# DSIE'10

## 5th edition

# Proceedings of the 5th Doctoral Symposium in Informatics Engineering

28-29 JANUARY 2010, PORTO, PORTUGAL

## 2010

EDITORS:

A. Augusto Sousa and Eugénio Oliveira

## DSIE'10 SECRETARIAT:

# FOREWORD

2010 Doctoral Symposium in Informatics Engineering - DSIE'10, already is the 5th edition of a scientific meeting organized by PhD students of the FEUP Doctoral Program in Informatics Engineering (ProDEI). These meetings have been held since the school year 2005/06 and their main objective is to provide a forum for discussion and demonstration of practical application of scientific research issues, particularly in the context of information technology, computer engineering and computer science. DSIE symposium comes out as a natural conclusion of a ProDEI course called Methodologies for Scientific Research (MSR).

The aim of that specific course is to teach students the processes, methodologies and best practices related to scientific research, particularly in the mentioned areas, as well as to improve their capability to produce adequate scientific texts. With a mixed format based on multidisciplinary seminars and tutorials, the course culminates with the realization of DSIE meeting, seen as a kind of laboratory test of the concepts learned by students. In the scope of DSIE, students play various roles, such as authors, both scientific and organization committee members, reviewers, duly accompanied by more senior lectures and professors.

DSIE is then seen as a "leitmotif" for the students to write scientific correct and adequate papers following the methods and good practices currently associated to outstanding research activities in the area. This year DSIE only admitted submissions from ProDEI students and, in many cases, papers already point to interesting research themes PhD students are willing to pursue. Although, still at an embryonic stage, and despite some of the papers are not yet enough mature or only report a state of the art, we already can find some interesting research work or an interesting perspective about future work. At this time, it was not essential, nor even possible, for the students in their PhD first year, to produce strong and deep research results. However, we hope that the basic requirements for presenting an acceptable scientific paper have been fulfilled.

DSIE'10 Proceedings include 23 articles accepted in the context previously defined. They have been put together according to the seven technical symposium sessions. These sessions group some different, although adjacent topics, once, as it was expected, the paper themes are very much heterogeneous. The sessions go from more theoretical to more applied

topics, from more focused to broader areas, but always including the algorithmic and scientific research methods flavor. The sessions are "Programming Fundamentals (3 papers), Machine Learning and Data Mining (4 papers), Multi-Agent Systems and Robotics (4 papers), Protocols and services (4 papers), Image Analysis (3 papers), E-learning Environments (2 papers), Information Technologies (3 papers).

The complete DSIE'10 meeting has a two days program that includes also two invited talks by outstanding Portuguese researchers that are also involved in non-academic important research activities related to engineering informatics.

MSR and ProDEI responsible professors, are proud to participate in DSIE'10 event and would like to acknowledge all the students who were deeply involved in the success of this event that, we hope, will contribute for a better understanding of the themes that have been addressed during the course, the best scientific research methods and the good practices for writing scientific papers.

***Eugénio Oliveira and A. Augusto Sousa***
(in charge of the MSR course – Methodologies for Scientific Research) January 2010

# PREFACIO

2010 Doctoral Symposium in Informatics Engineering- DSIE'10 representa a 5ª edição da série de encontros promovidos no contexto do Programa Doutoral em Engenharia Informática (Pro-DEI) da FEUP, desde o ano lectivo 2005/06. O seu principal objectivo é o de ser um fórum de discussão e de aplicação de práticas de investigação científica, nomeadamente no âmbito da informática, da engenharia informática e da Ciência da Computação. Organiza-se este simpósio, pelos estudantes, como conclusão natural de uma unidade curricular designada por Metodologias de Investigação Científica (MIC).

Esta unidade curricular, pertencente ao primeiro semestre da componente curricular do curso, pretende transmitir aos estudantes os processos, as metodologias e as boas práticas associados à investigação científica, sobretudo no âmbito das áreas referidas, assim como melhorar a sua capacidade de produção adequada de textos científicos. Com um formato misto baseado em tutoriais e seminários multidisciplinares, a unidade curricular culmina com a realização deste encontro que se destina a funcionar, figurativamente, como um laboratório de prova dos conceitos apreendidos pelos estudantes: estes desempenham vários papéis, como o de autores dos artigos e de membros das comissões de organização e científica, sempre devidamente acompanhados pelos docentes da unidade curricular e de outros docentes que aceitam o desafio de colaborar na revisão e, em alguns casos, na produção dos artigos submetidos.

