recognition stage. This thesis is focused on document image binarization, including both binarization techniques and evaluation methodologies. Specifically, according to the developed performance evaluation methodologies, the pixel-level ground-truth image is constructed using a semi-automatic procedure based on the edges and the skeleton of the characters. The new measures use (a) weights that start from the ground truth contour and (b) the local stroke width to limit the weights close to the character areas and to properly normalize those weights. Experimental results prove the validity and effectiveness of the new measures for document images, while other measures concern the image or signal processing area in general. Concerning binarization techniques, some improvements were initially proposed for the well-known technique of Yang&Yan. To further enhance the quality of binarization and be more robust against different types of degradations (e.g. faint characters, bleed-through and non-uniform background), a new binarization technique was developed that was based on background estimation and on the combination of selected global and local binarization techniques. Additionally, a binarization technique was developed for the binarization of the text areas captured from video content. This technique is also based on the Yang&Yan binarization technique and sets low and high values in its global parameter for the inside and outside area of the text. Initially, the definition of the text areas is based on the baselines of the text and at the final stage the text areas are better defined by the convex hulls of neighbouring textual components. Furthermore, through the document image binarization contests that we organized, a publicly available benchmark has been created that aids in the development of document image binarization techniques and evaluation methodologies.

**Keywords:** pre-processing, binarization, evaluation metrics, ground-truth image, historical document image processing

---

# 1 Introduction

Document image binarization (or thresholding) is the process that segments the grayscale or color document image into text and background by removing any existing degradations (such as bleed-through, large ink stains, non-uniform illumination and faint characters). It is an important pre-processing step of the document image processing and analysis pipeline that affects further stages as well as the final Optical Character Recognition (OCR) stage. This thesis is focused on document image binarization, including both binarization techniques and evaluation methodologies. Our core motivation for the binarization was to develop an easy to tune method that could be effective against characters of various sizes [1], as well as against many different degradation types [2] (e.g. faint characters and bleed-through). Apart from the binarization of document images, we developed a method for the binarization of textual content from video frames [3].

As far as the developed evaluation methodologies are concerned [4, 5], we were motivated by the fact that existing pixel-based evaluation measures concern the image or signal processing area in general, while for document image processing those measures do not always provide reliable results. Last but not least, using the ground-truth construction procedure of our methodology [5], we successfully organized Document Image Binarization Competitions (DIBCO) from 2009 to 2012 [6–10] and made publicly available the competition datasets. Therefore, we have created a benchmark which is widely used for the development of document image binarization techniques and evaluation methodologies.

In the following, in Section 2 we present the related works concerning the binarization methods along with the binarization methods developed through this thesis, in Section 3 we present the related works concerning the evaluation methodologies along with the evaluation methods developed through this thesis. In Section 4 we present the experimental results and finally, in Section 5, the conclusions are drawn.

# 2 Binarization Methods

## 2.1 Related Work

Many document image binarization methods have been proposed which are usually classified in two main categories, namely global and local. Reference points in binarization are considered the global thresholding method of Otsu [11] and the local adaptive methods of Niblack [12] and Sauvola et. al. [13] which are widely incorporated in binarization methods that followed, e.g. Kim et. al. [14], Gatos et. al. [15], Lu et. al. [16]. Certain document image binarization methods have incorporated background estimation and normalization steps, e.g. Gatos et. al. [15], Lu et. al. [16], as well as local contrast computations to provide improved binarization results, e.g. Su et. al. [17], Howe [18]. Other binarization

methods, aiming in an increased binarization performance, proposed combination methodologies of binarization methods, e.g. Gatos et. al. [19], Su et. al. [20].

As far as the video frame binarization is concerned, there exist several techniques that perform binarization on the textual content in video frames aiming in an improved OCR performance. Many techniques [21–23] incorporate modifications of well-known binarization techniques, such as the Logical Level Technique of [24], Otsu [11] and Sauvola et. al. [13]. Other techniques [25, 26] are based on training using mainly SVM (Support Vector Machine) classifier or convolutional neural network. In the most recent related work [27], the Canny edge detector [28] was used to specify the text boundaries on the image. Then, a flood fill algorithm was used to fill the edge contour and form the characters. However, the Canny edges can be very confusing since they also depict non-text objects. Especially, in videos with high background complexity the edges of text may connect with background edges and hence deform the actual contour of the characters.

## 2.2 Improvement of Yang&Yan Method

The method of Yang&Yan [29] assumes a single stroke width for the document image. The value of the stroke width determines the size of the windows that are used to calculate the threshold at each point. However, characters of various sizes may exist within a document (e.g. a newspaper with big titles). To adaptively define the stroke width and consequently the size of the windows, we rely on the binarization output of [15]. Then, we detect the contour points and the skeleton using skeletonization method [30]. Afterwards, the local stroke width is assigned to each skeleton point by measuring the distance of that skeleton point from the nearest contour point. Then, each remaining point inherits the value of the nearest skeleton point found. However, for machine-printed documents that may suffer from internal holes at their strokes, the maximum of the local stroke widths is considered. All the aforementioned stages are shown in Fig. 1. Another improvement is the modification of the local threshold by a factor $\beta$ ($T' = \beta \cdot T$). According to [1], this factor enhances the overall performance, especially for machine-printed documents. Representative results are shown in Section 4.

## 2.3 A Combined Binarization Approach

In degraded historical images, faint characters and bleed-through have quite similar characteristics. Thus, current methods are usually robust against one of the aforementioned degradation. In [2], we introduced a binarization method capable of achieving high performance in many different noise types. The main idea is to initially erase all the noisy components (false alarms) even if faint character parts are also removed. Then, perform binarization of high *Recall* such as Niblack [12] and perform combination at connected component level. In this way, noise is erased, the faint characters are completely detected, while the noise

**Fig. 1.** The stages of the adaptive stroke width detection: (a) initial binary image; (b) contour points along with the skeleton; (c) local stroke width is assigned to each skeleton point; (d) the character stroke width image; (e) the final stroke width map (used for handwritten documents); (f) the final stroke width map (used for printed documents).

levels are very low. All the aforementioned stages are detailed below and are shown in Fig. 2:

1. Niblack binarization (w=60x60, $k$=-0.2) and one iteration of dilation (3x3 element),
2. estimate the background, follow proposed inpainting [2] using the above Niblack result as inpainting mask,
3. normalize original image with the above estimated background (keep the range of the original image),
4. Otsu binarization and remove connected components of very small height,
5. calculate: (a) the stroke width map using the above binary image and (b) the global contrast,
6. Niblack binarization with window size and parameter $k$ based on the stroke width map and the global contrast, respectively,
7. combination at connected component level. Large Niblack components that correspond to only a few foreground pixels of Otsu are not considered,
8. enhance the final result using binary image of step (4) (before the components removal).

### 2.4 Thresholding of Video Text Areas

For the binarization of video frames, we assume that the text detection step has already been performed and we focus on the binarization step of the detected text boxes. We introduced in [3] a binarization technique that aims in improving the text/background separation. The main idea is to specify the main body of the text (Fig. 3a-3c) in order to extract valuable information concerning the textual content. The main body of the text is defined as the area which is limited by the upper and lower baselines. Then, within the main body of the text we detect the stroke width (SW) of the characters which is used in consecutive adaptive binarization steps that follow. At a next step, we perform adaptive binarization [29] with different valuation in parameters for the inside and outside area of the main body of the text (Fig. 3d). Hence, we remove most of the non-text information but in certain cases it results in the thinning and breaking of the textual

**Fig. 2.** The stages of the combined binarization approach: (a) original image; (b) estimated background; (c) normalized image; (d) Otsu binarization of (c) and small components removal; (e) Niblack binarization; (f) final combined result.

parts that are outside the main text body. Afterwards, we define the entire text body as the region inside the convex hulls of continuous connected components (Fig. 3e) and we perform the same adaptive binarization with different valuation in parameters for the inside and outside area of the entire text body (Fig. 3f).

## 3 Binarization Evaluation Methods

### 3.1 Related Work

Several efforts have been presented that strive towards evaluating the performance of document image binarization techniques. These efforts can be classified in three main categories (the human-oriented, the OCR-based and the pixel-based).

In the first category, evaluation is performed by the visual inspection of one or many human evaluators [31, 32]. For example, in [31], the amount of symbols that are broken or blurred, the loss of objects and the noise in background and foreground are used as visual evaluation criteria. In the second category, evaluation is addressed taking into account the OCR performance. The binarization

**Fig. 3.** The stages of the binarization method for video text areas: (a) original image; (b) binarization using [1] to detect the baselines; (c) main body defined by the baselines; (d) binarization [29] of (c) along with the convex hulls of neighbouring components; (e) main body defined by the convex hulls; (f) final bibinarization.

outcome is subject to OCR and the corresponding result is evaluated with respect to character and word accuracy [15, 27]. In the third category, pixel-based evaluation is used by taking into account the pixel-to-pixel correspondence between the ground truth and the binarized image. In this category, the evaluation is based either on synthetic images [33, 34] or on real images [35]. Ground truth images from real degraded images which correspond to real "challenging" cases for document image binarization were not publicly available. The Document Image Binarization (DIBCO) contests that were organized by us [6–10] made the datasets publicly available after each corresponding contest.

Concerning pixel-based evaluation, several measures have been used for the evaluation of document image binarization techniques, such as the F-Measure (Recall and Precision), the PSNR, the Negative Rate Metric (NRM) and the Misclassification Penalty Metric (MPM) [6], the chi-square metric [36], the geometric-mean accuracy [34], the normalized cross-correlation metric [35] and the DRD (Distance Reciprocal Distortion) [37]. Some researchers have stated the need for an improved pixel-based evaluation measure for document image binarization. For instance, in [35], wherein the ground truth generation from several users was studied, it was stated that there is a need for a weighted measure in relation to the ground truth borders in order to compensate the subjectivity of the ground truth.

### 3.2 Skeleton based Methodology

This method was presented in [4]. It consists of a semi-automatic procedure for the ground-truth construction and it also introduces the use of the skeleton of the characters for the evaluation of binarization output in terms of "Recall". However, the ground-truth construction procedure has certain issues which were resolved in the latest evaluation methodology presented in [5]. Thus, in this section we will focus on the evaluation stage and not at the ground-truth construction procedure.

The main novelty of this method was the use of a skeletonized ground-truth to measure the performance of binarization in terms of *Recall*. Due to the ambiguity in the boundary of the characters, which is mainly created by the digitization process, binarization methods are penalized when boundary pixels are missing. However, the loss of pixels is much more significant when character breaking occurs. In more details, taking into account a historical document with faint characters (Fig. 4), F-Measure (FM) could rank in a better position a binarized image with more broken characters and false alarms as in Fig. 4b (FM=94.37) than a better binarized image as in Fig. 4c (FM=93.69). For Fig. 4c that contains less broken characters, higher Recall is expected than Fig. 4b. However, the binarized image of Fig. 4c achieves lower Recall=89.78 compared to the Recall=93.77 of Fig. 4b, as a result of the more missing foreground pixels (false negatives) which are mainly situated along the borders of the characters, making their absence less obvious.



(a)  (b)  (c)

**Fig. 4.** Deviation between quantitative and qualitative evaluation using F-Measure (FM): (a) original image; (b) binarized image with broken characters and false alarms, FM=94.37 (Recall=93.77); (c) better binarized image, FM=93.69 (Recall=89.78).

However, the use of the skeletonized ground truth for the computation of Recall provides better evaluation results. For Fig. 4b, false negatives corresponding to broken characters are taken into account ($FM_{skel}$=95.29, $Recall_{skel}$=95.62), while false negatives situated near the contour as in Fig. 4c, are not considered at all ($FM_{skel}$=98.79, $Recall_{skel}$=99.64). However, the dual representation of the ground truth could mislead the evaluation results when the binarized image is deformed while the skeletonized ground truth can be completely detected, as shown in Fig. 5. In those cases, both $Recall_{skel}$ and Precision are 100 ($FM_{skel}$=100), leading to erroneous evaluation. Thus, we have greatly modified this evaluation method, as described in the following section.

### 3.3 Weighted Recall/Precision Methodology

The character boundary ambiguity, as we discussed in the previous section, suggests that a distance-based metric would compensate those errors, since the use of the skeletonized ground-truth have certain limitations. However, there are a

**Fig. 5.** Problematic cases concerning the skeletonized ground-truth: (a) original image; (b) ground-truth image along with the skeletonized ground-truth; (c) binarization output wherein the skeletonized ground truth is fully detected.

few factors that should be considered when penalization weights are based on the distance from the ground-truth contour. These factors are listed below:

- a breaking at a small/thin character part would have much less penalty than a bigger/thicker character part;
- noise inserted among the characters would have much less penalty than attached to a single character;
- noise attached to a big/thick character is less important than attached to a smaller/thinner character;
- noise far from the ground-truth that do not interfere with the textual content would be much grater penalized than noise among the characters that destroys the useful textual content.

In [5], we proposed proper weighting to minimize/diminish the effects of the aforementioned factors. In particular, to measure the amount of loss, the pseudo-Recall was introduced by which the distance-based weights are normalized according to local stroke width. In this way, each character breaking has the same importance regardless of the local thickness. Additionally, to measure the amount of the inserted noise, the pseudo-Precision was introduced according to which the weights are constrained within an area that extends to the background by the corresponding stroke width of each character. In this area the weights take values from 1 to 2, while outside this area the weights equal one. In this way, noise that is located among the characters hs higher significance, while noise far from the ground-truth does not get exaggerating penalty. Furthermore, the distance between the characters is also considered to handle the cases of noise among the characters that result in merging.

The metrics of Recall and Precision are combined into F-Measure. Hence, the proposed pseudo-Recall/Precision are combined into pseudo-FMeasure $F_{ps}$. After many test cases examined in [5], the proposed pseudo-FMeasure offers more reliable results and it also has greater consistency to the OCR results. Representative results are given through Fig. 6 and Table 1.

## 4 Experimental Results

In this section the experimental results for the binarization of document images are shown. In the following, Fig. 7 shows representative results of the developed

**Fig. 6.** (a) Original image; (b) ground-truth; (c) binarization where text is preserved but with background noise; (d) binarization with stains among the text.

**Table 1.** Comparison between existing and the proposed ($F_{ps}$) pixel-based measure. The OCR accuracy is also shown. Notice that for metrics MPM and DRD lower values denote higher performance.

|          | FM        | PSNR      | MPM      | DRD      | $F_{ps}$  | OCR accuracy |
|----------|-----------|-----------|----------|----------|-----------|--------------|
| Fig. 4b  | **94.37** | **22.89** | 0.85     | 1.72     | 95.02     | -            |
| Fig. 4c  | 93.69     | 22.56     | **0.07** | **1.53** | **98.24** | -            |
| Fig. 6c  | 90.30     | 16.89     | 7.33     | 3.85     | **93.38** | 92.86        |
| Fig. 6d  | **91.48** | **17.37** | **0.66** | **3.39** | 92.37     | 87.50        |

methods [1] and [2]. In Table 2 the detailed evaluation results are shown using the winning method of each DIBCO competition [6–10] as well as results from the current state-of-the-art methods [17, 18] that used the same DIBCO datasets. From Table 2, it is shown that the latest method presented in [2] achieves the highest performance for the majority of the evaluation metrics.

## 5 Conclusions

Though this thesis, we have thoroughly studied the research area of document image binarization by focusing not only at the development of novel binarization techniques but also at the corresponding evaluation methods and metrics. An initial binarization method was developed that is more robust for machine-printed documents, while it has poor performance in handwritten images. The latest binarization method achieves high performance in documents with many different degradations types and it also achieves higher performance than state-of-the-art methods or methods from the DIBCO contests. Furthermore, the idea of using the baselines and the convex hulls for binarization purposes seems promising for the video processing area. Additionally, the initial evaluation methodology revealed some benefits of using a skeletonized ground-truth for evaluation purposes but it also revealed some drawbacks. The latest evaluation methodology was developed on the premise that the effect of flipped pixels on the image should be considered, and not just the fact that pixels had been flipped, which

**Fig. 7.** (a) Original image; (b) ground-truth; (c)-(d) Ntirogiannis et. al. [1] and [2], respectively; (e)-(f) Ntirogiannis et. al. [1] and [2], respectively for Fig. 6.

leads to more reliable document-oriented evaluation. Last but not least, using the ground-truth construction procedure of the latest evaluation methodology, we made ground-truth from real degraded images and organized international document image binarization competitions. The datasets were made publicly available after each competition and have been widely used ever since.