O DSIE'10 surge assim como o mote para a produção de artigos com formato científico correcto, onde os autores colocam em prática os conhecimentos adquiridos ao longo da unidade curricular, sendo que nesta edição se limitou a participação aos estudantes da edição corrente do curso. Embora não de forma generalizada, as contribuições apresentadas apontam já, em muitos casos, para os temas que os estudantes pretendem seguir na componente de investigação do programa doutoral. Numa fase embrionária da investigação, é natural que alguns trabalhos se apresentem ainda algo incipientes, ficando-se por uma pesquisa de estado-da-arte numa área, ou por uma descrição pouco profunda de um trabalho apenas perspectivado e ainda a realizar no futuro.

Não é fundamental, nem tal seria possível, na totalidade dos casos, que os trabalhos fossem muito aprofundados, apresentando-se, no entanto, com o mínimo de requisitos que são os normalmente exigidos a um artigo científico.

O presente volume inclui os 23 artigos publicados nesse contexto, reunidos, de acordo com os assuntos versados, em 7 sessões técnicas do simpósio. Estas sessões agrupam e classificam os temas, expectavelmente heterogéneos, face à diversidade das áreas cobertas pelo ProDEI, dos artigos em publicação. Estas sessões incluem

tópicos que vão da teoria às aplicações, de tópicos focados a áreas mais abrangentes, mas sempre atendendo aos algoritmos e métodos de investigação em estudo: Fundamentos de Programação (3 artigos), Classificação e Extracção de Dados e Informação (4 artigos), Sistemas Multi-Agente e Robótica (4 artigos), Protocolos e Serviços (4 artigos), Análise de Imagem (3 artigos), Plataformas e Ambientes para Ensino (2 artigos) e Tecnologias da Informação em geral (3 artigos).

Num programa completo que se estende por dois dias, o DSIE'10 inclui ainda duas sessões preenchidas por oradores convidados que apresentam as respectivas experiências que são, simultaneamente, científicas e empresariais, em domínios da engenharia informática.

Os docentes de MIC do ProDEI, agradecem o empenho de todos quantos participaram nesta realização que, esperam, tenha contribuído para uma melhor apreensão dos temas tratados ao longo da unidade curricular, nos domínios das metodologias de investigação científica e da boa escrita de documentos relacionados.


***Eugénio Oliveira e A. Augusto de Sousa,***
(Docentes de MIC - Metodologias de Investigação Científica, Programa Doutoral
em Engenharia Informática)

# CONFERENCE COMMITTEES

## STEERING COMMITTEE

A. Augusto Sousa
Eugénio Oliveira

## SCIENTIFIC COMMITTEE CO-CHAIRS

Catarina Santiago
Isabel Margarido
Sofia Torrão

## SENIOR SCIENTIFIC COMMITTEE

A. Augusto Sousa
Ana Paiva
Ana Paula Rocha
André Oliveira Restivo
António Castro
António Coelho
António Lucas Soares
Eduarda Mendes Rodrigues
Eugénio Oliveira
Francisco Restivo
Gabriel David
Henrique Lopes Cardoso
João Canas Ferreira
João Correia Lopes

João Pascoal Faria
Jorge Alves da Silva
Jorge Barbosa
José Luis Borges
José Magalhães Cruz
Luis Paulo Reis
Luis Sarmento
Miguel Pimenta Monteiro
Paulo Portugal
Paulo Costa
Rosaldo Rossetti
Rui Maranhão
Rui Rodrigues
Sérgio Nunes

## SCIENTIFIC COMMITTEE

Adriano Kaminski Sanches
Altino Manuel Silva Sampaio
Bertil Maria Pires Marques
Bruno André Almeida Loureiro
Carlos Adriano de Oliveira Gonçalves
João Ferreira de Carvalho Castro Nunes
Joaquim António de Oliveira Meireles
Jorge Manuel Canelhas Pinto Leite
Lúcio Sanchez Passos
Luís Alexandre Moreira Matias

Marcelo Roberto Petry
Md. Anishur Rahman
Paula Alexandra Carvalho de Sousa Rego
Paulo Rogério Soares Proença
Pedro Brandão Neto
Raul Ramos Pollán
Robson Costa
Sílvio Costa Sampaio
Tiago Coelho dos Santos

## ORGANIZING COMMITTEE CO-CHAIRS

Bertil Maria Pires Marques
Jorge Manuel Canelhas Pinto Leite
Luís Alexandre Moreira Matias

## ORGANIZING COMMITTEE

Adriano Kaminski Sanches
Altino Manuel Silva Sampaio
Bruno André Almeida Loureiro
Carlos Adriano de Oliveira Gonçalves
Catarina Santiago
Isabel Margarido
João Ferreira de Carvalho Castro Nunes
Joaquim António de Oliveira Meireles
Lúcio Sanchez Passos
Marcelo Roberto Petry

Md. Anishur Rahman
Paula Alexandra Carvalho de Sousa Rego
Paulo Rogério Soares Proença
Pedro Brandão Neto
Raul Ramos Pollán
Robson Costa
Silvio Costa Sampaio
Sofia Torrão
Tiago Coelho dos Santos

# SPONSORS

2010 Doctoral Symposium in Informatics Engineering (DSIE'10) is sponsored by:

# WELCOME MESSAGE

**Eugénio Oliveira**
**Doctoral Program in Informatics Engineering (ProDEI) , Program Director**

**Dear Participants, dear Guests,**

It is my pleasure to open this Doctoral Symposium on Informatics Engineering to welcome you all, and to have the opportunity to address just a few words at this opening session.

I, firstly, would like to thank the presence of the FEUP Board representative as well as invited speakers, professors and, what is most important, all the crucial participants, you the PhD students.

This is the 5th meeting of this kind organized by PhD ProDEI (Doctoral Program in Informatics Engineering) students at FEUP, in the scope of the "Scientific Research Methods" course.

Main aims of this event simultaneously include the enhancement of scientific expertise in several Informatics related subjects, the use of adequate formats for both conveying ideas and to write papers, as well as the capability to organize the many different aspects that are implied by a scientific meeting organization.

This Symposium includes 7 sessions concerning different and broad subjects like: Programming fundamentals, Robotics and Multi-Agent Systems, Protocols and Services, Data Mining, Image Analysis, Learning Environments and other Information Technologies.

We are not here expecting outstanding scientific results although we do expect the rise of a genuine interest in scientific topics to be investigated further. I urge all of you to contribute to a fruitful discussion here at this symposium for the Informatics Engineering area, which encompasses both Computer Science and Computers Engineering.

I would like to finish by strongly acknowledge all the ProDEI PhD students and, if you allow me, to thank in particular those who took in charge the Chairs of both the Scientific and Organizing Committees.

Thank you and enjoy the opportunity.

# CONTENTS

## TECHNICAL PROGRAMME

**Invited Speaker – Prof. Henrique Madeira**
How to Measure Security and Trustworthiness

**Invited Speaker – Prof. Miguel Sales**

"Why Portugal is a top location for R&D software labs – the Microsoft experience

**Session 6 – Information Technologies**

**Session 7 – Learning Environments**

# Programming Fundamentals and Maturity Models

# Self-Protection Techniques in Malware

Tiago Santos

Faculdade de Engenharia da Universidade do Porto
R. Dr. Roberto Frias, 4200-465 Porto
pro09027@fe.up.pt

**Abstract.** This paper presents a survey over many of the methods of active and passive protection employed by malware software. All the information presented was obtained from official sources and from the author's personal experience. In the present paper it is showed the simplicity of several methods and how old techniques are continuously and successfully applied. Due to the ever-evolving of the current techniques and to the constant creation of new methods, it is impossible to assure a complete protection even with an up to date state of the art antivirus software.

**Keywords:** Malware, stealth, virus, worms, encryption, armoring, obfuscation, evasion

## 1 Introduction

Virus and malware in general, presents an extraordinary challenge not only for those who desire to study and create them but also for those who desire to fight them. Like in there's biological counterparts, these wonders of computation evolve, thanks to the improvements in software and hardware, making their ability to infect, spread and evade an ever growing problem. We can see an example of this evolution in the *Bagle* worm series, where it is possible to observe constants improvement along the entire series [1].

There are many publications about the subject, but they present or a very specific topic like in [2] or a wide range of issues but in great depth. Peter Ször with [3] is perhaps one of the bests works around, containing a very global and detailed information, but unfortunately many points are to lengthy or to short, and considering its publication year, the contents are many times outdated (which could be useful to understand the historical framework of the subject).

This paper was born from the necessity of an updated reference on the subject, and the goal it is to give a clear and elucidate brief overview over the main points on stealth, evasion and obfuscation used by the malware writers in general.

This paper will start with a quick and summarized presentation of the principal types of malware in Section 1, followed by the major techniques and methods em-

ployed by malware writers, in Section 3 and finalizing with some conclusions and remarks.