## References

1. Ntirogiannis, K., Gatos, B., Pratikakis, I.: A modified adaptive logical level binarization technique for historical document images. In: Proc. Int. Conf. on Document Analysis and Recognition. (2009) 1171–1175
2. Ntirogiannis, K., Gatos, B., Pratikakis, I.: A combined approach for the binarization of handwritten document images. Pattern Recognition Letters (2012) DOI: 10.1016/j.patrec.2012.09.026.
3. Ntirogiannis, K., Gatos, B., Pratikakis, I.: Binarization of textual content in video frames. In: Proc. Int. Conf. on Document Analysis and Recognition. (2011) 673–677
4. Ntirogiannis, K., Gatos, B., Pratikakis, I.: An objective evaluation methodology for document image binarization techniques. In: Proc. Int. Workshop on Document Analysis Systems. (2008) 217–224
5. Ntirogiannis, K., Gatos, B., Pratikakis, I.: Performance evaluation methodology for historical document image binarization. IEEE Transactions on Image Processing **22**(2) (2013) 595–609
6. Gatos, B., Ntirogiannis, K., Pratikakis, I.: ICDAR 2009 Document Image Binarization Contest - DIBCO 2009). In: Proc. Int. Conf. on Document Analysis and Recognition. (2009) 1375–1382
7. Gatos, B., Ntirogiannis, K., Pratikakis, I.: DIBCO 2009: Document Image Binarization Contest. International Journal on Document Analysis and Recognition **14**(1) (2011) 35–44

**Table 2.** Comparison of the proposed methods [1] and [2] to the winning method of each DIBCO contest as well as to methods [17] and [18].

| Method | FM | PSNR | NRM | MPM | DRD | FM$_{skel}$ |
|---|---|---|---|---|---|---|
| Winner DIBCO09 | 91.24 | 18.66 | 4.31 | 0.55 | - | - |
| Su [17] | 93.50 | 19.65 | 3.74 | **0.43** | - | - |
| Ntirogiannis [1] | 84.71 | 16.33 | 11.17 | 1.17 | - | - |
| Ntirogiannis [2] | **94.09** | **20.40** | **2.68** | 0.70 | - | - |
| Winner H-DIBCO10 | 91.50 | 19.78 | 5.98 | 0.49 | - | 93.58 |
| Su [17] | 92.03 | 20.12 | 6.14 | **0.25** | - | **94.85** |
| Ntirogiannis [1] | 70.64 | 15.22 | 22.16 | 1.44 | - | 84.22 |
| Ntirogiannis [2] | **94.49** | **21.72** | **3.18** | 0.30 | - | 94.32 |
| Winner DIBCO11 | 80.86 | 16.14 | - | 64.42 | 104.48 | - |
| Su [17] | 87.80 | 17.56 | - | 5.17 | 4.84 | - |
| Howe [18] | 91.70 | 19.30 | - | **3.87** | 3.48 | - |
| Ntirogiannis [1] | 80.39 | 15.47 | - | 5.78 | 6.68 | - |
| Ntirogiannis [2] | **92.64** | **19.93** | - | 5.12 | **3.13** | - |
| Winner H-DIBCO12 | 89.47 | 21.8 | - | - | 3.44 | 90.18 |
| Ntirogiannis [1] | 76.96 | 16.21 | - | - | 7.77 | 87.28 |
| Ntirogiannis [2] | **95.12** | **22.29** | - | - | **1.89** | **94.84** |

8. Pratikakis, I., Gatos, B., Ntirogiannis, K.: H-DIBCO 2010 - Handwritten Document Image Binarization Competition. In: Proc. Int. Conf. on Frontiers in Handwriting Recognition. (2010) 727–732

9. Pratikakis, I., Gatos, B., Ntirogiannis, K.: ICDAR 2011 Document Image Binarization Contest (DIBCO 2011). In: Proc. Int. Conf. on Document Analysis and Recognition. (2011) 1506–1510

10. Pratikakis, I., Gatos, B., Ntirogiannis, K.: H-DIBCO 2012 - Handwritten Document Image Binarization Competition. In: Proc. Int. Conf. on Frontiers in Handwriting Recognition. (2012) 813–818

11. Otsu, N.: A thresholding selection method from hray-level histogram. IEEE Transactions on Systems, Man and Cybernetics **9**(1) (1979) 62–66

12. Niblack, W. In: An Introduction to Digital Image Processing. Englewood Cliffs, NJ: Prentice-Hall (1986) 115–116

13. Sauvola, J., Pietikainen, M.: Adaptive document image binarization. Pattern Recognition **33**(2) (2000) 225–236

14. Kim, I.K., Jung, D.W., Park, R.H.: Document image binarization based on topographic analysis using a water flow model. Pattern Recognition **35**(1) (2002) 265–277

15. Gatos, B., Pratikakis, I., Perantonis, S.J.: Adaptive degraded document image binarization. Pattern Recognition **39**(3) (2006) 317–327

16. Lu, S., Su, B., Tan, C.L.: Document image binarization using background estimation and stroke edges. International Journal on Document Analysis and Recognition **13**(4) (2010) 303–314

17. Su, B., Lu, S., Tan, C.L.: A robust document image binarization technique for degraded document images. IEEE Transactions in Image Processing **22**(4) (2013) 1408–1417

18. Howe, N.R.: Document binarization with automatic parameter tuning. International Journal on Document Analysis and Recognition (2012) DOI: 10.1007/s10032-012-0192-x.

19. Gatos, B., Pratikakis, I., Perantonis, S.J.: Improved document image binarization by using a combination of multiple binarization techniques and adapted edge information. In: Proc. Int. Conf. on Pattern Recognition. (2008) 1–4

20. Su, B., Lu, S., Tan, C.L.: Combination of document image binarization techniques. In: Proc. Int. Conf. on Document Analysis and Recognition. (2011) 22–26

21. Kwak, S., Chung, K., Choi, Y.: Video caption image enhancement for an efficient character recognition. In: Proc. Int. Conf. on Pattern Recgnition. (2000) 606–609

22. Wolf, C., Jolion, J.M., Chassaing, F.: Text localization, enhancement and binarization in multimedia documents. In: Proc. Int. Conf. on Pattern Recognition. (2002) 1037–1040

23. Merler, M., Kender, J.R.: Semantic keyword extraction via adaptive text binarization of unstructured unsourced video. In: Proc. Inter. Conf. on Image Processing. (2009) 261–264

24. Kamel, M., Zhao, A.: Extraction of binary character-graphics images from grayscale document images. CVGIP: Computer Vision Graphics and Image Processing **55**(3) (1993) 203–217

25. Li, J., Tian, Y., Huang, T., Gao, W.: Multi-polarity text segmentation using graph theory. In: Proc. Int. Conf. on Image Processing. (2008) 3008–3011

26. Saidane, Z., Garcia, C.: Robust binarization for video text recognition. In: Proc. Int. Conf. on Document Analysis and Recognition. (2007) 874–879

27. Zhou, Z., Li, L., Tan, C.L.: Edge based binarization for video text images. In: Proc. Int. Conf. on Pattern Recognition. (2010) 133–136

28. Canny, J.: A computational approach to edge detection. IEEE Transactions on Pattern Analysis and Machine Intelligence **8**(6) (1986) 679–698

29. Yang, Y., Yan, H.: An adaptive logical method for binarization of degraded document images. Pattern Recognition **33**(5) (2000) 787–807

30. Lee, H.J., Chen, B.: Recognition of handwritten chinese characters via short line segments. Pattern Recognition **25**(5) (1992) 543–552

31. Trier, D., Taxt, T.: Evaluation of binarization methods for document images. IEEE Trans. Pattern Anal. Mach. Intell. **17**(3) (1995) 312–315

32. Kavallieratou, E., Stathis, S.: Adaptive binarization of historical document images. In: Int. Conf. on Pattern Recognition. Volume 3. (2006) 742–745

33. Stathis, P., Kavallieratou, E., Papamarkos, N.: An evaluation survey of binarization algorithms on historical document images. In: Proc. Int. Conf. on Pattern Recognition. (2008) 1–4

34. Paredes, R., Kavallieratou, E., Lins, R.D.: ICFHR 2010 Contest: Quantitative evaluation of binarization algorithms. In: Proc. Int. Conf. on Frontiers in Handwriting Recognition. (2010) 733–736

35. Barney Smith, E.H.: An analysis of binarization ground truthing. In: Proc. Int. Workshop on Document Analysis Systems. (2010) 27–33

36. Badekas, E., Papamarkos, N.: Automatic evaluation of document binarization results. In: Proc. Iberoamerican Congress on Pattern Recognition. (2005) 1005–1014

37. Lu, H., Kot, A.C., Shi, Y.Q.: Distance-reciprocal distortion measure for binary document images. IEEE Signal Process. Lett. **11**(2) (2004) 228–231

# Landmark Detection for Unconstrained Face Recognition

Panagiotis B. Perakis [*]

National and Kapodistrian University of Athens
Department of Informatics and Telecommunications
p.perakis@di.uoa.gr; takis@antinoos.gr

**Abstract.** In this dissertation a novel method for 3D landmark detection and pose estimation, suitable for both frontal and side 3D facial scans, is presented. It exploits 3D and 2D information by using local shape descriptors to extract candidate interest points that are subsequently identified and labeled as anatomical landmarks. Additionally, a novel generalized framework for combining facial feature descriptors that can be used for landmark detection is introduced, and several feature fusion schemes are proposed and evaluated. However, feature detection methods which use general purpose shape descriptors cannot identify and label the detected candidate landmarks. To this end, a 3D *Facial Landmark Model* (FLM) of facial anatomical landmarks is introduced. Candidate landmarks, irrespectively of the way they are generated, can be identified and labeled by matching them with the FLM. Finally, a novel method for unconstrained face recognition is introduced. It employs the 3D landmark detector to provide an initial pose estimation and to indicate occluded areas with missing data for each facial scan. Subsequently, a 3D *Annotated Face Model* (AFM) is registered and fitted to the scan using facial symmetry to complete the occluded areas. Using a biometric signature resulted from the wavelet representation of the fitted AFM, the proposed method can perform comparisons among interpose facial scans, unlike previously proposed methods that require frontal scans.

**Keywords:** Biometrics, Face Recognition, Landmark Detection, Shape Models, Shape Descriptors, Feature Extraction, Feature Fusion, Pose Estimation, Partial Matching, Deformable Models.

## 1 Introduction

*Biometrics* is the science of establishing the identity of a person based on the physical (e.g., fingerprints, face, hand geometry, and iris) or behavioral (e.g., gait, signature, and keyboard dynamics) attributes associated with an individual [1].

*Face recognition* is the procedure of recognizing an individual from their facial attributes or features and is one of the primary biometric modalities. Face

---

[*] Dissertation Advisor: Theoharis Theoharis, Professor.

recognition has several advantages over other biometric technologies: it is non-intrusive, since the facial region is generally exposed, and potentially easy to use [2]. Thus, research and development in face recognition followed naturally.

The performance of face recognition systems has improved significantly since the first automatic face recognition system was developed by Kanade [3] in 1973. Furthermore face recognition can now be performed in "realtime" for images captured under constrained situations. Although progress in face recognition has been encouraging, it has also turned out to be a difficult endeavor, especially for unconstrained tasks where view point, illumination, inter-object occlusions, facial expressions and facial accessories vary considerably [2].

## 2 Challenges & Motivation

Face recognition has proved to be a very challenging task due to the numerous sources of variation in 2D and 3D facial data. These variations can be environment-based (illumination conditions, occlusions by other objects or accessories), subject-based (pose and expression variations) and acquisition-based (image scale, distortion, noise, spikes and holes).

The main reason for using information from 3D data in face recognition systems is that the data acquired by 3D acquisition devices are invariant to pose and lighting conditions, these being the major challenges with which face recognition algorithms must cope [4].

With the increase in the availability of 3D data, several 3D face recognition approaches have been proposed. These approaches aim to overcome the limitations of 2D face recognition by offering pose invariance. However, although they claim pose invariance, they mostly utilize frontal 3D scans assuming that the entire face is visible to the sensor (see the surveys of Bowyer *et al.* [5] and Chang *et al.* [6]). This assumption is not always valid in real-world applications, since unconstrained acquisition may lead to facial scans with extensive occlusions that result in missing data due to pose variations.

Thus, existing 3D face recognition methods, fail to address large pose variations and to confront the problem of missing facial areas in an automatic way. The main assumption of these methods is that even though the head can be rotated with respect to the sensor, the *entire* face is always visible. However, this is true only for "almost frontal" scans or "reconstructed" complete face meshes aligned to frontal pose. *Side scans usually have large missing areas, due to self-occlusion, that depend on pose variations.* These scans are very common in realistic scenarios such as uncooperative subjects or uncontrolled environments. Therefore, to take advantage of the full pose invariance potential of 3D face recognition, the problem of missing data must be addressed. Thus, in a face recognition system, an initial registration step, based on landmark points' correspondence, is necessary in order to make the system pose invariant [7, 8].

However, facial landmark detection also suffers from the same sources of variation in 2D and 3D facial data that face recognition does [9–13]. Both 2D and 3D facial landmark detection suffer from occlusion, pose and expression variations.

In addition, 2D facial landmark detection also suffers from illumination variations. Thus, a landmark detection algorithm must be pose-invariant to address the problem of missing facial areas and, at the same time, expression-invariant in order to allow the registration of the various instances of the face liable to expression variations.

## 3  Aim & Methodology

The uncontrolled conditions of real-world biometric applications pose a great challenge to any 3D face recognition approach. The unconstrained acquisition of data from uncooperative subjects may result in facial scans with significant pose and expression variations.

*In this dissertation, an integrated novel method is proposed, in order to automatically detect landmarks on 3D facial scans that exhibit pose and expression variations, and hence consistently register and compare any pair of facial datasets subjected to missing data due to self-occlusion in a pose- and expression-invariant face recognition system.*

The proposed landmark detection and face recognition system employs an automatic pose- and expression-invariant landmark detector, using local facial feature descriptors and a deformable 3D *Facial Landmark Model* (FLM) to ensure global topological consistency of the detected landmarks [14, 8, 15, 16].

### 3.1  Training of Facial Landmark Models and Feature Templates

At the training phase, a Facial Landmark Model (FLM) is created by first aligning the training landmark sets and calculating a mean landmark shape using *Procrustes Analysis*, and then applying *Principal Component Analysis* (PCA) to capture the shape variations [17–19]. The FLM serves as a 3D geometric model of the landmark points. Also, templates for each shape descriptor that represents each landmark point are calculated from training facial datasets [14, 8, 15, 16].

The shape templates serve as feature descriptors for each landmark point. The feature descriptors that have been used, depending on the case, include the *Shape Index* [20], a continuous map of principal curvature values of a 3D object's surface, the *Spin Image* [21], a local descriptor of the object's 3D point distribution, the *Extruded Points* [15], a local descriptor of a 3D object's points that extrude most and the *Edge Response* [22] descriptor, a local descriptor of the 2D texture gradient of a 3D object.

### 3.2  Facial Landmark Detection

At the detection phase, the algorithm first detects candidate landmarks on the queried facial datasets according to the similarity of the extracted facial features
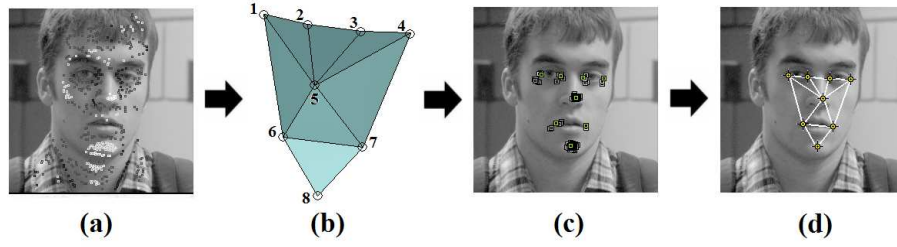
**Fig. 1.** Process pipeline of landmark detection: (a) extracted candidate landmarks using feature descriptors; (b) Facial Landmark Model (FLM); (c) landmark sets consistent with FLM; (d) resulting optimal landmark set.

with the feature templates. The extracted candidate landmarks are then filtered out and labeled by matching them with the FLM (Fig. 1) [14, 8, 15, 16].

During the research conducted under this dissertation, several versions of the presented generalized framework for facial landmark detection were applied. The most important are summarized in the following:

**SISI–NPSS METHOD** To locate landmark points, shape index target values for each landmark class (eye outer corner, eye inner corner, nose tip, mouth corner and chin tip) were searched for on the shape index map. Subsequently, the candidate landmark points of the five landmark classes that are obtained from the shape index map were further filtered out according to the similarity of their spin images with the spin image templates representing each landmark class. The resulting candidate landmark points of the five landmark classes were subsequently filtered out according to their consistency with the FLM. To find the optimum landmark set, the product of *normalized Procrustes distance ×* *(1–mean spin similarity)* was used as a distance metric between the candidate landmark sets and the FLM.

**Fusion METHOD** In this method, fusion schemes for combining landmark features were incorporated into the landmark detection pipeline. To locate landmark points the shape index map, the spin image map and the edge response map were fused into a resultant similarity map, each for every landmark class. The candidate landmarks for each landmark class were searched on the corresponding resultant similarity map. Subsequently, the candidate landmark points of the five landmark classes were filtered out according to their consistency with the FLM. To find the optimum landmark set, the product of *normalized Procrustes distance × (1 – resultant similarity)* was used as a distance metric between the candidate landmark sets and the FLM.

**Fig. 2.** Interpose matching using the proposed method: (a) and (b) opposite side facial scans with extensive missing data and detected landmarks; (c) generic Annotated Face Model (AFM); (d) and (e) registered and deformed AFM for each scan (facial symmetry used); (f) and (g) extracted geometry images.

### 3.3 Partial Face Recognition

The landmark detector provides an initial pose estimation (frontal, right, left) and indicates occluded areas with missing data for each facial scan resulting from pose variations. Facial landmark detection is a crucial first step for the registration of the facial datasets that have to be compared [8, 15].

Subsequently, a generic *Annotated Face Model* (AFM) [7] is registered and fitted to each facial probe scan, using a subdivision-based deformable model framework. During fitting, facial symmetry is used to complete the occluded areas of the face [8, 15]. Signature metadata are extracted using a *wavelet transformation* on the *geometry* and *normal images* of the fitted AFM (Fig. 2). A similarity measure between signature metadata of probe and gallery facial datasets provide the face recognition results.

## 4 Experimental Results

### 4.1 Landmark Detection

**Test Databases** For the performance evaluation of the proposed landmark detector, the largest publicly available 3D face and ear databases were combined. To evaluate the performance of the method against yaw variations, frontal, semi-profile and profile facial datasets were used. To evaluate the tolerance of the

method against expression variations, subjects with varying degrees of expressions were included. To have a measure of the landmark detection error, the used facial datasets were manually annotated at the queried landmark points. For frontal facial scans, the FRGC v2 database [23, 24] was used. For side facial scans, the Ear Database from the University of Notre Dame (UND) [25] was used.

For the conducted experiments, the following collections of facial datasets were created:

- DB00F: Contains 975 frontal facial scans obtained from 149 different subjects, selected from the FRGC v2 database, including subjects with varying degrees of expressions (45.44% "neutral", 36.41% "mild" and 18.15% "extreme"), acquired under varying illumination conditions (e.g. half of the face shaded).
- DB00F45RL: a composite frontal-to-profile database with the datasets of 39 common subjects found in the FRGC v2 database and in the UND Ear database. This database contains 117 (3x39) facial scans having three poses, frontal (39 scans) and 45° left (39 scans) and right (39 scans).
- DB45L and DB45R: two semi-profile databases with 118 left and 118 right 45° side datasets, which come from 118 different subjects, obtained from the UND Ear database.
- DB60L and DB60R: two profile databases with 87 left and 87 right 60° side datasets, which come from 87 different subjects, obtained from the UND Ear database.

**Landmark Detection Evaluation** The performance evaluation of a landmark detector is generally presented by computing the following values, which represent the localization accuracy of the detected landmarks:

**Absolute Distance Error**: The Euclidean distance in physical units (e.g., $mm$) between the position of the detected landmark and the manually annotated landmark, which is considered ground truth.

**Detection Success Rate**: The percentage of successful detections of a landmark over a test database. Successful detection is considered as the detection of a landmark with Absolute Distance Error under a certain threshold (e.g., $10\ mm$).

Summary results for METHOD SISI–NPSS on all tested databases are presented in Table 1. The results clearly indicate that the proposed method exhibits high accuracy and robustness both to yaw and expression variations. The mean error is under 6.3 $mm$, with standard deviation under 2.6 $mm$ on all tested facial scans. Also note that the mean error is under 10 $mm$ for at least 90.4% of the tested facial scans and the facial side was correctly estimated on over 98.9% of the tested facial scans.