## 2   Malware Overview

While we can classify the malware software into categories, some of them or have a thin line between classifications or are a combination of several. For example, in terms of payload the *Trojan.Peed.Gen* and *Trojan.Downloader* are not only virus but also worms and trojans. So a correct classification is sometime difficult.

*Virus:* A computer virus, in its infection form is not a program by itself, but a piece of code that attaches itself to a file or resource and replicates itself, spreading to other resources. The offspring will then repeat the process. One should not confuse virus with *germs*, *droppers* and *injectors*. The first ones, are the first generation of a virus, i.e., the primary ancestor. Dropper are considered as the *installers* for the first generation viruses. Injectors belongs to the dropper family, but is specialized to install virus code in memory.

*Worm:* The distinction between worms and virus arises basically from their spread method. While the virus needs a file to attach itself in order to spread, the worm it is a self contained program or autonomous agent, that self propagates through a network thanks to flaws in security policies, services and programs [4] (when dealing with contagion worms [5], a thin line between worm and virus exists). The classification of the Worm must be done in terms of how they operate, their spread methods, target discovery and selection, carrier methods, activation, their payload and who uses them.

*Logic bombs:* A logic bomb is a piece of code, installed intentionally or not, in a legitimate application. Denominating a logic bomb a malware is sometimes incorrect, since they might be used by legitimate programmers with legitimate purposes like copy protection. This problem may occur in close code software and in large projects, where revision is difficult. One example of a malicious logic bomb can be seen in the game Mosquitos on Nokia Series 60 phones [3]. Another example, but a non one malicious, can be seen in some Microsoft products, where some programmers have inserted hidden function that are activated with some special keys combinations, showing pages giving credits to team members of the project. Spite the fact that they are not malicious, one should ask oneself: what other occult functions and process might exist?

*Trojans:* According to [6] a Trojan is a computer program disguised as a useful application but containing some additional or hidden functions that exploit the legitimate authorizations of the invoking process to the detriment of security.

*Rootkits* Rootkits are not malware, but it is relatively easy to use them for that purpose. They can be used to give the attacker access and root control of a system.

Between their capabilities, the most dangerous is the possibility of changing the kernel's behavior.

## 3 Stealth, Evasion, Obfuscation

In this section the main self-protection protection's techniques used by malware are presented in a summarized manner.

One point must be made clear: no self protection method is perfect or enough. But the same can be said about all the detection and cleaning methods.

Self-defense techniques can be classified in 2 classes: passive and active. As the name suggests the first one implies methods that protect the malware by creating a sort of a shield (polymorphism, encryption, compression, etc.). Active protection implies a more aggressive form of defense, like attacking or disrupting the thread (rootkits, anti-debugging, retro-virus, etc.).

Since all extra routines and functions have a cost in efficiency, the malware writers must always balance the gains and loses in applying each method.

### 3.1 Code Obfuscation

[2] states "obfuscated codes are generally too slow and of too large a size to be efficiently used by undetectable malware", but the author loses the main point of this technique. This technique is intended not only to escape detection but mostly to enhance the difficulty in the malware's code analysis.

Take the following code:

```
LEA EAX, DWORD PTR [040200H]
JMP EAX
```

Its complexity can be increased if optimization is ignored. By just inserting unneces-sary steps, or more easily, by inserting garbage and changing a few registers:

| *Without junk code* | *With junk code* |
|---|---|

```
Without junk code          With junk code
MOV   AX, 0200H            PUSH  EBX
MOV   BX, 0004H            XOR   EAX, EAX
AND   EAX, 0XFFFF          MOV   EAX, 040200H
SHL   EBX, 16             PUSH  EAX
OR    EAX, EBX            ADD   EAX, EBX
JMP   EAX                 POP   EBX
                          POP   EAX
                          JMP   EBX
```

The problem with the above method is that the code tends to increase significantly, but there are methods like using checksums, opcodes confusion and undocumented instructions, that allows obfuscating the code without substantially increasing the code's size.

## 3.2   Entry-Point Obfuscation

With this method, the virus hides its location, causing many scanning tools to fail [7]. This technique is effective, since unlike traditional virus, the Entry-Point Obfuscation (EPO) doesn't change the application's entry point, but instead it changes any jump or call in the infected file. Therefore, there are potentially an unlimited potential points where the virus can start running and since the antivirus (AV) can't check them all, misses are frequent.

## 3.3   Oligomorphism

Oligomorphic virus are capable of changing their decryptors in every new generation [8]. It does this by selecting in a random manner each piece of the decryptor from several predefined.

Its weakness resides on the small number of possible different decryptors and on the fact that their signatures are usually too common [3]. Interesting fact, is that, AV software tends to detect this type of virus by trying to dynamically decrypt the body [8].

## 3.4   Polymorphism

Polymorphism gives the ability to the malware to create a complete different copy of itself in each infection.

To keep the size of the malware small, it is impractical to make its body fully polymorphic [9], due its high number of instructions. The approach of applying polymorphism to the entire code is not easy to implement and for this reason the malware writers opt to only use polymorphism in the malware's decryptor section and encrypt the main body (*Win32/Coke*, *Win95/Marburg* and *Win95/HPS*). The routine responsible for the polymorphism is called "*mutation engine*".
    [10] states:

>     "With no fixed signature to scan for, and no fixed decryption routine,
> no two infections look alike. The result is a formidable adversary."

Today's AV software uses virtual machine methods to catch these types of malware, but are useless if the malware possesses anti-emulation systems (see Sec. 3.12). Other method used by the AVs is the *Heuristic-Based Generic Decryption*,

which is faster that the use of virtual machines. This method uses a generic set of rules to differentia-te malware and non-malware behavior. Once again, if the malware has anti-heuristics methods (see Section 3.12), the AV will be fooled again.

### 3.5 Metamorphism

This one is a more effective way to avoid detection than polymorphism, since it can many times by itself fool the emulation of the AV software.

It differs from polymorphism by not using encryption, which means that the entire body is altered, i.e., it gives the malware the ability to rewrite itself completely, creating an entire new copy of itself in each infection, without affecting its functionality.

Metamorphism if achieved using a metamorphic engine, which usually uses techniques like reordering instructions, inserting NOPs, changing blocks of instructions or the program flow. Since all the code mutates, the metamorphic engine is able to create multi-platform malware.

A metamorphic malware uses code obfuscation techniques to fight against static analysis but also use behavior modification to fight dynamic analysis.

Normally a metamorphic malware does not use encryption, unless it is for obfuscation purposes. Thanks to this, each mutation is different, performing that way evasion from the AV. The detection and evasion methods that were mentioned to the polymorphic case, can also be applied here.

Due the complexity of the metamorphic engine, false infections and bugs are a norm.

### 3.6 Compression and Packing

The use of compressions gives the malware the ability of stealth, reduces significantly its size and makes the code analysis much more difficult and tedious. The malware might have compression functions but a more sophisticate process and a way to reduce the malware's size is to search and use any compression and encryption software, like pkzip and ASPack, in the infected system.

Using a packer for compression it is possible to change the entire packed code by just changing one single byte [11]. Another advantage of using packers, resides in the possibility that some of them, incorporates in the compressing file some anti-debugging capabilities.

### 3.7 Retrovirus

Retrovirus (RV) is a computer virus that actively attacks the AV software, trying to bypass or to block the AV's operations, personal firewall or other security programs [3].

Since the virus writers have access to AV software, it is not hard to develop routines able to attack them, in a way that the AV developer didn't expected [12]. Some techniques are:

– direct modification of the AV code (file or in memory) or its operation (e.g. modifying command-line parameters or the configuration file, modify the alerts by laun-ching false alerts, deleting the database, etc.)
– deleting the AV (this may alert the user)
– detecting the presence of a AV and hide itself or launch a destructive payload
– changing the infected system in a way that affects the AV
– the RVs might use code that causes problems to the AV
– exploitation of weakness in AV software (e.g. if the AV uses emulators, once the retrovirus is executed, it detects the emulator and escapes from the protected environment, extending its infection or launching its payload)
– use of methods that difficult or make dangerous to perform any action against the RVs

Thanks to RVs, modern AV software are required to have extra self-protection in order to prevent these type of attacks.

## 3.8 Sparse Infection

Sparse infection consists in an extreme control over the malware propagation and infection. For example, a virus may infect other resources only if certain conditions are met, like date, the size of a file or available memory, applications installed, etc. By doing this, the malware reduces the risk of detection and capture.

## 3.9 Tunneling

This stealth and resilient method is used mostly by memory-resident virus. The goal to bypass any AV monitoring programs by being the first of a series of calls, installing himself before any other resident program [3].

AV software stays in the background, searching for specific malware actions. It works by intercepting the Operating System (OS) before the virus is executed, but applying this technique, the malware is able to intercept directly the interrupt handlers of the OS or BIOS, making evasion possible.