**Table 1.** Summary results for METHOD SISI–NPSS

| Database | Mean Error | | | Side |
| | mean (mm) | stdev (mm) | ≤ 10 (mm) | Detection Rate |
| --- | --- | --- | --- | --- |
| DB00F | 5.00 | 1.85 | 97.85% | 99.90% |
| DB00F-neutral | 4.52 | 1.51 | 99.32% | 100.00% |
| DB00F-mild | 4.95 | 1.46 | 99.72% | 100.00% |
| DB00F-extreme | 6.28 | 2.60 | 90.40% | 99.44% |
| DB00F45RL | 4.97 | 1.92 | 97.44% | 100.00% |
| DB45R | 5.03 | 1.92 | 96.61% | 100.00% |
| DB45L | 4.75 | 1.91 | 97.46% | 100.00% |
| DB60R | 4.95 | 1.80 | 96.55% | 98.85% |
| DB60L | 5.30 | 2.49 | 93.10% | 100.00% |

**Evaluation of Fusion Schemes** The evaluation of the performance of the proposed distance to similarity mappings and fusion schemes for landmark detection is not a straight-forward task, since there are many factors that characterize performance. As already stated, fusion techniques are expected to improve system's *accuracy*, *efficiency* and *robustness*. An equally important characteristic of a fusion scheme is that of *monotonicity*, i.e., the addition of a new feature descriptor should improve prior results. A qualitative performance evaluation of the proposed fusion schemes according to the aforementioned characteristics is presented in Table 2.

**Table 2.** Qualitative evaluation of proposed fusion schemes

| | Accuracy | Efficiency | Robustness | Monotonicity |
| --- | --- | --- | --- | --- |
| L−L1 | Fair | **High** | Fair | Fair |
| L−L2 | Fair | Low | Fair | Fair |
| L−Lg | **High** | Fair | Fair | Fair |
| Q−L1 | **High** | **High** | Fair | Fair |
| Q−L2 | **High** | **High** | **High** | **High** |
| Q−Lg | **High** | Fair | Fair | Fair |
| G−L1 | **High** | **High** | **High** | **High** |
| G−L2 | **High** | **High** | Fair | Fair |
| G−Lg | **High** | Fair | Fair | Fair |
| L−Lmax | Low | Low | Low | Low |
| Q−Lmax | Low | Low | Low | Low |
| G−Lmax | Low | Low | Low | Low |
| L−Lmin | Unreliable | Fair | Fair | Low |
| Q−Lmin | Unreliable | Fair | Fair | Low |
| G−Lmin | Unreliable | Fair | Fair | Low |

Current experimental results show that, in general, the Quadratic (Q) and Gaussian (G) mappings behave better than the Linear (L) mapping of distance measure to similarity measure. For the Linear mapping the product rule (Lg) behaves better than other rules. For the Quadratic mapping the rms rule (L2) behaves better than other rules. For the Gaussian mapping the sum rule (L1) behaves better than other rules. Quadratic and Gaussian mappings have almost the same performance.

Accuracy improvement is more dramatic when the information fused is correlated. In correlated features the performance of one descriptor predicts to some extent the performance of the other and strengthens the results. On the other hand highly uncorrelated features have similarity peaks that do not coincide and degrade the results. Efficiency improvement is achieved by excluding obvious non-matches, reducing the number of candidate landmarks, for each landmark class. Fusion, also, reduces system sensitivity to sample-specific, poor-quality or erroneous descriptors.

We can thus deduce that the best performance in terms of accuracy is exhibited by the Q-L2 and G-L1 fusion schemes, with the Q-L2 exhibiting a slight better performance than the G-L1 in landmarks' likelihood area reduction. Q-L2 and G-L1 also exhibit high robustness in yaw, expression and illumination variations, and strong monotonicity.

Also landmark localization using the Q-L2 fusion scheme improved the accuracy and robustness of the landmark detector (with $3.5-5.5 \ mm$ mean landmark localization error), indicating the superiority of the fusion approach.

### 4.2   Partial Face Recognition

**Test Databases**

**Combined UND Databases:** To evaluate the performance of the proposed partial face recognition method, a combination of the largest publicly available 3D face and ear databases was used. For frontal facial scans, the FRGC v2 database [23, 24] was used. For side facial scans, the Ear Database from the University of Notre Dame (UND) [25] was used.

For the conducted experiments the following collections were defined:

 – UND45LR: Contains 45° side scans from 118 subjects. For each subject, the left scan is considered gallery and the right is considered probe. *Total: 236 scans.*
 – UND60LR: Contains 60° side scans from 87 subjects. For each subject, the left scan is considered gallery and the right is considered probe. *Total: 174 scans.*
 – UND00LR: Gallery set has one frontal scan for each of the 466 subjects. Probe set has two 45° side scans (left and right) from 39 subjects and two 60° side scans (left and right) from 32 subjects. *Total: 608 scans.*

In all cases there is only one gallery scan per subject. Also, all subjects present in a probe set are also present in the gallery set (the opposite is not always true).

**UH Databases:** In addition to the UND databases a database with data collected at the University of Houston was used. The database contains 1,075 left and 1,075 right scans of 281 subjects. The novelty of this database is that each pair of left and right side scans was acquired simultaneously.

For the conducted experiments the following collection is defined:

– UHDB7LR-M: Contains multiple left and right side scan pairs from 281 subjects. For each subject, one left and one right scan are considered gallery and the rest are considered probes (1–6 left and 1–6 right scans per subject). *Total: 2,150 scans.*

In all cases there is one pair of gallery scans per subject. Also, all subjects present in a probe set are also present in the gallery set (the opposite is not always true).

The proposed method tackles the problem of matching arbitrary facial scans (left, right or frontal). This is considerably harder than matching only frontal scans, since a lot of the facial information is missing and it is not known a priori whether each scan is left, right or frontal.

**Matching facial scans of arbitrary side** In this experiment, the performance of the proposed partial face recognition method, using scans of arbitrary sides for gallery and probe sets, was evaluated. This is a realistic scenario, as the side scans (with extensive occlusions that lead to missing data) are very common in real world applications with unconstrained acquisition. The proposed method can match any combination of left, right or frontal facial scans with the use of facial symmetry. Moreover, the proposed method automatically detects the side of the scan by using the automatic landmark detector. For this experiment we utilized the UND45LR, UND60LR, UND00LR and UHDB7LR-M databases and the rank-one rates are given in Table 3.

**Table 3.** Rank-one Recognition Rate between facial scans of arbitrary side

|           | *Rank-one Rate* |
|-----------|-----------------|
| UND45LR   | 86.4%           |
| UND60LR   | 81.6%           |
| UND00LR   | 76.8%           |
| UHDB7LR-M | 89.1%           |

In the cases of UND45LR and UND60LR, for each subject, the gallery set contains a single left side scan while the probe set contains a single right side scan. Therefore, facial symmetry is always used in order to perform identification. As expected, the 60° side scans yield lower results as they are considered more challenging compared to the 45° side scans.

In the case of UND00LR, the gallery set contains a frontal scan for each subject, while the probe set contains left and right side scans. This scenario is

very common when the enrollment of subjects is controlled but the identification is uncontrolled. Compared to UND45LR and UND60LR, there is a decrease in the performance of the proposed method in UND00LR. One could argue that since the gallery set consists of frontal scans (without missing data), there should be an increase in performance. However, UND00LR has the largest gallery set, making it the most challenging database in current experiments.

In the case of UHDB7LR-M, for each subject, the gallery set contains a left and right side scan pair, while the probe set contains multiple left and right side scan pairs. As expected, since the gallery set has two scans per subject, the performance on this database is the highest among all databases. The performance difference is substantial compared to UND00LR (89.1% versus 76.8% rank-one). This indicates that one pair of left and right side scans is more descriptive than one frontal scan.

## 5 Conclusion

In this thesis an automatic facial landmark detection methodology has been proposed. It offers pose invariance and robustness to large missing (self-occluded) facial areas with respect to large yaw variations and high tolerance to large expression variations. The proposed approach consists of methods for landmark localization that exploit the 3D facial geometry and the modeling ability of trained landmark models. It has been evaluated using the most challenging 3D facial databases available, which contain scans with yaw variations of up to $80°$ and strong expressions. In these databases it achieved state-of-the-art accuracy (with $4.5 - 6.3\ mm$ mean landmark localization error), significantly outperforming existing methods.

Also, a novel generalized framework of fusion methods and their application to landmark detection has been presented. The proposed fusion scheme transforms features to similarities and then combines them to generate a resultant feature similarity, which is considered as the matching score for the detection of the queried landmarks. The proposed feature fusion framework is easily extensible to new feature-components, offers significant dimensionality reduction and works equally well for features extracted from 3D or 2D facial data.

For the proposed fusion scheme different distance to similarity mappings and different fusion rules have been evaluated. The results indicate that the quadratic distance to similarity mapping in conjunction with the rms rule for fusion (Q-L2) exhibits the best performance. Landmark localization using this fusion scheme achieved state-of-the-art accuracy (with $3.5 - 5.5\ mm$ mean landmark localization error), indicating the superiority of the fusion approach.

Finally, a novel 3D face recognition method suitable for real-world biometric applications was proposed. Unlike most previous methods that require frontal scans, the proposed method can perform partial matching among interpose facial scans, even when extensive data are missing. It exploits the 3D landmark detector to provide an initial pose estimation and to indicate occluded areas with missing

data for each facial scan. By using facial symmetry to complete missing facial data, it can handle seamlessly frontal and side facial scans.

The presented method for partial face recognition is extensively evaluated against a variety of 3D facial databases, achieving state-of-the-art performance (with average rank-one recognition rate 83.7%), considerably outperforming existing methods, even when tested on the most challenging data, which contain scans with yaw variations up to 80° and strong expressions.

The proposed system is suitable for real-world scenarios as the only requirement is that half of the face is visible to the sensor, and its computational cost is low. Using a standard Intel Core 2 Duo 2.2 $GHz$ PC, 18 $sec$ on average are required to process a facial scan: 9 $sec$ to localize the facial landmarks plus 9 $sec$ to extract the biometric signature (geometry and normal images). The biometric signatures can be matched at a rate of 15,000 $matches/sec$.

## Ackowledgement

## References

1. Ross, A., Jain, A.K.: Biometrics, overview. In Li, S., ed.: Encyclopedia of Biometrics. Springer, New York, NY (2009) 168–172
2. Li, S., Jain, A.K. In: Handbook of Face Recognition. Springer (2005) 1–11
3. Kanade, T.: Picture Processing by Computer Complex and Recognition of Human Faces. PhD thesis, Kyoto University (1973)
4. Kakadiaris, I., Passalis, G., Toderici, G., Perakis, P., Theoharis, T.: 3D-based face recognition. In Li, S., ed.: Encyclopedia of Biometrics. Springer, New York, NY (2009) 329–338
5. Bowyer, K., Chang, K., Flynn, P.: A survey of approaches and challenges in 3D and multi-modal 3D+2D face recognition. Computer Vision and Image Understanding **101**(1) (Jan. 2006) 1–15
6. Chang, K., Bowyer, K., Flynn, P.J.: An evaluation of multi-modal 2D+3D face biometrics. IEEE Transactions on Pattern Analysis and Machine Intelligence **27**(4) (Apr. 2005) 619–624
7. Kakadiaris, I., Passalis, G., Toderici, G., Murtuza, M., Lu, Y., Karampatziakis, N., Theoharis, T.: Three-dimensional face recognition in the presence of facial expressions: An annotated deformable model approach. IEEE Transactions on Pattern Analysis and Machine Intelligence **29**(4) (Apr. 2007) 640–649
8. Passalis, G., Perakis, P., Theoharis, T., Kakadiaris, I.: Using facial symmetry to handle pose variations in real-world 3D face recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence **33**(10) (Oct. 2011) 1938–1951

9. Lu, X., Jain, A.: Multimodal facial feature extraction for automatic 3D face recognition. Technical Report MSU-CSE-05-22, Michigan State University (Oct. 2005)

10. Colbry, D., Stockman, G., Jain, A.: Detection of anchor points for 3D face verification. In: Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, San Diego, CA (Jun. 20-25 2005) 118

11. Lu, X., Jain, A.: Automatic feature extraction for multiview 3D face recognition. In: Proc. $7^{th}$ International Conference on Automatic Face and Gesture Recognition, Southampton, UK (Apr. 10-12 2006) 585–590

12. Lu, X., Jain, A., Colbry, D.: Matching 2.5D face scans to 3D models. IEEE Transactions on Pattern Analysis and Machine Intelligence **28**(1) (2006) 31–43

13. Colbry, D.: Human Face Verification by Robust 3D Surface Alignment. PhD thesis, Michigan State University (2006)

14. Perakis, P., Passalis, G., Theoharis, T., Kakadiaris, I.: 3D facial landmark detection under large yaw and expression variations. IEEE Transactions on Pattern Analysis and Machine Intelligence **35**(7) (July 2013) 1552–1564

15. Perakis, P., Passalis, G., Theoharis, T., Toderici, G., Kakadiaris, I.: Partial matching of interpose 3D facial data for face recognition. In: Proc. $3^{rd}$ IEEE International Conference on Biometrics: Theory, Applications and Systems, Arlington, VA (Sep. 28-30 2009) 439–446

16. Perakis, P., Theoharis, T., Passalis, G., Kakadiaris, I.: Automatic 3D facial region retrieval from multi-pose facial datasets. In: Proc. Eurographics Workshop on 3D Object Retrieval, Munich, Germany (Mar. 30 - Apr. 3 2009) 37–44

17. Dryden, I., Mardia, K.: Statistical Shape Analysis. Wiley (1998)

18. Stegman, M., Gomez, D.: A brief introduction to statistical shape analysis. Technical report, Technical University of Denmark (Mar. 2002)

19. Cootes, T., Taylor, C.: Statistical models of appearance for computer vision. Technical report, University of Manchester (Oct. 2001)

20. Dorai, C., Jain, A.K.: COSMOS - a representation scheme for 3D free-form objects. IEEE Transactions on Pattern Analysis and Machine Intelligence **19**(10) (Oct. 1997) 1115–1130

21. Johnson, A.E.: Spin Images: A Representation for 3-D Surface Matching. PhD thesis, Robotics Institute, Carnegie Mellon University, Pittsburgh, PA (Aug. 1997)

22. Harris, C., Stephens, M.: A combined corner and edge detector. In: Proc. $4^{th}$ Alvey Vision Conference. (1988) 147–151

23. Phillips, P., Flynn, P., Scruggs, T., Bowyer, K., Chang, J., Hoffman, K., Marques, J., Min, J., Worek, W.: Overview of the Face Recognition Grand Challenge. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition, San Diego, CA (2005) 947–954

24. Phillips, P., Scruggs, T., O'Toole, A., Flynn, P., Bowyer, K., Schott, C., Sharpe, M.: FRVT 2006 and ICE 2006 large-scale experimental results. IEEE Transactions on Pattern Analysis and Machine Intelligence **32** (2010) 831–846

25. UND: University of Notre Dame Biometrics Data Sets. http:// www.nd.edu/ ∼cvrl/ CVRL/ Data_Sets.html (2012)

# Modeling the regulatory intervention in the telecommunications market

Markos Tselekounis[*]

National and Kapodistrian University of Athens
Department of Infromatics and Telecommunications
markos@di.uoa.gr

**Abstract.** This thesis discusses the role of sector-specific regulators in the rapidly changing telecommunications industry. In particular, it studies the access pricing policy which provides the optimal balance between static and dynamic efficiency that better reflects the changing regulatory goals in a highly variable economic and technological environment. In fact, there are three distinct phases in the evolution of the telecommunications markets which directly affect the optimal mixture of regulatory policy. These phases are: (i) the migration from a state monopoly market to a competitive telecommunications industry which reflects the past regulatory goal of achieving static efficiency; (ii) the migration from service-based competition over copper access networks to service-based competition over NGA networks which reflects the current regulatory goal of inducing NGA investments without distorting competition; and (iii) the migration from service-based to facilities-based competition over NGA networks which reflects the future regulatory goal of promoting dynamic efficiency. It is obvious that a different regulatory policy is required to be implemented in each migration phase in order to fulfill the desirable investment and competition outcomes. This thesis models the regulatory intervention in the telecommunications market and derives the access pricing policy that achieves the efficiency goals of each migration phase.

## 1    Introduction

The liberalization of the telecommunications markets in the United States (US) and Britain in the early 1980s and in Europe in the late 1990s was the result of the conventional wisdom that competition serves consumers and social welfare better than the former state monopoly, both from a short-term perspective, where entry and investment decisions are taken as given as well as from a long-term perspective, where these are treated as endogenous. However, the migration from a state monopoly market to a competitive telecommunications industry required the existence of a sector-specific regulator for the restructuring process of the telecommunications sector be-

---

[*] Dissertation Advisor: Dimitris Varoutas, Assistant Professor.

cause the incumbent upstream monopolist was also a supplier of the final services, and hence, there was the obvious danger that this integrated firm would seek to exclude competing providers by setting high access prices.

This fact led to a fierce debate about the terms and conditions on which competitors would have unbundled upstream access to the historical operators' local loop facilities. The reason is that regulators should achieve too many goals with only one instrument: the determination of the access charge (i.e. the price that new entrants should pay to the incumbent in order to have access to its local loop facilities). The main goals of the regulatory policy are: (i) the achievement of (static) economic efficiency, with a particular focus on improving consumers' surplus, which is achieved through low prices and high quality; and (ii) the achievement of dynamic efficiency so that investment incentives give rise to socially optimal investment decisions.

Static efficiency concerns the short-run regulatory goal to reduce the incumbent's market power in order to enable alternative operators (new entrants) to enter the market and compete effectively with the incumbent in the downstream (retail) market. Unbundling of the local loop facilitates entry by allowing new entrants to have the right to use the same network as the incumbent. As a result, both incumbent and entrants have significant incentives to invest in innovative, differentiated services. Such service-based competition promotes productive efficiency (i.e. existing assets are utilized efficiently) and allocative efficiency (i.e. existing resources are efficiently allocated to the economy). Therefore, service-based competition ensures that firms behave in a competitive manner, and hence, consumers enjoy the welfare gains from static efficiency (lower prices, better quality and extended variety of services).

On the other hand, dynamic efficiency concerns the long-run goal of access regulation to induce firms to undertake the socially optimal (efficient) investment decisions in terms of both timing of investments and the extent of network deployment. According to Bourreau and Doğan [1], facilities-based competition, which requires investments in new competing infrastructures from the incumbents and (especially) entrants, leads to efficient investment decisions and adoption of better technologies. In particular, facilities-based competition is regarded as the only means to achieve sustainable competition since it creates a level playing field between the incumbent and entrants [2-4]. Facilities-based competition achieves the full benefits of competition, and hence, consumers enjoy the full welfare gains from dynamic efficiency (maximum market growth, minimized costs, innovative technologies and services).