## 3.10 Self-Disinfection

It is a particularity interesting method, where the malware eliminates itself after the delivery of his payload. This process allows the malware to avoid been analyzed and captured for a relatively long time.

### 3.11  Encryption

This section some of the major encryption algorithms used by malware writers and some code implementations are presented. The goals of using encryption are to hide the malware's fixed signature [10], making it unrecognizable to a malware scanner and to difficult the analysis by AV developers. An encrypted malware is made of 2 parts: the decryptor and the encrypted body. To infect, the malware's decryptor part gains control and decrypts the body, passing then the control to it [10], allowing the launch of the payload or the infection of other resources.

One important weakness in malware programs, is the use of constants or strings (e.g. ".exe", "windows/system32" or email address to where worms might send the information they collected). In this case the encryption application eliminates this serious Archiles' heel. Encryption not only hides this type of information but also makes the code extremely hard and tedious to analyze.

During the infection process, there's the possibility that a malware decrypts itself badly, causing a possible bad infection. In this case [7] the author states:

> "... virus writers don't have to achieve perfect infection. It isn't crucial that infection attempts sometimes fail, or if a virus can't reliably tell whether a file is already infected. Viruses are messy and an imperfect virus can spread quite well in the real world. Conversely, an antivirus product must be extremely reliable. Unreliable detection, either by identifying a benign file as infected or vice versa, is a fatal flaw."

**Typical Encryption Algorithms.**  *Simple encryption:* Simple encryption (also known as substitution encryption) is one of the most basic form of encryption. In this type, each character has one correspondence, given by a mathematical operation between one key and the text to encrypt. In order to decrypt, it is necessary to save the key. The problem with this type of encryption is that is only really effective against casual or lazy people. In computer sciences, this type of encryption, is usually done using *xor*.

*Sliding key encryption:* As the name suggests, in this type of encryption the key changes after the encryption of each byte or block. Since the key is smaller than the block to encrypt, the mathematical key is still extremely easy to discover.

*Long key encryption:* In this case, the key is longer than the block to codify. With a long key, the key discovery is slightly more complex.

*Transposition encryption:* This method is basically the reordering or the scramble of the target code. The major problem is the quantity of code and data necessary to perform the scrambling and the reconstruction of the code.

*Variable length transposition:* This is a variation of the transposition encryption, and works relatively well when used more than once and combined with other method. Like the previous method, the objective is to scramble bytes. In this case,

this process is done in chunks of bytes, instead of applying to an entire block. So it can be first applied for groups of 4 bytes, then 2, then 8, etc. The important to remember is to perform the decryption in the reverse order.

*Boundary Scrambling:* Simple and effective. The goal is to destroy the bytes' bounda-ries, by scrambling them with the boundaries of the predecessor's or the antecedent's bytes.

*RDA:* RDA stands for Random Decryption Algorithm. First implemented in the *RDS.Figther* virus. The strength of RDA, resides in not knowing the key, i.e., after encryption, the only way to decrypt it, is by brute force. The key can be anything from local DNS addresses to current system time and date.

*Integrity-Dependent Encryption/Decryption:* Very simple and useful encryption / decryption method. This method allows to verify the integrity of the data, allowing the deployment of the payload in case of tempering and corruption. The method uses the value of the previous decrypted byte to decrypted the next. It uses the checksum, to allow to check if the decryption is correct or not. During the encryption a checksum is created, using it as the key. This means that the key is always changing during the encryption. Initially the checksum is zero and the value that is going to be encrypted is added before its encrypted. The key must be written before the encrypted code as the key value.

*RCX encryption:* The RC encryption algorithm family raises the bar in terms of complexity relatively to the previous methods. These algorithms are extremely simple to implement and the time to decipher them without the correct key is pro-portional to the key's size. These algorithms are excellent when the key is delivered by the virus writer through the network. After the end of encryption, the key is lost and to decrypt, the key must be obtained again from the network.

*Public key encryption:* In a public key or asymmetrical key encryption, there are two keys, a private and a public. These are chosen in a way that the two cannot be deduced by each other [13]. To encrypt, the encryptor uses the public key. The encrypted data cannot be decrypt using the public key, only with the private key. The use of this algorithm to encrypt of parts of a malware is not recommended and it is preferable to use for example the RC4 that is usually faster than for example the RSA algorithm [13]. But this type of encryption has its utility in extortion schemes.

Following are several simple implementations of encryption algorithms, but please keep in mind that for clarity, the code is clean and not optimized:

*Simple encryption*

```
MOV  ESI,STARTPOINT
MOV  EDI,DESTINATION
MOV  ECX,CODESIZE/4
@LOOP:
  LODSD
  XOR  EAX,05060708H
  STOSD
  LOOP @LOOP
```

*Sliding key*

```
MOV  ESI,STARTPOINT
MOV  EDI,DESTINATION
MOV  ECX,CODESIZE/4
XOR  EDX,EDX
@LOOP:
  LODSD
  XOR   EAX,EDX
  ADD   EDX,05H
  STOSD
  LOOP  @LOOP
```

*Long Key encryption*

```
MOV  ESI,STARTPOINT
MOV  EDI,DESTINATION
MOV  ECX,CODESIZE
MOV  EDX,KEY
MOV  EBX,0
@LOOP:
  LODSB
  XOR   AL,BYTE PTR [EDX+EBX]
  STOSB
  INC   EBX
  CMP   EBX,KEYSIZE
  JB    @JUMP
  MOV   EBX,0
@JUMP:
  LOOP  @LOOP
```

*Transposition Encryption*

```
; invert all code
MOV  ESI,STARTPOINT+CODESIZE-1
MOV  EDI,DESTINATION
MOV  ECX,CODESIZE
STD
@LOOP:
  LODSB
  MOV   BYTE PTR [EDI], AL
  INC   EDI
  LOOP  @LOOP
```

*Boundary Encryption*

```
MOV  ESI,STARTPOINT
MOV  EDI,DESTINATION
MOV  ECX,CODESIZE
@LOOP:
  LODSW
  DEC   SI
  MOV   BX,AX
  AND   AX,0FF0H
  XCHG  AH,AL
  SHR   AH,4
  SHL   AL,4
  AND   BX,0XF00F
  OR    AX,BX
  STOSW
  DEC   DI
  LOOP  @LOOP
```

*Integrity-Dependent Encryption*

```
MOV  ESI,STARTPOINT
MOV  EDI,DESTINATION
MOV  ECX,CODESIZE/4
XOR  EDX,EDX   ; checksum=0
@LOOP:
  LODSD
  ADD   EDX,EAX
  XOR   EAX,EDX
  STOSD
  LOOP  @LOOP
MOV  EAX,EDX
STOSD  ; save checksum
```

**Encryption Strength.** When the encryption or the decryption sections are visible, the encryption algorithm is weak, making the decryption of the malware relatively easy. Many of the methods described above are relatively simple to decrypt, but the important is to increase the times between the four temporal points: infection, detection, analysis and contra-measures. The strength increases with the application of multiple encryption algorithms. The major problem with all the encryption presented here is that all of them (except for algorithms like RDA and public key) have private keys. (One way malware writers use is making the malware to ask the writer for a key through the network. When that key is received encrypts or decrypts the code and then the key is lost. Obviously to increase the security in this communications channels, the keys and address must be continuously altered, e.g., 100 keys and 1 address for a group of 100 infections.

### 3.12 Armoring

In [2] the author demonstrated that an effective armoring is possible. Many of processes described above can also be included in this section. Armoring belongs to the class of anti reversing methods, including: anti-disassemble anti-emulation and others. This techniques are not only applied to malware but also to licenses, copy protections and digital rights management [14]. These methods are also applied in the applications' binaries, and they are able to increase dramatically their complexity. They are design specifically to counteract the technologies used by the reverse people.

*Anti-disassemble:* In order to prevent the analysis of the malware with a disassembler, the malware writers applies techniques to trick the disassembler by generating an incorrect disassembled code. The process by which this is done, is to add mix-up blocks of code so that the disassembler gets confused in the identification of the correct data byte without affecting the execution of the code.

Starting with the *Linear Sweep Disassemblers*, they assume that every byte in the code section belongs to the executable code, making extremely easy to trick them, by just adding data in the middle of the code and a local jump over the data so do not affect the code execution. The previous method does not work with *Recursive Traversal Disassemblers* (like IDA Pro and Ollydbg), since they operate based on the control flow of the application [14]. In this case an easy trick is the use of opaque predicates, which are false conditions statements.