It is obvious that since the promotion of efficient entry is a short-run goal in the transition from state monopoly to private and competitive market structures, access regulation should indisputably aim at fostering service-based competition. In practice, both in the United States and in the European Union a light regulation with unregulated retail prices combined with *ex ante* regulation of the upstream access component has become dominant. The Telecommunications Act of 1996 [5] administered by the Federal Communications Commission (FCC) as well as the European Commission's Regulation on Local Loop Unbundling [6] mandated unbundled access to the metallic local loops of incumbent operators at cost-based prices. According to Armstrong [7], the chief benefits of cost-based access charges are two-fold. Firstly, there is no need for information about the demand for the final services. In particular, the only information needed is the cost of providing the access which is needed for all access pric-

ing policies. Secondly, cost-based access regulation is the only access pricing policy that gives the correct make-or-buy signals to entrants when bypass is a possibility.

Indeed, cost-based access regulation has led to improved service-based competition in many European countries, and hence, it seems that consumers enjoy the welfare gains from static efficiency. However, this expectation lacks of theoretical justification since academic research has focused on studying the impact of access prices on an entrant's incentives to undertake the productively efficient make-or-buy decision. This thesis contributes to the related literature by studying the conditions under which (cost-based) access prices induce the entrant to undertake the efficient make-or-buy decision in terms of both productive and allocative efficiency. This thesis shows that when the only goal of regulators is to achieve static efficiency, they should simply set the input prices at the incumbent's marginal cost of producing the upstream input since cost-based access prices lead the entrant to undertake the productively efficient make-or-buy decision which is also socially optimal [8, 9].

However, in the last decade the number of internet users as well as the capacity they demand has increased dramatically making the traditional access copper networks incapable of providing end-users with the demanded bandwidth. On the contrary, the transmission capabilities of fibre are theoretically unlimited. For this reason, the deployment of fibre access infrastructures, the so-called Next Generation Access (NGA) networks, has received significant interest among all operators since they are regarded as the only future proof solution capable to handle future demand [10]. In addition, investment in NGA networks has also attracted the interest of national governments since higher speed broadband services increase the positive impact of broadband on economic growth, productivity at the firm level, employment growth and consumers' welfare [11–14]. However, investment in NGA networks not only requires a huge initial fixed cost, but also is mainly sunk once the investment has been made. This implies that potential investors are reluctant to invest in NGA networks unless they are reimbursed for the risk they incur when investing in such networks. In other words, cost-based access prices, which are limited to boost entry and promote service-based competition within one network, discourage both incumbents and entrants to invest in new facilities as well as result in a substantial deviation from the socially desirable investment outcomes implying losses in dynamic efficiency [15].

Perhaps the most challenging task for academics, governments and policy makers is to design a regulatory policy that encourages investments in NGA networks and promotes sustainable competition. This implies that regulators aim at facilitating the migration from service-based competition over copper access networks to service-based competition over NGA networks. Therefore, the related literature focuses on proposing alternative regulatory practices which aim at achieving the current regulatory two-fold goal.

A first proposal concerns the deviation from cost-based access prices by implementing investment-contingent access charges. This thesis points out that the related literature fails to take into account the fact that it is uncertain whether the regulator will set an investment-contingent or a welfare-maximizing access price after the NGA deployment. In particular, this thesis models this fact in order to study the impact of regulatory uncertainty on an incumbent's incentives to undertake the socially optimal investments in NGA networks [16]. It is found that when the slope of the marginal investment cost function is not particularly steep in relation to the impact of invest-

ments on demand, the incumbent underinvests compared to the socially optimal investment level. On the contrary, in the more realistic case when the impact of investments on demand is low in relation to the slope of the marginal investment cost function, the incumbent may overinvest or underinvest depending on the probability of incorporating an access markup into the access price.

A second proposal concerns the deviation from the permanent regulation of access by implementing alternative regulatory regimes such as "regulatory holidays" and "sunset clauses". Particular attention has received the implementation of the "regulatory holidays" access regime, under which the investor is not imposed to any regulatory constraints for a pre-determined period of time. The reason for such particular attention is the implementation of "access holidays" in the US broadband markets and the dispute between the German government and the European Commission (EC) about the power of national legislation (which envisioned the provision of "access holidays" to the German incumbent operator) to limit the discretionary powers of the national regulator in its exclusive right to assess whether markets should be regulated or not under EU rules [17].

Obviously, such a regulatory policy provides significant investment incentives but also ambiguous outcomes in terms of social welfare. This thesis contributes to the debate about the effectiveness of "regulatory holidays" to provide efficient outcomes by studying: (i) the impact of geographic price discrimination on an unregulated monopolist's incentives to deploy a larger NGA network and on the subsequent social welfare outcomes [18]; and (ii) the optimal decision of an unregulated operator to deploy different quality NGA technologies in geographic areas which differ in their population density [19]. It is found that: (i) the regulator should allow the monopolist to geographically price discriminate as long as the investment cost is not extremely low since in this case the monopolist chooses the socially optimal pricing regime; and (ii) although a geographically differentiated NGA investment provides the unregulated monopolist with incentives to install a nationwide NGA deployment, the monopolist underinvests compared to the socially optimal levels of both quality and geographic coverage.

Even though service-based competition over NGA networks increases both static and dynamic efficiency, the full benefits of competition are only achieved by facilities-based competition. This explains why the ultimate goal of regulators is to promote dynamic efficiency which results in maximum welfare gains, maximum market growth and minimum production costs.

This thesis reviews the proposed regulatory approaches which aim to encourage access seekers to invest in their own fibre-based access networks when an initial investor has already deployed an NGA network and sustainable service-based competition has been established. In addition, a comparison of these regulatory approaches with the current regulatory framework in the European NGA market as described by the EC Recommendation shows that the proposed regulatory approaches not only fail to reflect the basic principles of the EC Recommendation, but also fail to take into account the fact that the regulatory policy applied in this phase has a direct impact on the initial investor's incentives to invest in NGA networks [20].

For this reason, this thesis also presents an innovative theoretical approach that not only reflects the current regulatory framework in the European NGA market, but also encourages the initial investor (which is assumed to be the incumbent) to invest in

NGA networks, although at the same time it incentivizes the entrants to gradually invest in their own NGA infrastructures. It is shown that the proposed approach, which is based on the basic principles governing a Credit Default Swap (CDS), provides an effective migration path towards facilities-based competition over NGA networks [20].

It is thus obvious that this thesis studies the interplay between the continuously evolving scope of telecommunications regulation and technological development by modeling the regulatory intervention and deriving the access pricing policy that achieves the efficiency goals of each migration phase.

## 2    Main thesis contributions

### 2.1    On the social optimality of make-or-buy decisions

Many economists argue that cost-based access prices encourage the right amount of entry, and hence, lead to service-based competition in the downstream market. On the contrary, Sappington [21] shows that input (or access) prices are irrelevant for an entrant's decision to make or buy an input required for downstream production when the competition between the providers in the downstream market is described by the standard Hotelling model. According to Sappington, the reason for this striking result is that previous studies fail to take into account the impact of a new entrant's make-or-buy decision on subsequent retail price competition. When the incumbent sells an upstream input to the new entrant, the incumbent faces an opportunity cost of expanding its retail output. The incorporation of this opportunity cost into the incumbent's total cost makes the incumbent act as if its upstream cost of production were equal to the specified input price. Therefore, regardless of the input price, the entrant will choose to buy (respectively, make) the upstream input whenever the incumbent (respectively, entrant) has an innate upstream cost advantage. Hence, the entrant's decision always minimizes industry costs and ensures efficient entry and utilization of the telecommunications infrastructure. Thus, the entrant always undertakes the productively efficient make-or-buy decision.

In addition, Tselekounis, Varoutas and Martakos [8] complement the work of Sappington by studying the effectiveness of input prices on inducing the entrant to undertake the socially optimal make-or-buy decision. They show that input prices do not have an impact on social welfare. The reason is that a marginal increase (decrease) in the input price causes a unit increase (decrease) in the incumbent's profits and a unit decrease (increase) in consumer surplus. As social welfare is the unweighted sum of industry profits and consumer surplus, it is thus not affected by a marginal change in input prices. Therefore, input prices are irrelevant not only for the entrant's efficient make-or-buy decision, but also for the regulator's goal to maximize social welfare. In particular, they show that regardless of the established price of the upstream input, the entrant's decision to buy (respectively, make) the upstream input from the incumbent is socially optimal when the incumbent (respectively, entrant) is the least-cost supplier of the input. As a result, in the equilibrium of the Hotelling model, the entrant un-

dertakes the efficient make-or-buy decision in terms of both productive and allocative efficiency regardless of the regulated input price.

However, these results are found to be strongly dependent on the particular model of downstream competition. Gayle and Weisman [22] consider the impact of input prices on the entrant's incentives to undertake the productively efficient make-or-buy decision under alternative downstream interactions. They show that input prices are not necessarily irrelevant in the Bertrand vertical differentiation model and are not irrelevant in the Cournot model. In addition, cost-based input prices always result in the productively efficient outcome. This implies that departure from cost-based input prices may distort the efficiency of the entrant's make-or-buy decision.

Tselekounis, Varoutas and Martakos [9] study the robustness of the result concerning the irrelevance of input prices to the entrant's incentives to undertake the productively and allocatively efficient make-or-buy decision when the downstream competition is not characterized by the Hotelling model but downstream interactions are better described by the Cournot or the Bertrand vertical differentiation competition model. They find that the social optimality of the entrant's make-or-buy decision is affected by two crucial factors: (i) the particular level of the price of the upstream input; and (ii) the cost differential between the incumbent's and the entrant's unit costs of producing the upstream input. For this reason, they obtain the range of input prices and upstream cost differential that induce the entrant to undertake the socially desirable decision. They conclude that the entrant's productively efficient make-or-buy decision is socially optimal for the set of input prices that induce the entrant to undertake the efficient decision in the case of Cournot competition and is not necessarily socially optimal in the Bertrand vertical differentiation model.

It is thus obvious that the particular model that describes the competition in the downstream market as well as each provider's efficiency in producing the upstream input have a significant impact on the social optimality of the entrant's (efficient) make-or-buy decision. This implies that regulators should have perfect information about each provider's unit cost of producing the upstream input and the way that the two providers compete in the downstream market in order to draw their optimal access pricing policy. However, when the only goal of regulators is to achieve static efficiency (e.g. in the transition from state monopoly to private and competitive market structures), they should simply set the input prices at the incumbent's marginal cost of producing the upstream input since the results of [9] show that regardless of the type of competition, cost-based access prices lead the entrant to undertake the productively efficient make-or-buy decision which is also socially optimal.

## 2.2 Investments in Next Generation Access infrastructures under regulatory uncertainty

The related literature discusses the effectiveness of two different regulatory approaches on the regulator's current goal to achieve the socially efficient investment level when it sets the access price after the investment decision of the incumbent. The first approach supports that the regulator sets a particular investment-contingent access price, which compensates the incumbent for the investment risks, in order to

provide significant investment incentives. On the contrary, the second approach argues that the regulator deviates from such *ex ante* known access price (once the investments are in place) by setting the access price at the marginal cost of providing the access in order to maximize social welfare.

Tselekounis and Varoutas [16] modeled the more realistic case in which the regulator sets the access price at the marginal cost of providing the access with some probability and gives an access markup, which equals the average cost of the investments, with the complementary probability. Therefore, it is uncertain which of the two assumptions made in the related literature will prevail when the new access infrastructures are in place.

A non-commitment setting is used in order to take account for regulatory uncertainty. In addition, the retail (downstream) market is characterized as an unregulated duopoly market in which the incumbent and the entrant choose quantities simultaneously and independently (i.e. firms compete á la Cournot). The level of NGA investment undertaken by the incumbent leads to an outward parallel shift in the demand, and hence, NGA investments have a positive impact on the demand for the new fibre-based services. Furthermore, the incumbent faces a quadratic NGA investment cost function with respect to the investment level implying that the slope of the marginal investment cost function is linear and increasing in the investment level.

The privately and the socially optimal investment levels are derived as a function of the probability $\alpha \in [0,1]$ of incorporating into the access price an access markup, which equals the average cost of the investments, in order to fully compensate the incumbent for the NGA investment risk. A first significant finding is that a marginal increase in such probability positively affects the private investment incentives and negatively affects the socially optimal investments. The comparison of the privately and the socially optimal investment levels show that there is a unique positive value $\tilde{\alpha}$ of the probability of incorporating into the access price an access markup which induces the incumbent to undertake the socially optimal investments. If $\alpha > \tilde{\alpha}$ (respectively, $\alpha < \tilde{\alpha}$), the NGA investment level chosen by the incumbent is higher (respectively, lower) than the socially optimal one. This implies that any deviation from the socially optimal investment level leads to welfare losses.

A second significant result is that the derived value of $\tilde{\alpha}$ is significantly affected by the impact of the investments on demand and the slope of the marginal investment cost function. In particular, the value of $\tilde{\alpha}$ is positively affected by an increase in the impact of investments on demand and negatively affected by an increase in the slope of the marginal investment cost function (ceteris paribus). This implies that, for a given slope, higher consumers' valuation for the NGA services results in higher $\tilde{\alpha}$, which in turn leads to higher efficient investment levels. In other words, higher consumer consumers' valuation for the NGA services makes the investments more socially desirable, and hence, the socially optimal investment level is achieved for a higher probability of compensating the incumbent for the investment risks. This result positively affects the incumbent's investment incentives, and hence, the achieved efficient investment level increases as well.

On the contrary, for a given positive impact of the investments on demand, a steeper slope of the marginal investment cost function leads to lower values of $\tilde{\alpha}$. This implies that as the NGA investment becomes marginally more expensive, the society is better off by a lower NGA deployment which is achieved by a higher probability of

setting the access price at the marginal cost of providing the access. Therefore, the efficient NGA investment level is achieved for lower values of $\tilde{\alpha}$.

Combining the two aforementioned significant results leads to the main result of Tselekounis and Varoutas [16]:

(i) When the slope of the marginal investment cost function is not particularly steep in relation to the positive impact of investments on demand, the incumbent always underinvests compared to the socially optimal investment level. The reason is that the critical value of the probability of including an access markup into the access price ($\tilde{\alpha}$) is higher that 1. This implies that the socially desirable outcome cannot be achieved even if the regulator commits to an access price scheme that includes an access markup equal to the average cost of the investments. In this case, a higher access markup which leads to $\tilde{\alpha} \leq 1$ seems to be socially desirable.

(ii) On the contrary, in the more realistic case when the impact of investments on demand is low in relation to the slope of the marginal investment cost function, the incumbent may overinvest or underinvest depending on the probability of incorporating an access markup into the access price. In this case $\tilde{\alpha} \in (0,1)$, and hence, the incumbent overinvests for high probability of incorporating an access markup into the access price and underinvests for low probability values. As a result, the optimal social welfare outcome cannot be achieved with the incumbent's profit maximizing investment level when $\alpha \neq \tilde{\alpha}$. This implies that regulatory uncertainty significantly affects the incumbent's incentives to undertake the socially optimal investments in NGA networks.

## 2.3 A CDS approach to induce facilities-based competition over NGA networks

Tselekounis, Varoutas and Martakos [20] propose an innovative approach that reflects the current regulatory framework in the European NGA market as described by the EC Recommendation. In particular, the proposed approach models the basic principles of the EC Recommendation and then assesses its effectiveness on inducing facilities-based competition over NGA networks. This implies that this approach can be included in the literature that departs from assessing the efficiency outcomes of the regulation of the copper access networks. The aim of the proposed approach is to meet the current and the future regulatory goals by tackling the initial trade-off between encouraging the incumbents to invest in NGA networks and fostering competition, while incentivizing the entrants to gradually climb the ladder of investment when the NGA investment is proven to be successful. Therefore, the proposed approach provides a theoretical approach to encourage the deployment of a nationwide NGA network (i.e. maximize the potential investment outcome in terms of geographic coverage) with the ambition to finally reflect the socially desirable choice as reflecting in an effective migration path towards facilities-based competition over NGA networks.

The structure and the implementation of the proposed approach are based on the basic principles governing a Credit Default Swap (CDS). A CDS contract is an agreement between two parties, the protection buyer and the protection seller. The first party to the contract, the protection buyer, wishes to insure against the possibility of default on a bond issued by a particular company. The company that has issued the bond is called the reference entity. The second party to the contract, the protection

seller, is willing to bear the risk associated with default by the reference entity. The protection buyer of the CDS makes a series of payments (the CDS "fee" or "spread") to the protection seller and, in exchange, receives a payoff in the event of a default by the reference entity. If a default does not occur over the life of the contract, the contract expires at its maturity date, and hence, the protection seller does not make any payments to the protection buyer.

In an NGA context, the incumbent, which invests in NGA networks, and the regulator agree on a business plan that allows the incumbent to recover the investment in a nationwide NGA deployment during a certain period of time. If the investment has not been recovered at the end of this period, the regulator commits itself that it will compensate the incumbent for the unrecovered part of the investment. After the end of this period, no regulatory remedies will be imposed to the incumbent (sunset clause). In exchange, the incumbent should make periodic payments to the regulator. However, the regulator chooses to subtract this amount from the payments that an access seeker makes to the incumbent in order to have access to the NGA networks. This implies that the incumbent does not pay a periodic fee to the regulator but he subtracts this amount from the access payments he receives. If, however, the investment has been recovered before the end of the clause, the regulator does not make any payment to the incumbent, the incumbent stops making indirect periodic payments to the access seeker and no remedies imposed to the incumbent. In such contract, the incumbent is the protection buyer and the regulator is the protection seller which will compensate the incumbent in the case of a default event (i.e. if the investment has not been recovered at the end of the pre-determined period).

In addition, the model proposes that the contract commits the regulator to apply a certain policy during the whole pre-determined period. This policy, which concerns the derivation of the access pricing formula as well as its evolution over time, is known to the incumbent *ex ante*. In particular:

At time $t=0$ the incumbent and the regulator agree on a business plan that allows the former to have recovered the investment in NGA networks at time $t=T$ with a given probability. Or, in other words, they estimate the probability of default ($P_0$) as well as the corresponding unrecovered part of the investment ($X_0$) at the end of the pre-determined period. The subscript "0" denotes the values of the parameters $P_t$ and $X_t$, $t \in [0,T]$, at the time that the estimation takes place (i.e. $t=0$ in this case). Based on the estimated values of $X_0$ and $P_0$, they assess the amount of the periodic payments ($K_0$) that the incumbent should make to the regulator. This implies that if the estimated demand parameters at $t=0$ coincide with the actual ones during the whole predetermined period $T$, the total amount of the periodic payments will be $TK_0$. However, the regulator chooses to not receive such payments but to subtract this amount from the access payments. Therefore, the reduced access payments that the entrant will finally make to the incumbent from $t=0$ to $t=T$ are given by:

$$AP_{0-T} = W_C \sum_{t=0}^{t=T} Q_t^E - TK_0 \tag{1}$$

or

$$AP_{0-T} = (W_C - R_o) \sum_{t=0}^{t=T} Q_t^E \tag{2}$$

where $W_C$ denotes the cost-based access price, $Q_t^E$ represents the estimated number of consumers served by the entrant at time $t$ and $R_0$ is a regulatory parameter such that

$$R_0 \sum_{t=0}^{t=T} Q_t^E = TK_0$$

(3)

Therefore, the access price that the entrant pays to the incumbent during the $T$ years is given by $W_0 = (W_C - R_o)$. However, the regulator reviews this access price at pre-determined periods. In each periodic review the regulator may increase or decrease the access price according to whether the NGA investment (at the time of each review) is more successful (i.e. an upside case) or less successful (i.e. a downside case) than the initial estimations.

It is shown that in an upside (respectively, downside) case, the implementation of the basic principles governing a CDS contract requires a proper increase (respectively, decrease) in the access price through a proper decrease (respectively, increase) in the regulatory parameter $R_0$. Therefore, an endogenous access pricing rule encourages the entrants to climb the ladder of investment in each upside case. On the contrary, such endogenous access pricing rule provides the entrants with disincentives to invest in each downside case. However, in the latter case, the regulator's goal is to increase the total demand rather than to incentivize the entrant to invest in NGA networks. The reason is that the entrant invests in NGA networks only when the NGA investment is successful. Therefore, the regulator should first promote the success of the NGA investment and then encourage the entrant to invest in its own facilities. It is obvious that in the downside cases the proposed approach fulfills in enhancing the diffusion process since a lower access price facilitates service-based competition over NGA networks. As a result, such an access pricing policy increases the probability of an upside case in the next regulatory review.

Therefore, the proposed approach will eventually lead to the recovery of the NGA investment at the end of the pre-determined period or even earlier. This implies that although its limitations and its potential implementation shortcomings, the proposed approach, which is based on the basic principles governing a Credit Default Swap (CDS), tackles the initial trade-off between encouraging the incumbent to invest in NGA networks and fostering competition, while it incentivizes the entrant to gradually climb the ladder of investment. As a result, the proposed approach represents an effective path towards facilities-based competition over NGA networks.


## 3    CONCLUSIONS

The telecommunications industry is the most rapidly evolving network industry since it has undergone extensive changes in recent decades. Although these changes are mainly related to technological advancements, the regulatory policy has played a significant role in the promotion of competition and innovation. This thesis models the regulatory intervention in the telecommunications market and derives the access pricing policy that achieves the efficiency goals in each of the three major migration phases in the structure of the telecommunications industry.

Firstly, the framework during the migration from a state monopoly market to a competitive telecommunications industry is modeled in order to study the impact of

access prices on the entrant's incentives to undertake the efficient make-or-buy decision in terms of both productive and allocative efficiency. It is found that the particular model of competition that describes the competition in the retail market significantly affects the effectiveness of access prices to achieve static efficiency. However, cost-based access regulation, which has been widely adopted by the regulatory authorities, is found to promote both productive and allocative efficiency regardless of the competition conditions. Therefore, theoretical modeling shows that usage cost-based prices achieve the past regulatory goal concerning the promotion of static efficiency.

Secondly, this thesis reviews the research articles which study the effectiveness of cost-based access prices on achieving the current regulatory goal to promote service-based competition over NGA. The related literature concludes that mandating access to NGA networks at usage cost-based prices discourages both incumbents and entrants to invest in such networks. Therefore, the research focuses on studying alternative regulatory schemes that may promote both investments and competition. The most significant deviation from the permanent regulation of access at usage cost-based prices concerns the implementation of non-cost-based access prices.

The related literature concludes that investment-contingent access prices can induce the incumbent to undertake to socially optimal investments in NGA networks (i.e. promote both static and dynamic efficiency) under certain conditions concerning the demand and cost structure. However, these studies do not take into account the fact that regulators have significant incentives to deviate from such schemes once NGA networks have been deployed by setting a cost-based access price in order to maximize social welfare. This thesis models this fact in order to study the impact of regulatory uncertainty on an incumbent's incentives to undertake the efficient investments in NGA networks. It is found that the feasibility of the socially optimal outcome is not only affected by the demand and cost structure, but also by the perceived regulatory uncertainty.

Thirdly, this thesis points out that the current and the future regulatory goal of promoting dynamic efficiency through facilities-based competition are closely related, and hence, a combined regulatory policy should be applied. As a result, it proposes an innovative regulatory approach which is based on the basic principles governing a CDS contract. It is shown that under quite general but plausible assumptions about demand and cost factors, the proposed approach can induce an efficient migration towards facilities-based competition over NGA networks. It is thus obvious that this thesis not only discusses the past, the present and the future state of telecommunications networks, but also significantly contributes to the literature which studies the optimal access pricing policy that achieves the past, the current and the future regulatory goals.

## References

1. M. Bourreau and P. Doğan, "Regulation and innovation in the telecommunications industry," Telecommunications Policy, vol. 25, no. 3, pp. 167–184, 2001.
2. P. De bijl and M. Peitz, Regulation and Entry into Telecommunications Markets. Cambridge University Press, 2002.

3. A. Oldale and J. Padilla, "From state monopoly to the '"investment ladder"': competition policy and the NRF. Swedish Competition Authority Series: The Pros and Cons of Antitrust in Deregulated Markets," Stockholm, 2004.

4. M. Cave, "Encouraging infrastructure competition via the ladder of investment," Telecommunications Policy, vol. 30, no. 3–4, pp. 223–237, Apr. 2006.

5. FCC (Federal Communications Commission), "1996 Telecommunications Act," Pub. LA. No. 104-104, 110 Stat. 56,1996.

6. European Parliament and Council, "Regulation (EC) No 2887/2000 of the European Parliament and of the Council of 18 December 2000 on unbundled access to the local loop," 2000.

7. M. Armstrong, "The theory of access pricing and interconnection," in Handbook of Telecommunications Economics: Vol. 1, structure, regulation and competition, no. 15608, M. Cave, S. K. Majumdar, and I. Vogelsang, Eds. Amsterdam: North Holland, 2002, pp. 295–386.

8. M. Tselekounis, D. Varoutas, and D. Martakos, "On the irrelevance of input prices from a regulatory perspective," in 5th International conference on competition and regulation (CRESSE), 2-4 July 2010, Chania, Greece.

9. M. Tselekounis, D. Varoutas, and D. Martakos, "On the social optimality of make-or-buy decisions," Journal of Regulatory Economics, vol. 41, no. 2, pp. 238–268, Sep. 2012.

10. P. W. Shumate, "Fiber-to-the-Home: 1977–2007," Journal of Lightwave Technology, vol. 26, no. 9, pp. 1093–1103, May 2008.

11. N. Czernich, O. Falck, T. Kretschmer, and L. Woessmann, "Broadband Infrastructure and Economic Growth," The Economic Journal, vol. 121, pp. 505–532, 2011.

12. ITU, "The Impact of Broadband on the Economy: Research to Date and Policy Issues," 2012.

13. R. Katz, S. Vaterlaus, P. Zenhäusern, and S. Suter, "The impact of broadband on jobs and the German economy," Intereconomics, vol. 45, no. 1, pp. 26–34, Feb. 2010.

14. T. Reynolds, "The role of communications infrastructure investment in economic recovery (DSTI / ICCP / CISP ( 2009 ) 1 / FINAL)," Paris, 2009.

15. J. Bouckaert, T. van Dijk, and F. Verboven, "Access regulation, competition, and broadband penetration: An international study," Telecommunications Policy, vol. 34, no. 11, pp. 661–671, Dec. 2010.

16. M. Tselekounis and D. Varoutas, "Investments in next generation access infrastructures under regulatory uncertainty," Telecommunications Policy, 2013. http://dx.doi.org/10.1016/j.telpol.2013.06.001

17. ITU, "ICT Regulation Toolkit. Module 2: Competition and Price Regulation." [Online]. Available: http://www.ictregulationtoolkit.org/2. [Accessed September 7, 2013].

18. M. Tselekounis, D. Maniadakis, and D. Varoutas, "NGA Investment Incentives under Geographic Price Discrimination," in 40th EARIE Conference, 30 August-1 September 2013, Évora, Portugal.

19. M. Tselekounis, E. Xylogianni, D. Varoutas, and D.Martakos, "Geographically Differentiated NGA Deployment," accepted in 24th European Regional Conference of the International Telecommunications Society (ITS), 20 - 23 October 2013, Florence, Italy.

20. M. Tselekounis, D. Varoutas, and D. Martakos, "A CDS approach to induce facilities-based competition over NGA networks," submitted to Telecommunications Policy (under 3rd round revision), 2013.

21. D. E. M. Sappington, "On the Irrelevance of Input Prices for Make-or-Buy Decisions," American Economic Review, vol. 95, no. 5, pp. 1631–1638, Dec. 2005.

22. P. G. Gayle and D. L. Weisman, "Are input prices irrelevant for make-or-buy decisions?," Journal of Regulatory Economics, vol. 32, no. 2, pp. 195–207, Apr. 2007.

# Transfer and Management of Broadband Traffic in Wireless Telecommunication Networks

Nicholas Vaiopoulos*

National and Kapodistrian University of Athens

Department of Informatics and Telecommunications

nvaio@di.uoa.gr

**Abstract.** The present thesis aims to study resource allocation techniques in fixed wireless networks in order to optimize their spectral efficiency as well as to investigate alternative broadband transfer methods. A brief overview of broadband wireless access networks with their characteristics is given and the fundamental resource management techniques proposed in the open technical literature are also referred. Then, a typical wireless broadband network model is presented and a comprehensive review of the key resource management techniques proposed in the literature is given. Following this, the proposed resource management techniques are presented and compared using proper simulation results.

On the other hand, the analytical model of a wireless broadband traffic model over a terrestrial wireless optical link is analyzed. Its performance is extracted using analytical expressions of the average outage probability and the average error probability metrics. Appropriate simulation results are depicted as well. Furthermore, a network architecture comprising of several high amplitude platforms communicating with each other using optical links, is introduced in order to transfer broadband traffic over long distances. The outage probability performance is examined using either one-hop or multi-hop scenarios and suitable numerical results are provided. Finally, concluding remarks are summarized and suggestions for further research are indicated.

**Keywords**: Resource management, interference management, free space optics, fading, turbulence, pointing error.

## 1    Dissertation Summary

The vision of 'broadband to all' has necessitated the deployment of fixed wireless access (FWA) as an alternative technology in order to provide broadband access to geographical areas where the cost of wired infrastructure is extremely high. Such a technology is the Worldwide Interoperability for Microwave Access, commonly known as WiMAX, which is based on the IEEE 802.16 standard. The continuously

---

* Dissertation Advisor: Dimitris Varoutas, Assistant Professor

growing demand for higher data rates and bandwidth-consuming services becomes a driving force toward larger channel widths and wider spectrum block allocation consequently. Traditional frequency planning with high reuse patterns wastes the limited available spectrum, especially in the lower regions of the WiMAX frequency operation band. Therefore a full frequency reuse in each sector of each cell is a very attractive alternative to fulfill this challenge.

However, such an approach results in high co-channel interference (CCI), which arises from concurrent intracell and intercell transmissions, and affects significantly the users' quality of service (QoS). Several approaches have been proposed in the literature for CCI reduction. An increased effort is focused on the corresponding time domain radio resource allocation techniques (RRA) [1–4] in order to organize the total amount of interference.

Another critical issue is the transfer of broadband traffic (e.g. WiMAX based traffic) in remote terminal stations in order to improve coverage, flexible access and reduce the cost of deployment. An effective solution is the transportation of radio signals between a central base station and multiple radio access units in optical form and the transmission through a fiber optic. Radio over Fiber (RoF) transmission [5] has a number of advantages but installation cost may be prohibited. Hence, it is not always feasible its deployment in practice. In this case, the transmission of radio signals on Free Space Optics (FSO) links combines the benefits for ease deployment in wireless links and high capacity enabled by fiber optic technologies. Several Radio over FSO (RoFSO) systems studies have been recently presented in the literature [6-9].

Obviously, as the terminal stations are located in increasingly greater distances satellite communications becomes a dominant alternative for broadband traffic transferring. Nevertheless, excessive high power requirements and high installation costs are serious disadvantages and an alternative high challenging solution is to employ a network consisting of high altitude platforms (HAPs). HAPs combine some of the most distinctive characteristics of terrestrial wireless and satellite communication systems, e.g., broad service areas, great capacity, low transmission delay, adequate power consumption, etc. They are located in the stratosphere, approximately 25km above the ground, and remain stationary, maintaining, thus, the same behavior as geostationary satellites. However, the short distance between HAPs and ground stations, lead to lesser power demands and much smaller round trip delays, making this technology quite attractive for broadband services in next-generation wireless communications [10].

The overall performance of outdoor FSO systems depends upon the climatological conditions and the general characteristics of the transmission paths. Furthermore, the transfer of broadband traffic over FSO links has not been adequately investigated in the literature. That was one of the motivations for this thesis and the contributions include detailed performance analysis of WiMAX RF signals transferring through two alternative FSO technologies. At first, an FSO terrestrial link is considered where turbulence effects for the FSO subchannel and composite fading effects for the RF subchannel are adopted [11,12]. Then, a multi-hop HAP network is examined [13]. These HAPs take the role of terrestrial base stations and collect the WiMAX traffic from the area they cover. They have transparent transponders that convert the WiMAX signals to optical ones and the reverse. The optical signals are transmitted from the source to the destination HAP through inter-HAP links and the traffic is delivered

by this way to the end users after RF conversion. In addition, the optimization of spectral efficiency for the downlink segment of a FWA network is examined and three effective resource allocation methods are proposed and analyzed. The basic idea of the first proposed algorithm is the avoidance of the major interferers by using different allocation schemes for the even and odd sectors [14, 15]. This algorithm is enhanced with the adoption of multi-mode modulation schemes [16]. Finally, the use of alternate polarization allocation (PA), as a means to decrease the interference into the desired signal in conjunction with a time domain RRA technique is proposed in [17, 18]. The proposed architectures were evaluated with realistic sets of parameters and closed form expression and simulation results are presented.

## 2 Results and Discussion

### 2.1 A RRA scheme for FWA systems with avoidance of major interferers

As it has been proposed in [3], ESRA examines a packet switched broadband wireless network using time division multiple access (TDMA) technique and time division duplexing (TDD) with full frequency reuse. The service area is divided in hexagonal cells and sectors are labeled from 1 to 6, counter-clockwise, in such a way that there are no adjacent sectors bearing the same label (Fig. 1).
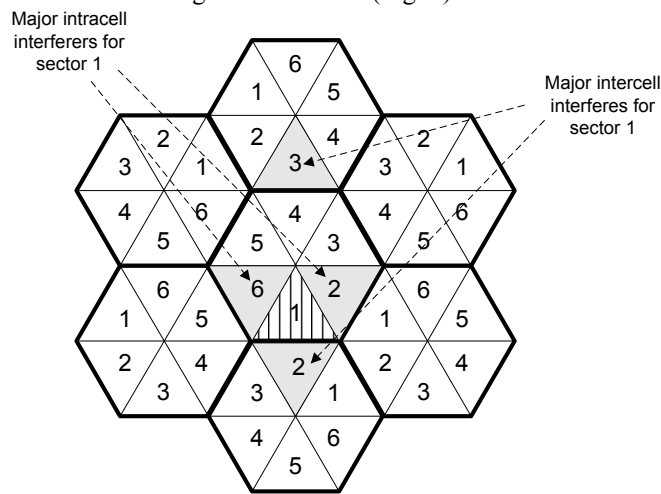


**Fig. 1.** Major interferers for the downlink direction of the hexagonal cell layout.

The frame is divided into six subframes (SBs), which are further divided into mini-frames (MF) labeled from 1 to 6. Each sector schedules packets for transmissions in available MFs of each SB following the staggered order of Fig. 2a.
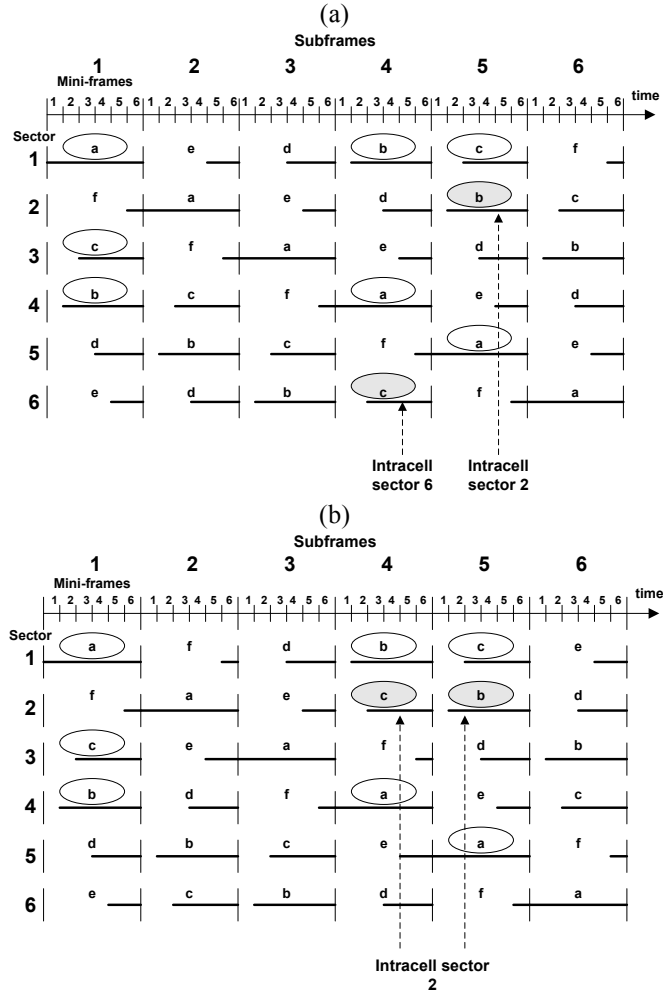
(a)

(b)

**Fig. 2.** Three concurrent transmissions for sector 1 according to (a) ESRA and (b) proposed method.

According to this order, sector 1 schedules packets for transmissions in SB 1 (denoted by a). For further packet transmissions, exploiting Base Station (BS) directional antennas, it uses SB 4, which is the first SB of the opposite sector. The next options for sector 1 will be the available MFs of the first SBs of the other two opposite sectors (i.e. sectors 5 and 3) clockwise and the last two options will be the available MFs of the first SBs of the adjacent intracell sectors (i.e. sectors 2 and 6) clockwise too. The same procedure is applied to the other sectors. Terminals are classified into six classes according to the number of maximum tolerable concurrent transmissions following the staggered order of Figure 2a. For example, a class 3 terminal of sector 1 tolerates the following three concurrent transmissions: sectors 1-4-3 in SB 1, sectors 4-1-6 in

140

SB 4 and sectors 5-2-1 in SB 5 (Fig. 2a). The degree of concurrent transmissions of each MF is defined on its label. More precisely, MF 1 allows only 1 packet transmission; MF 2 allows 2 concurrent packet transmissions and so on. Terminal classification is based on the reception quality of their location that depends on antenna characteristics (3dB beamwidth and Front-To-Back (FTB) ratio), shadowing and distance from serving BS. The reception quality may be improved through a macrodiversity procedure (i.e. selection of serving sector). Each MF with the same label has the same number of timeslots in each SB (MF $i$ has $n_i$, $i = 1, 2,...,6$ timeslots in each SB). The size of MFs is chosen to match the expected traffic load and is a mechanism to increase the overall throughput.

For a typical FWA system, throughput enhancement is achieved by upgrading the number of higher classes' terminals through advanced methods of major intercell and intracell interferers' avoidance. As shown in Fig 1, the major interferers for sector 1, under a simple path loss model, comes from the shadowed intracell sectors 2 and 6 (due to overlapping sector antenna patterns) as well as from shadowed intercell sector 3 (because of the front lobe of sector antenna 3 that point directly to the terminals of the tagged sector) and the opposite shadowed sector 2. Following the ESRA staggered order (Figure 2a) for higher terminal classes (i.e. classes 4, 5 and 6) the major intracell interferers are appearing not only solely but also together.

The proposed method [14, 15] is based on a different allocation scheme for odd and even sectors, which is examined and proposed as follows. The sector 1 schedules packets for transmission in SB 1, as shown in Fig. 2b. If there are more packets for transmission, it uses the available MFs of SB 4, which are the first SB of the opposite sector, in order to exploit the BS directional antennas and the low level of interference. All sectors follow this procedure for the first two SBs. However, the next two options for sector 1, according to the staggered order will be the available MFs of the first SBs of the other two opposite sectors (i.e. sectors 5 and 3) clockwise and the last two options will be the available MFs of the first SBs of the two adjacent intracell sectors (i.e. sectors 6 and 2) counter-clockwise. This concept is applied to the odd labeled sectors (i.e. sectors 3 and 5). On the contrary, for sector 2, the next two options will be the available MFs of the first SBs of the other two opposite sectors (i.e. sectors 4 and 6) counter-clockwise while the last two resorts will be the available MFs of the first SBs of the two adjacent intracell sectors (i.e. sectors 3 and 1) clockwise. This procedure is repeated for the even labeled sectors (i.e. sectors 4 and 6). This situation is exploited by the proposed method in order to increase the number of class 3 terminals due to the fact that each sector is interfered only by one adjacent intracell sector in contrast with ESRA where each sector is interfered by both adjacent intracell sectors one after the other. More explicitly, in the case of three concurrent transmissions (class 3 terminals) of sector 1, Fig. 2a presents the allocation scheme based on ESRA and Fig. 2b the proposed scheme. Sector 1, following the ESRA staggered order, is interfered in SB 4 by sector 6 (one dominant intracell interferer) and in SB 5 by sector 2 (the other dominant intracell interferer). On the other hand, following the proposed staggered order, sector 1 is interfered in SBs 4 and 5 by sector 2 (i.e. the same dominant intracell interferer). As a result, the proposed method upgrades the

fraction of terminals that tolerate three concurrent transmissions enhancing the maximum throughput per sector.

## 2.2 Performance improvement of FWA systems using multi-mode modulation schemes

The classification procedure presented in the previous section is based on the worst case scenario, so that all users are guaranteed a minimum QoS for a given signal to interference ratio threshold ($SIR_{thr}$). However, this procedure produces a different signal to interference ratio (SIR) for each terminal in each SB. Indeed, each class-$i$ terminal has $i$ respective SIRs in $i$ different SBs. Therefore, a SIR margin becomes available for each terminal and i different SIR margins may be expected from SB to SB. It is well known that a higher SIR supports a higher modulation mode and therefore a larger number of bits per symbol. Approaches proposed in [3,14,15] ignore these SIR margins and adopt a single modulation mode according to the specific $SIR_{thr}$ used in the classification procedure. The motivation of [16] is to examine the possibility of utilizing the above SIR margins by adopting higher modulation modes. A performance enhancement in each sector is then expected.

The adoption of different modulation modes affects the TDMA frame structure (Fig. 3). The time sharing of MF among terminals using the same modulation mode induces its further partition into $N_{modes}$ micro-frames, where $N_{modes}$ is the number of the adopted modulation modes. It is noticed that one mode is assigned to each micro-frame. In [16] it is shown that the adoption of multiple modulation modes induces a maximum throughput per sector enhancement exceeding 50%.
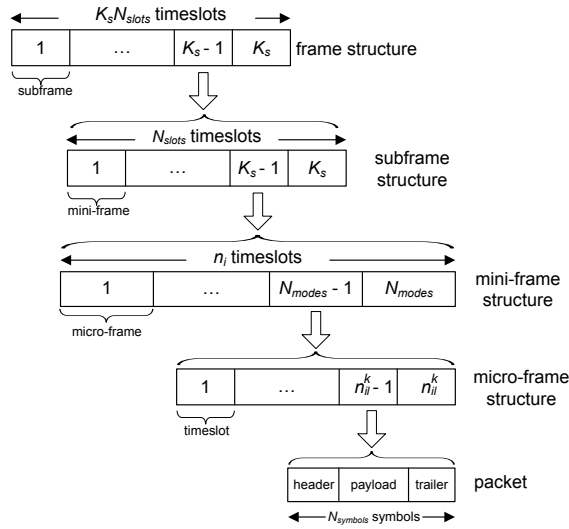


**Fig. 3.** TDMA frame and proposed MF structures.

### 2.3 Performance improvement of FWA systems by conjunction of dual polarization and time domain RRA technique

The CCI reduction is a major challenge and the proposed RRA algorithms significantly improve the maximum throughput per sector. However, a further improvement would be achieved by jointly utilizing a scheduling algorithm and an alternate polarization allocation (PA) scheme as presented in [17, 18].

The service area is divided in hexagonal cells and triangular sectors. Each sector is equipped with its own directional antenna and labeled from 1 to 6 counter-clockwise, in such a way that there are no adjacent sectors bearing the same label (Fig. 4). Each sector is connected with the IP backbone network through a switching module. The labels of three adjacent cells are rotated by 120°, in respect one another, creating a cluster (contained in the heavy line in Fig. 4) whose pattern is repeated across the entire service area. Adjacent sector antennas use alternate polarization so as to maximize isolation among them and increase the communication quality. Users use rooftop directional antennas pointing to their respective sector antenna following the polarization pattern. The beamwidth of a sector antenna has to be wide enough in order to serve the entire sector, whereas the beamwidth of each terminal antenna must be more directional to lower interference.
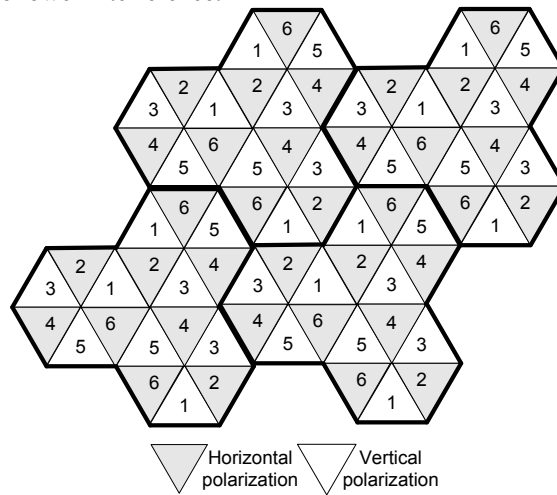


**Fig. 4.** A representative part of the hexagonal cell layout consisting of 4 clusters.

Considering an FWA network without a PA pattern, it is clear that the most significant amount of intracell interference, for each sector, is originating from its two adjacent sectors due to the overlapping sector antenna patterns. However, under the proposed framework, each sector antenna is transmitting with orthogonal polarization with respect to its adjacent sector antennas. Hence, the amount of intracell interference is significantly reduced and the maximum throughput per sector is further increased.

## 2.4    Transfer WiMAX signals via terrestrial optical wireless links

In [11, 12] we consider a terrestrial FSO link which is used to deliver WiMAX traffic from one geographic region to another. The overall system configuration is composed of the optical and the wireless subsystems. We assume that the WiMAX traffic from heterogeneous networks reaches the optical transmitter through an access gateway. The transmitter (Fig. 5b) converts the electrical signal to laser. It is composed of a modulator, a laser driver, a light-emitting diode (LED) or laser, and a telescope as a whole. The laser propagates through the atmosphere to the receiver assuming a Gaussian beam wave model.

The receiver (Fig. 5b) uses a direct detection scheme and includes a telescope, a filter, a positive–intrinsic–negative (PIN) photodetector, and a trans-impedance amplifier. Depending on cost restrictions and reliability requirements, a tracking and pointing subsystem may be implemented in both sides of the communication link to maintain transmitter–receiver alignment.

The electrical signal is guided to a WiMAX base station and delivered to the users located there. The WiMAX standard is based on the OFDM standard utilizing a large number of closely spaced orthogonal subcarriers.
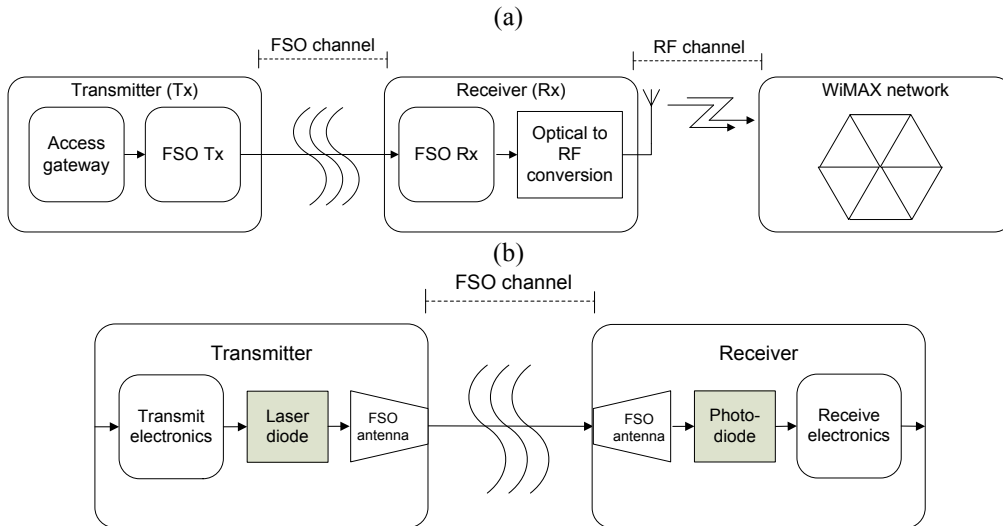


**Fig. 5.** (a) Optical and RF subsystem (b) Optical subsystem.

An appropriate channel model is adopted, which entails some of the most critical impairments of the optical channel, i.e., attenuation, turbulence, pointing error effects, as well as of the RF channel, i.e., path loss, shadowing, and fast fading, is taken into account. The overall link budget and a closed-form of the outage probability of the system are deduced. Several analytical results are depicted using a realistic set of parameter values, to lend a helpful insight to the performance of the proposed architecture.

## 2.5 Transfer WiMAX signals using an FSO multi-hop HAP network

In [13], we present a novel HAP network architecture which aims at delivering WiMAX services to extremely far distances on Earth using multi-hop routing. HAPs in the network take the role of terrestrial base stations and collect the WiMAX traffic from the area they cover. They have transparent transponders that convert the Wi-MAX signals to optical ones and the reverse. The optical signals are transmitted from the source to the destination HAP through inter-HAP links and the traffic is delivered by this way to the end users after RF conversion. In such an architecture, we determine the WiMAX quality of service (QoS) by minimizing the outage probability for the network configuration.

The source HAP (Tx HAP in Fig. 6) communicates with the destination HAP through $R_i$, $i=1,2,\ldots,N$ -1 optical transceivers, which act as relays-nodes all being in equidistance $d_O$, i.e., there are $N$ point-to-point propagation links before the laser signals arrive to the destination. Relay assisted transmission is a common technique in wireless RF communication systems since it provides a broader and more efficient coverage and can be used as a fading mitigation tool. Every intermediate node in a multi-hop network acts as a router that forwards traffic towards its destination. In [13], we consider N-1 relays where each one has knowledge of the channel state of the previous hop. We assume the use of amplify-and-forward (AF) relays which just amplify and forward the incoming signal without performing any sort of decoding. These relays use less complex circuitry compared to decode and forward (DF) ones, which decode the signal and then transmit the detected version to the destination ones.

The destination HAP (Rx HAP in Fig. 6) receives the optical power using a telescope. Optical light can be concentrated by using lenses and mirrors, or any combination of them. A filter helps to remove the background radiation from entering the Rx which generally creates shot noise and saturates the detector. The optical signal from the output of the filter propagates to the detector, which converts it to an RF electrical signal using a photodetector [19]. Finally, the RF electrical signals are delivered to the end users located in the HAP Rx area.
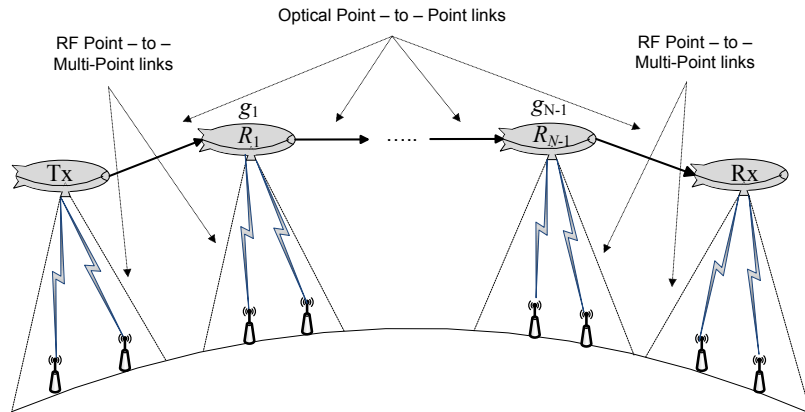


**Fig. 6.** HAP network configuration.

The overall performance is examined by using a channel model which incorporates laser path loss as well as pointing error effects. A closed-form of the outage probability of the system is extracted and several analytical results are depicted adopting a realistic set of parameter values.

## 3    Conclusions

In the context of this dissertation we examined advanced RRA methods for the downlink of a FWA network as well as the transfer of broadband traffic via optical wireless links.

More precisely, we presented a RRA which improves the ESRA method in terms of throughput per sector by increasing the number of terminals that tolerate more concurrent transmissions. For a typical radio environment, a 10% increment of the maximum throughput per sector with respect to ESRA method is achieved. It's worth mentioning that the proposed RRA scheme performs better under realistic transmission conditions and with various types of antennas and exploits better low performance antennas, which present large beamwidths and FTB ratios.

Furthermore, the SIR margin, induced by the terminal classification procedure, in a TDMA FWA system is examined and analyzed. An advanced frame structure is proposed, facilitating the utilization of multiple modulation mode schemes. It has been shown that a significant throughput improvement is achieved, for various $SIR_{thr}$ values, used in terminal classification procedure.

The third contribution includes an integrated time domain RRA technique of concurrent transmissions and polarization alternation pattern for the downlink direction of an FWA system. The proposed scheme presents an enhanced performance since the PA pattern incorporation reduces significantly the impact of dominant interferers. It must be noticed that the proposed scheme performs better under worst propagation conditions and exploits better low-performance antennas, which present large beamwidths and small FTB ratios.

Moreover, we constructed a simple but adequate architecture to investigate WiMAX transmission over terrestrial FSO channels. The channel model considers the laser link and the WiMAX communication system parameters used in practice. Specifically, some of the most critical impairments of the optical channel, i.e., path loss, turbulence, pointing error effects, as well as of the RF channel, i.e., path loss, shadowing, and fast fading were taken into account, and an analytical derivation of the outage probability was obtained. The feasibility of the proposed architecture was further evaluated with a realistic set of parameter values and depicted using proper graphs. The present architecture may constitute the outset of adopting and evaluating more complicated and, at the same time, more realistic RoFSO deployment scenarios. In this vein, the incorporation of forward error correction schemes in order to increase the overall performance seems to be quite challenging and such an extension is a subject of ongoing research.

Finally, we presented an alternative method to deliver WiMAX services at extremely far distances on Earth by using a HAP network. A source HAP collects the

traffic from the multi-hop routing through a number of intermediate HAPs. HAPs communicate with each other using laser links. After reaching the destination HAP, the traffic is transformed in RF form and delivered to the end users on Earth. The optical HAP transceiver and the laser inter-HAP channel, that takes account of the pointing error statistics, were analytically described. At first, we considered the case where the source directly communicates with the destination and then we generalized to a relayed scenario. We particularly focused on the outage probability at the ground users and presented proper graphical results for the performance evaluation of the network. The obtained results can serve as a guideline for designers to predict and evaluate a HAP network ability to deliver broadband services in practice. The analysis conducted assuming a typical parameter set mainly used in practical systems. However, a thorough investigation is necessary in order to find the appropriate set for better performance. The choice for example of optimum gains and the beam divergence angle is a crucial point before the implementation process. Random angular jitter affects the overall performance and therefore proper tracking systems need to be considered. The present analysis can be extended in a number of ways. For instance, it would be interesting to consider nonlinear laser diodes and examine the effect of intermodulation distortion which is often present in practical laser links. Another extension is the consideration of turbulence in optical links for large inter-HAP distances. Fading may also be included in the RF downlink link. These additions would increase the mathematical complexity of the model but on the other hand make it much more reliable. Since practical HAP networks are currently working using microwave links, the idea of evaluating the performance of RF multi-hop links would be of particular help for comparison reasons. Moreover, different types of relays may be used, e.g., DF, all-optical relays, etc. Finally, the investigation of a coexistence scenario between heterogeneous HAP and satellite networks for increased overall performance appears quite challenging.

# 4    References

1.  X. Qiu and K. Chawla, "Resource Assignment in a Fixed Broadband Wireless System," IEEE Commun. Lett., pp. 108-110, July 1997.
2.  T. K. Fong, P. S. Henry, K. K. Leung, X. X. Qiu and N. K. Shankaranarayanan, "Radio resource allocation in fixed broadband wireless networks," IEEE Trans. on Commun., pp. 806-818, June 1998.
3.  K. K. Leung and A. Srivastava, "Dynamic allocation of downlink and uplink resource for broadband services in fixed wireless networks," IEEE J. on Sel. Areas in Commun., pp. 990-1006, May 1999.
4.  V. Tralli, R. Veronesi and M. Zorzi, "Power-shaped advanced resource assignment (PSARA) for fixed broadband wireless access systems," IEEE Trans. on Wireless Commun., pp. 2207-2220, Nov. 2004.
5.  H. Al-Raweshidy and E. S. Komaki, Radio over Fiber Technologies for Mobile Communication Networks, Artech House, 2002.
6.  N. Cvijetic and T. Wang, "WiMAX over free-space optics evaluating OFDM multi-subcarrier modulation in optical wireless channels," in Proc. of IEEE Sarnoff Symp., Princeton, NJ, USA, 2006.

7. N. Cvijetic and T. Wang, "A MIMO architecture for IEEE 802.16d (WiMAX) heterogeneous wireless access using optical wireless technology," in Proc. of Next Generation Teletraffic and Wired/Wireless Advanced Networking (NEW2AN), St. Petersburg, Russia, 2006.

8. A. Bekkali, C. B. Naila, K. Kazaura, K. Wakamori and M. Matsumoto, "Transmission analysis of OFDM-based wireless services over turbulent radio-on-FSO links modeled by gamma-gamma distribution," IEEE Photon. J., p. 510–520, June 2010.

9. S. Arnon, "Minimization of outage probability of WiMAX link supported by laser link between a high-altitude platform and a satellite," J. Opt. Soc. Am. A, p. 1545–1552, July 2009.

10. D. Grace and M. Mohorcic, Broadband Communications via High Altitude Platforms, John Wiley & Sons, 2011.

11. N. Vaiopoulos, H.G. Sandalidis, D. Varoutas "WiMAX on FSO: Outage Probability Analysis", IEEE Trans. Commun., vol. 60, no. 10, pp. 2789–2795, Oct. 2012.

12. N. Vaiopoulos, H.G. Sandalidis, D. Varoutas, "Transferring WiMAX Signals via Terrestrial Optical Wireless Links" submitted for publication.

13. N. Vaiopoulos, H.G. Sandalidis, D. Varoutas "Using a HAP Network to Transfer WiMAX OFDM Signals: Outage Probability Analysis", to appear in IEEE/OSA J. Optical Comm. Netw.

14. N. Vaiopoulos, A. Vavoulas, D. Varoutas and T. Sphicopoulos, "A Radio Resource Allocation Scheme for Fixed Broadband Wireless Access Systems with Avoidance of Major Interferers," Wireless Personal Commun., Springer, vol. 40, No. 4, pp.479-487, March 2007.

15. A. Vavoulas, N. Vaiopoulos, D. Varoutas and G. Stefanou, "'Throughput enhancement of Fixed Broadband Wireless Access systems through a novel Radio Resource Allocation approach," in Proc. of European Wireless Conference (EWC2006), Athens, Greece, April 2006.

16. N. Vaiopoulos, A. Vavoulas, D. Varoutas and G. Stefanou, "Performance improvement of Fixed Cellular Networks using multi-mode modulation schemes," in Proc. of 18th IEEE International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC'07), Athens, Greece, September 2007.

17. A. Vavoulas, N. Vaiopoulos, D. Varoutas and G. Stefanou, "Co-channel interference reduction through integrated scheduling and polarization allocation for Fixed Cellular Systems," in Proc. of 18th IEEE International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC'07), Athens, Greece, September 2007.

18. A. Vavoulas, N. Vaiopoulos, D. Varoutas, A. Chipouras, G. Stefanou "Performance improvement of Fixed Wireless Access networks by conjunction of dual polarization and time domain radio resource allocation technique", Int. J. Commun. Syst., vol. 24, No. 4, pp. 483 – 491, April 2011.

19. S. Arnon, "Optical wireless communications," in Encyclopedia of Optical Engineering, 2003.

# Connectivity Issues for Optical Wireless Networks

Alexander Vavoulas*

National and Kapodistrian University of Athens

Department of Informatics and Telecommunications

`vavoulas@di.uoa.gr`

**Abstract.** Optical Wireless Communications (OWC) networks are becoming more and more popular for delivering broadband traffic since they are introducing significant advantages against the other alternative technologies. Operating wavelengths range from ultraviolet (UV) to the infrared (IR) portion of the electromagnetic spectrum and present significant attenuation from channel impairments. As a result, the transmission range is significantly reduced when a single hop is used.

Therefore, multi-hop operation, which is a common technique in wireless radiofrequency (RF) communication systems, is adopted in order to increase the effective distance between transmitter and receiver. To improve their reliability connectivity issues need to be investigated. Connectivity has been investigated in RF ad hoc networks (either one or two dimension) in contrast with OWC networks. The present dissertation aims to examine this research area by connecting the minimum transmission range ensuring connectivity with a plethora of parameters such as the adopted modulation and/or coding format, the transmitted power, the supported data rate and the error probability. Analytical expressions are extracted and the derived results are depicted using appropriate figures. The outcomes constitute a valuable tool to design such networks in practice.

**Keywords**: $k$-connectivity, node isolation probability, multi-hop networks, UV-C transmission, NLOS propagation, weather effects, underwater optical wireless networks.

## 1    Dissertation Summary

Optical wireless has been launched as an attractive candidate technology to provide broadband communications. The way optical wireless transceivers operate is more or less the same as fiber optics ones; however, since laser signals are now transferred through the atmosphere, the path loss between the transmitter and the receiver is getting raised due to a plethora of pernicious factors that appear. In this dissertation three types of optical wireless networks have been examined: non-line-of-sight (NLOS) optical transmission in the ultraviolet UV-C spectral region, free space optics (FSO)

---

* Dissertation Advisor: Dimitris Varoutas, Assistant Professor

in the infrared spectral region and underwater optical wireless transmission in the visible spectral region. More details on these networks are given in the following.

Transmission in the ultraviolet UV-C spectral region (particularly from 200nm to 280nm), also known as the solar-blind band, exhibits some unique characteristics. Firstly, most of the solar radiation is getting absorbed by the ozone in the upper atmosphere, leading to almost negligible background noise at the Earth's surface. Secondly, the UV-C light generated from terrestrial sources is strongly scattered due to the presence of suspended particles in the atmosphere. The insignificant background solar radiation and the strong atmospheric scattering enable the activation of NLOS communication links with large field-of-view (FoV) receivers which allow a large amount of scattered light collection. Recent advances in hardware have led to the emergence of low cost semiconductor laser diodes and miniaturized LEDs at UV-C frequencies making this new technology a quite promising solution for short-range communications [1].

On the other hand, FSO transfers broadband services via line of sight (LOS) links and the common candidates operating wavelengths are 780nm, 950nm, and 1550nm. Weather, propagation distance, scattering, absorption, turbulence, pointing error effects, laser wavelength, and data rates are some of the deterministic and random elements that contribute to the overall performance of an optical wireless link. Even if all the deteriorating factors play a significant role, designers and implementers should particularly take meteorological phenomena into consideration when they intend to deploy a robust FSO network in practice. Fog, snow, and rain cause the scattering of laser signals in the atmosphere. Scattering makes a portion of the light beam traveling from a source deflect away from the intended receiver. Another atmospheric effect under clear weather conditions is the turbulence induced by random changes in the atmospheric refractive index. As a result, random phase and irradiance fluctuations (scintillation) of the optical signals are observed at the receiver [2]. Furthermore, the FSO links also depend on the pointing error performance. Pointing errors occur due to mechanical misalignment or errors in tracking systems. Among all phenomena, fog brings about the greatest repercussions since it is constituted of small water droplets having dimensions near the size of infrared wavelengths. Snow and rain also influence the FSO performance, though their impact is significantly less than that of fog. Note that these weather phenomena rarely occur concurrently.

Finally, underwater optical communications were based on simple point to point links delivering composite applications such as monitoring the ocean environment, mapping of the sea floor, and sensing purposes [3]. These systems are usually operating in the blue/green spectrum region where absorption is minimum compared to other wavelengths. The aquatic medium contains almost 80 different elements, dissolved or suspended in pure water, with different concentrations, as well as phytoplanktons, zoo planktons and many marine organisms and plants. These components may redirect the transmitting light or transform it into heat due to the two fundamental physical processes, namely absorption and scattering. Obviously, the light transmittance is sensitive to high wavelengths since these two physical processes have a highly spectral dependence.

The deteriorating factors that are introduced by the propagation channel reduce both the maximum obtainable data rate and the transmission range and an efficacious

solution to mitigate these impairments is to employ relay-assisted techniques. In a multi-hop transmission, the total transmission path is divided into smaller distances between relays or hops which suffer from less loss. At each relay, the received optical field is processed and forwarded to the next one. In that way, an effective serial relayed network can provide services at far distances. Transmission through relays is quite a common practice in wireline and RF wireless communication systems. A plethora of such relevant studies has appeared recently; see, e.g., the newly published books of Uysal [4] or Dohler and Li [5].

Multi-hop networks can properly operate if connectivity between their nodes is satisfied. A fully connected network contains a path from any node to another. When there is no path between at least one source–destination pair, the network is disconnected. Obviously, connectivity plays a critical role for wireless networking. Connectivity was mainly studied for RF wireless networks where the propagation suffers from various severe, random in nature impairments such as path loss, multi-path fading, and shadowing [6], [7]. Another metric of network connectivity is the node isolation probability which can be defined as the probability that a random node cannot communicate with any other nodes. Some of the studies on this topic for one-dimensional relayed networks are as follows. In [8], the probability of having a wireless network composed of at most C clusters is extracted. In [9], an ad hoc network consisting of nodes and base stations is considered, and the probability of node to base station connectivity is derived. Analytical expressions for the probability that a wireless network is connected are presented in [10], as well. Finally, Miorandi and Altman [11] obtained exact results for the coverage probability, the node isolation probability, and the connectivity distance for various node placement statistics.

In the context of optical wireless networking there appears an absence of similar works in the open technical literature mainly due to the fact that the interest on this specific field is nowadays in season. That was the key motivation for this dissertation and the outcomes were the first papers in the areas of ultraviolet communications ([12], [13]), free space optics ([14]) and underwater optical communications ([15]). The numerical results of this dissertation are of significant value for telecom researchers working toward a flexible optical wireless network deployment in practice.

## 2 Results and Discussion

### 2.1 k-connectivity issues

The multi-hop optical wireless network network can be represented as an undirected graph $G$ with a set of vertices $V$ and a set of edges $E$ [16]. The set of vertices has cardinality $n$ and represents the set of nodes, while the set of the edges corresponds to the OWC links between the nodes. The node degree, $d(u)$, is defined as the number of links of a node (i.e., the number of neighbor nodes within its range). An isolated node has a null node degree. The minimum node degree, $d_{min}$, of $G$ is defined as the minimum value among the node degrees. Since the nodes are placed in fixed positions, the existence of isolated nodes is undesirable. In terms of communication networks, the probability that each node in a multi-hop OWC network has a minimum

node degree $d_{min} \geq k$ depends on the node density, $\rho$, as well as the transmission range, $R$ and is given by [17]

$$P(d_{min} > k) = \left(1 - \sum_{i=0}^{k-1} \frac{(\pi\rho R^2)^i}{i!} e^{-\pi\rho R^2}\right)^n \tag{1}$$

A path is defined as a sequence of successive edges on $G$. The existence of a path between two nodes denotes that they are connected. $G$ is connected when a path exists between all pairs of nodes. Similarly, we say that $G$ is $k$-connected ($k \geq 1$) when $k$ mutually independent paths exist between all pairs of nodes [16]. The conditions (1) ensure that every node has at least one neighboring node within its range. However, the event $d_{min} > 0$ is not a sufficient condition for ensuring the connectivity of the network. It can be proven that, if $n \gg 1$, then

$$P(G \text{ is } k\text{-connected}) = P(d_{min} \geq k) \tag{2}$$

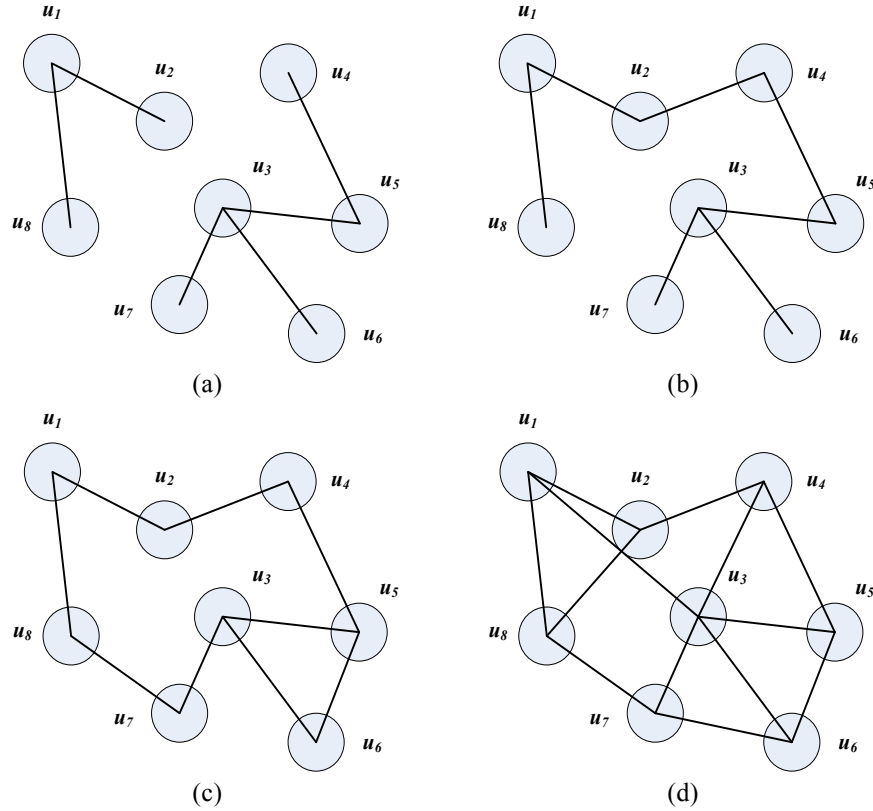for $P(d_{min} \geq k)$ almost 1 [16].



**Fig. 1.** Graph connectivity: (a) non connected, (b) 1-connected, (c) 2-connected and (c) 3-connected.

## 2.2 k-connectivity issues for ultraviolet UV-C two-dimensional multi-hop networks

In [12] we consider a multi-hop network configuration consisting of several NLOS UV-C communication sensors. A typical UV-C link between a transmitter (Tx) and a receiver (Rx) is shown in Fig. 2. Both the Tx and Rx face vertically upwards; i.e., they have 90° apex angle. In this scenario, the Tx transmits a signal vertically upwards having a beam divergence angle $\theta_T$. The cone produced by the Tx beam intersects the Rx FOV cone of $\theta_R$ degrees. The separation between Tx and Rx is $r$, while the distances from the common volume V to the Tx and Rx are $r_1$ and $r_2$, respectively. A communication link is established when the optical power is backscattered by particles inside the volume produced by the intersection of the two cones and reaches the Rx node.

Next, we assume that a number of $n$ nodes are distributed at fixed positions on a service area $A$. Each node is independently placed on the service area according to a homogeneous Poisson point process. Assuming large values of $n$ and $A$, a constant node density, $\rho=n/A$, can be obtained. Under this assumption, the homogeneous Poisson point process can be obtained as the limiting case of the uniform distribution. Consider, now, the case where all the nodes have the same transmission range $R$, i.e., homogeneous range assignment. This means that every node covers a circular region with area $A'=\pi R_0^2$. Every source node forwards traffic towards one or more destination nodes provided that their cones are intersected. If this does not happen, the node is getting isolated. If the transmission range is short, the probability of having isolated nodes increases. In contrast, assuming a large range, interference problems may appear. Apparently, if a node can communicate with more than one neighbor, the network robustness significantly increases; hence a proper selection of the transmission range is a critical parameter for the network connectivity robustness.
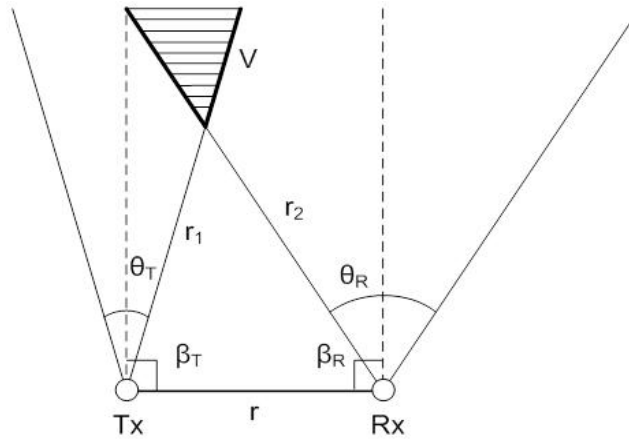


**Fig. 2.** UV-C NLOS link geometry.

The determination of an appropriate path loss model is critical for the k-connectivity investigation. Chen *et al.* in [20] proposed an empirical channel path loss model based on a set of extensive measurements. In their study, the authors demonstrated a communication test-bed and collected path loss measurements, for various combinations of Tx and Rx apex angles. On the basis of these measurements, a simple power decay model was proposed in [21] and adopted in [12].

As a result, analytical expressions for k-connectivity have been obtained, assuming the most fundamental modulation formats, i.e., on–off keying (OOK) [27] and pulse position modulation (PPM) [27] and adopting noise modeling for daily (Gauss model) and nightly (Poisson model) operation. For example, eq. (3):

$$P_{OOK,P}(d_{min} > k) = \left(1 - \sum_{i=0}^{k-1} \frac{\exp\left(-\rho\pi\left(\sqrt[a]{-\frac{\eta\lambda P_t}{hc\xi R_b \ln(2P_e)}}\right)^2\right)\left(\rho\pi\left(\sqrt[a]{-\frac{\eta\lambda P_t}{hc\xi R_b \ln(2P_e)}}\right)^2\right)^i}{i!}\right)^n \tag{3}$$

presents the network k-connectivity for OOK assuming the Poisson noise model. In (3), $P_t$ denotes the transmitted power, $R_b$ is the supported data rate, $P_e$ is the error probability, $\lambda$ is the operating wavelength, $h$ is the Planck's constant, $c$ is the speed of light, $\eta$ is the quantum efficiency of the optical filter and photodetector, $a$ is the path loss exponent and $\xi$ is the path loss factor.

The aim of this study is to address the following question: Given a homogeneous range assignment $R$, for a given modulation scheme, what is the minimum node density $\rho$ required to achieve a k-connected network with probability close to 1?

## 2.3 Node isolation probability for serial ultraviolet UV-C multi-hop networks

In [13] we focus on the node isolation probability evaluation of a serial multi-hop UV-C network and derive analytical expressions assuming the OOK and PPM modulation formats. The network consists of n transceivers (nodes), deployed at fixed positions on a service interval with length $l$, operating under NLOS conditions (Fig. 3). Every node is independently placed on the service interval according to a homogeneous one-dimensional PPP. Assuming large values of n and $l$, a constant node density, $\rho = n /l$, can be obtained. Under this assumption, the homogeneous PPP can be obtained as the limiting case of the uniform distribution [22]. Every node is equipped with a Tx and a Rx, with elevation angles of $\beta_T$ and $\beta_R$ degrees, respectively. The distance between a transceiver and its first neighbor is a random variable following a generalized Gamma distribution [23]. The Tx produces a cone, which has a beam divergence angle of $\theta_T$ degrees, and intersects the Rx FOV cone of $\theta_R$ degrees. A communication link is established when the optical power is backscattered by particles inside the common volume, generated by the intersection of the two cones, and reaches the Rx node.

We consider the case of homogeneous range assignment, i.e., all the nodes have the same transmission range $R$. Every source node forwards traffic toward its first

neighbor node provided that their cones are intersected. If this does not happen, the node becomes isolated. If the transmission range is short, the probability of having isolated nodes increases. On the contrary, assuming a large value of $R$, the interference level for each node may be significantly increased, thus degrading the quality of the communication link.
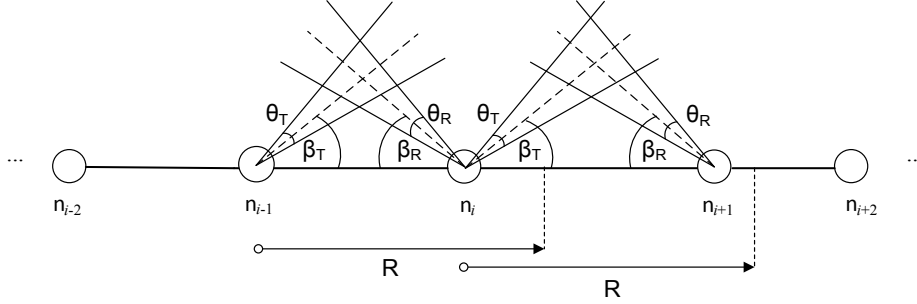


**Fig. 3.** Serial UV-C multi-hop network geometry.

A single scattering model is adopted, where each photon is assumed to be scattered at most once through its propagation from Tx to Rx. This model has been derived by Xu *et al.* in [19], as a fine approximation of the one introduced by Luettgen *et al.* in [18], which is presented in an integral form. Therefore, analytical expressions for the node isolation probability, $P_{iso}$, can be obtained in terms of the system parameters (i.e., transmitted power, supported data rate, and probability of error), as well as the geometrical configuration parameters (i.e., Tx and Rx elevation angles, Tx full beam divergence angle, and Rx FoV), and the node density. For example, eq. (4)

$$P_{iso,\text{OOK}} = \exp\left(-2\rho\,\frac{\sin\beta_S}{k_e\beta_1}\,W_0\left(\frac{\eta P_t k_s H(\mu) A_r \theta_T^2 \theta_R k_e \beta_1 (12\sin^2\beta_R + \theta_R^2\sin^2\beta_T)}{96\sin\beta_T\sin^2\beta_R\left(1-\cos\frac{\theta_T}{2}\right)\sqrt{N_0 R_b}Q^{-1}(P_e)}\right)\right) \qquad (4)$$

presents the node isolation probability for OOK. In (4), $H(\mu)$ denotes the composite phase function, $k_e$ and $k_s$ are the extinction and scattering coefficients, respectively, $\eta$ is the quantum efficiency of the optical filter and photodetector, $A_r$ is the area of the receiving aperture and $N_0$ is the white noise power spectral density. Obviously, an appropriate tradeoff between the node density, $\rho$, and the transmission range, $R$, is required to ensure a minimum number of isolated nodes. The presented results not only investigate this trade off but also considering the impact of the transceiver geometrical configurations (elevation angles, divergence angles) and the interaction of several parameters such as the supported data rate, the transmitted power and the network length.

## 2.4 Weather effects on node isolation probability for serial FSO multi-hop networks

In [14] we consider a serial network architecture composed of $n$ relays, i.e., $n$ FSO transceiver nodes uniformly distributed in a service interval of length $l$ according to a binomial point process (BPP) model (Fig. 4). The distance between a node and its $k$-th neighbor follows a generalized beta distribution given by eq. (1) in [14].
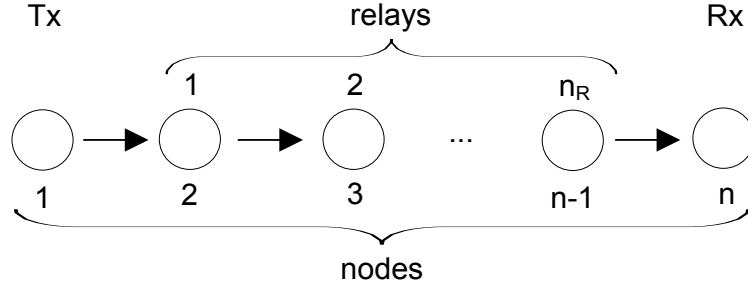


**Fig. 4.** Serial FSO multi-hop network geometry.

In recent years, significant effort has been devoted to the development of a channel model to predict weather effects on FSO transmission [24]. A quite effective model for the link budget evaluation is described in [25]. According to this model, the received power is related to the atmospheric attenuation which in turn depends on fog, haze, rain, or snow appearance.

Consequently, the node isolation probability is extracted. For example, eq. (5)

$$P_{iso,fog} = \left( 1 - \frac{2W_0\left(\frac{a_{fog}}{2\Theta}\sqrt{\frac{P_t A_r}{P_r}}\right)}{\ell a_{fog}} \right)^n \tag{5}$$

gives the node isolation probability when fog is present. In (5) $\Theta$ denotes the beam divergence and $\alpha_{fog}$ the attenuation coefficient due to fog. From this equation, we can conclude that an adequate node isolation probability for a given length, $l$, depends on various network parameters and the number of nodes, $n$.

## 2.5 k-connectivity issues for underwater optical wireless multi-hop networks

In [15] we consider a three-dimensional (3D) uOWC network consisting of nodes floating at different depths over a service aquatic volume. Each node is equipped with six sensors with FoV of 60º in order to cover all directions and allow transmission in three dimensions. A possible way to deploy such a network was proposed in [26]. Every sensor is placed at the bottom of the sea by an anchor and is connected with a buoy that can be inflated by a pump. Sensors are pushed towards the surface by the buoy. The depth adjustment can be controlled by arranging the wire length connecting the sensor to the anchor. By properly adjusting the depth, we can also construct a one-

dimensional uOWC network where nodes are arranged over a service length into the aquatic medium. The one dimensional arrangement can be considered as a special case of the general three dimensional node configuration.

The determination of the minimum achievable transmission range of each node is critical for the design and the connectivity investigation of the network. It is found, that this parameter is a function of network parameters (transmitted power, supported data rate, error probability), the operating wavelength and the chlorophyll concentration as shown in (6):

$$R = \frac{2\cos(\theta)}{c(\lambda)} \times W_0 \left[ \frac{c(\lambda)}{\cos(\theta)} \frac{1}{\sqrt{\frac{2\pi(1-\cos(\theta_0))TR_b hc}{\eta\lambda P_t \eta_t \eta_r A_r \cos(\theta)}\left(\sqrt{r_{dc}+r_{bg}}+\sqrt{\frac{2}{T}}\mathrm{erfc}^{-1}(2P_e)\right)^2 - r_{dc}+r_{bg}}} \right] \quad (6)$$

In (6), $c(\lambda)$ denotes the extinction coefficient, $\theta_0$ is the Tx beam divergence angle, $\theta$ is the angle between the perpendicular to the Rx plane and the Tx-Rx trajectory, $T$ is pulse duration, $\eta$ is the detector counting efficiency while $r_{dc}$ and $r_{bg}$ are the sources of additive noise due to dark counts and background illumination, respectively.

A set of analytical results is presented in order to study the connectivity behavior of an underwater multi-hop sensor network in the 300nm–700nm spectral region for two case studies. At first, a one–dimensional scenario is considered, where $n = 100$ nodes are uniformly distributed at a given service length. Secondly, a three–dimensional scenario is considered, where $n = 100$ nodes are uniformly distributed at a given service volume.


## 3    Conclusions

In the context of this dissertation we examined connectivity issues for three types of optical wireless communication networks.

At first, in order to extend the coverage region for a NLOS UV-C optical communication network, the network operation via multiple node-to-node hops is proposed and the system performance in terms of the k-connectivity property is investigated. The effect of several network parameters, such as the node density and the probability of error, were examined for uncoded OOK and PPM modulation schemes, assuming Poisson and Gaussian noise models. Various numerical results were illustrated, showing a useful outcome for telecom system designers for constructing a reliable UV-C network. Apparently, better results could be obtained by using more efficient modulation schemes, e.g., subcarrier intensity modulation, as well as incorporating coding, e.g., repetition or convolutional codes. Moreover, the use of different deployment geometries, e.g., use of directional beams, can obviously enhance the performance; however, the connectivity analysis in this case is getting much more complicated and is a subject of ongoing research.

Secondly, analytical expressions for the node isolation probability of a serial NLOS UV-C network were presented where transceivers are distributed statistically

on a given service interval. An effective path loss model is used and transmission with both OOK and PPM modulation schemes was considered assuming Gaussian noise. Several illustrative examples were depicted to show the interaction between various parameters, including the node density, the data rate, the required amount of power to achieve a certain error probability floor, etc. Different geometrical transceiver configurations were examined in order to obtain the node density required to achieve $P_{iso} \approx 0$ as well. The adoption of other path loss models and the consideration of more effective modulation and/or coding schemes are some of the topics for further research.

Thirdly, we focused on the node isolation probability of a serial FSO network where transceivers are placed on a given path-link following an one-dimensional BPP. We used an effective path loss model and considered operation under the most critical weather phenomena, e.g., fog, haze, rain, and snow. Proper design scenarios were presented in order to reveal the interaction between the number of required nodes, the length of the service interval, and the weather condition parameters (visibility, rainfall/snowfall rate) so that to achieve $P_{iso} \approx 0$, as well. In the worst case scenario, i.e., having thick fog with a 50m visibility, a network operator needs 15 nodes to cover a service length of 1km [13]. The work can be improved in a number of ways, e.g., using different path loss models, considering other modulation formats, using forward error correction schemes, etc.

Finally, we considered a underwater optical wireless communication (uOWC) configuration consisting of uniformly distributed nodes communicating with each other using IM/DD with OOK modulation scheme. We, then, investigated the interaction between the node density and various parameters such as error probability, wavelength, transmitted power, data rate, etc, in order to achieve connectivity. As an example, assuming the one dimensional scenario with a service length of 1km and $P_e = 10^{-6}$ we find from that the required number of nodes in order to deploy a fully connected serial multi-hop uOWC network is 336 for $R_b$=10Mbps and 184 for $R_b$=100kbps [14]. Obviously, these values would be significantly reduced by adopting more efficient modulation and/or coding schemes. This work can be expanded in several ways. At first, the adopted path loss model is only an approximation of the received intensity at each node and does not capture the scattering component of the transmitted beam. Several channel models presented in the literature, e.g., ([28] or [29]) could be assumed and more realistic results would be obtained. Furthermore, alternative models to describe the absorption and scattering coefficients could be adopted. As an example, the distinction between "small" and "large" particles according to Haltrin's model is determined by the index of refraction. However, this distinction could also be determined relative to the wavelength of illuminating resolution [30]. Moreover, the role of channel sharing/MAC/collisions should be clarified in case the non-isolated nodes do not operate successfully at a quoted bit rate.

# 4    References

1. Z. Xu and B. M. Sadler, "Ultraviolet communications: Potential and state-of-the-art," IEEE Commun. Mag., vol. 46, no. 5, pp. 67–73, May 2008.

2. S. Arnon, J. Barry, G. Karagiannidis, R. Schober, and M. Uysal, Eds., Advanced Optical Wireless Communication Systems. Cambridge University, 2012.

3. B. Cochenour and L Mullen, "Free-space optical communications underwater," in S. Arnon, J.R. Barry, G.K. Karagiannidis, R. Schober, and M. Uysal, "Advanced Optical Wireless Communication Systems", Cambridge University Press, 2012, pp. 273–302.

4. M. Uysal, Ed., Cooperative Communications for Improved Wireless Network Transmission: Framework for Virtual Antenna Array Applications. Information Science Reference, 2010.

5. M. Dohler and Y. Li, Cooperative Communications: Hardware, Channel and PHY. Wiley, 2010.

6. D. Miorandi and E. Altman, "Coverage and connectivity of ad hoc networks in presence of channel randomness," in IEEE Proc. of INFOCOM, Miami, FL, USA, Mar. 13–17, 2005, vol. 1, pp. 491–502.

7. L. Zhang, B.-H. Soong, Y. Zhang, M. Ma, and Y. Guan, "An analysis of k-connectivity in shadowing and Nakagami fading wireless multi-hop networks," in Proc. IEEE VTC Spring, Singapore, May 11–14, 2008, pp. 395–399.

8. A. Ghasemi and S. Nader-Esfahani, "Exact probability of connectivity in one-dimensional ad hoc wireless networks," IEEE Commun. Lett., vol. 10, no. 4, pp. 251–253, Apr. 2006.

9. A. Behnad and S. Nader-Esfahani, "Probability of node to base station connectivity in one-dimensional ad hoc networks," IEEE Commun. Lett., vol. 14, no. 7, pp. 650–652, July 2010.

10. M. Desai and D. Manjunath, "On the connectivity in finite ad hoc networks," IEEE Commun. Lett., vol. 6, no. 10, pp. 437–439, Oct. 2002.

11. D. Miorandi and E. Altman, "Connectivity in one-dimensional ad hoc networks: A queueing theoretical approach," Wireless Networks, vol. 12, no. 5, pp. 573–587, Oct. 2006.

12. A. Vavoulas, H.G. Sandalidis, D. Varoutas "Connectivity issues for ultraviolet UV-C networks", IEEE/OSA J. Opt. Commun. Netw., vol. 3, No. 3, pp. 199–205, March 2011.

13. A. Vavoulas, H.G. Sandalidis, D. Varoutas "Node isolation probability for serial ultraviolet UV-C multi-hop networks", IEEE/OSA J. Opt. Commun. Netw., vol. 3, No. 9, pp. 750–757, September 2011.

14. A. Vavoulas, H.G. Sandalidis, D. Varoutas "Weather effects on FSO network connectivity", IEEE/OSA J. Opt. Commun. Netw., vol. 4, no 10, pp. 734–740, October 2012.

15. A. Vavoulas, H.G. Sandalidis, D. Varoutas "Underwater optical wireless networks: a k-connectivity analysis", to appear in IEEE Journal of Oceanic Engineering.

16. M. D. Penrose, "On k-connectivity for a geometric random graph," Rand. Struct. Alg., vol. 15, no. 2, pp. 145–164, 1999.

17. C. Bettstetter, "On the minimum node degree and connectivity of a wireless multihop network," in Proc. ACM MobiHoc, Lausanne, Switzerland, June 9–11, 2002, pp. 80–91.

18. M. Luettgen, J. Shapiro, and D. Reilly, "Non-line-of-sight single-scatter propagation model," J. Opt. Soc. Am., vol. 8, no. 12, pp. 1964–1972, Dec. 1991.

19. Z. Xu, H. Ding, B. M. Sadler, and G. Chen, "Analytical performance study of solar blind non-line-of-sight ultraviolet short-range communication links," Opt. Lett., vol. 33, no. 16, pp. 1860–1862, Aug. 2008.

20. G. Chen, F. Abou-Galala, Z. Xu, and B. M. Sadler, "Experimental evaluation of LED-based solar blind NLOS communication links," Opt. Express, vol. 16, no. 19, pp. 15059–15068, Sept. 2008.

21. G. Chen, Z. Xu, H. Ding, and B. M. Sadler, "Path loss modeling and performance trade-off study for short-range non-line-of-sight ultraviolet communications," Opt. Express, vol. 17, no. 5, pp. 3929–3940, Mar. 2009.

22. C. Bettstetter and C. Hartmann, "Connectivity of wireless multihop networks in a shadow fading environment," Wireless Networks, vol. 11, no. 5, pp. 571–579, Nov. 2005.

23. M. Haenggi, "On distances in uniformly random networks," IEEE Trans. Inf. Theory, vol. 51, no. 10, pp. 3584–3586, Oct. 2005.

24. S. Bloom, E. Korevaar, J. Schuster, and H. Willebrand, "Understanding the performance of free-space optics [Invited]," J. Opt. Netw., vol. 2, no. 6, pp. 178–200, June 2003.

25. S. S. Muhammad, P. Kohldorfer, and E. Leitgeb, "Channel modeling for terrestrial free space optical links," in Proc. of ICTON, July 2005, pp. 407–410.

26. I.F. Akyildiz, D. Pompili, and T. Melodia, "Underwater acoustic sensor networks: research challenges," Ad Hoc Networks, vol. 3, no. 3, pp. 257—279, Mar. 2005.

27. Q. He, B. M. Sadler, and Z. Xu, "Modulation and coding tradeoffs for non-line-of-sight ultraviolet communications," Proc. SPIE, vol. 7464, pp. 74640H1–74640H12, 2009.

28. B. Cochenour and L Mullen, "Free-space optical communications underwater," in S. Arnon, J.R. Barry, G.K. Karagiannidis, R. Schober, and M. Uysal, "Advanced Optical Wireless Communication Systems", Cambridge University Press, 2012, pp. 273–302.

29. S. Jaruwatanadilok, "Underwater wireless optical communication channel modeling and performance evaluation using vector radiative transfer theory," IEEE J. Sel. Areas Commun., vol. 26, no. 9, pp.1620-1627, Dec. 2008.

30. C.D. Mobley, Light and Water: Radiative Transfer in Natural Waters, Academic Press, 1994.