Take the following code created by a programmer where the method just described is used:

```
MOV  BX, 0xF0F0
CMP  BX, 0
JE   @OPAQUE
JNE  @REALCODE
@OPAQUE:
  DB 50H
@REALCODE:
  ADD  CX,BX
  XOR EAX, EAX
...
```

Analyzing the previous code in a disassembler[1]:

```
401000  MOV  BX,0F0F0
401004  CMP  BX,0
401008  JE   SHORT 0040100C
40100A  JNZ  SHORT 0040100D
40100C  PUSH EAX       ; incorrect
40100D  XOR  CX,BX     ; incorrect
401010  ADD  EAX,EAX   ; incorrect
...
```

The advantage of using this last technique is that it works both in the *Recursive Traversal Disassemblers* and in the *Linear Sweep Disassemblers*.

*Anti-Heuristics:* First of all, one needs to explain how heuristics in AV works. There are 2 types of heuristics: static and dynamic. The static relies in the analysis of the file format and common code fragments [3]. The dynamic, on the other hand, uses code emulation to impersonate the processor and the OS, detecting the malware actions in a controllable environment. There are countless ways to trick many of the heuristics methods employed by AV, so it will be explained just a few of them.

Taking for example the infection if PE files. Currently a virus might use an old method of infection, where he appends sections of code in the end of PE files. Since this is a well know trick, the AV can detect them without difficulty, in view of that all they have to do is to check if the entry point of the PE file points to the last section of the file. Obviously, this might signal some false positives. Another heuristic test is performed in this cases, a compression check, seeing that the malware writers uses many times legal packers like ASPack to perform this compression.

New methods of code insertion were created to evade heuristic methods, like the insertion in multiple sections, very common with win32 viruses. This process is effective because the entry point will no longer point to the last section, but to other points, in sections like .data and .text, bypassing the heuristic methods from the previous paragraph.

---

[1] Ollydbg v1.10

*Anti-Emulation and Virtual Machine Detection:* Ever since the virus writers realized that some AV used emulation for virus detection, they immediately started to create methods to deceive them. There are many ways for the malware to determine if it is running inside an emulator, from checking systems tables to analyzing the set of ins-tructions of the processor or co-processor like MMX and FPU instructions (since some emulators do not implement them).

*Anti-Debugging:* The anti-debugging goal is to entail difficulty to the analyzer in using a debugger in the malware, by detecting if it is running inside one. Since not only the software permits the use of these techniques but also the hardware (e.g. using interrupt vectors), it creates an entire of rather exquisite set of ticks. Of course when malware writers use methods based on hardware, he creates a platform-specific dependency.

These techniques can be founded in virus like *Whale*, *Cryptor*, *W32/Cabanas* and *CIH*. It must be said that these methods are also used by legal software, especially by commercial executable protectors and packers. Among the current major anti-debugging techniques there is the *API Based Anti-Debugging* (e.g. *isDebuggerPresent* ) and the *Exception Based Anti-Debugging*. Protecting the code through API functions is almost useless, due the easiness in identifying and overcome them. Figure 1 shows an example, of how an anti-debugging API its showed by a disassembler[2]. Running the program inside the debugger, the application doesn't run normally, but by changing or adding a simple jump or other instruction in the code is sufficient to overcome this protection.



| 00402030 | .rdata | Import | user32.GetWindowTextA |
| 00402000 | .rdata | Import | comctl32.InitCommonControls |
| 00402010 | .rdata | Import | kernel32.IsDebuggerPresent |
| 0040202C | .rdata | Import | user32.LoadCursorA |
| 00402020 | .rdata | Import | user32.LoadIconA |
| 00402024 | .rdata | Import | user32.MessageBoxA |
| 00401000 | .text | Export | <ModuleEntryPoint> |
| 00402028 | .rdata | Import | user32.PostQuitMessage |

**Fig. 1:** isDebuggerPresent API in the import API area of the disassembler

Among the interrupts, we can find INT 03h or INT 2Dh. INT 03h is the most common BP (breakpoints) used to reverse engineer, making it in a priority target. INT 03h is an instruction represented by the opcode 0xCC that is automatically inserted by the disassembler to cause a breakpoint The disassembler does not shows the code altera-tion, and to see it, it is necessary to check it indirectly (see Fig. 2). The detection of this instruction is extremely easy, all that it is necessary is to verify the presence of the opcode 0xCC in the current instruction position instead of the original code.

INT 2D is a very interesting instruction, since can be used has a general purpose debugger detector. If the application is not been run inside a debugger, an exception

---

[1] Ollydbg v1.10

```
00401310  .  E8 2B050000        CALL 00401840
00401315  .  E8 A6000000        CALL 004013C0
0040131A  .  C70424 00304000    MOV DWORD PTR SS:[ESP],00403000
00401321  .  E8 7A050000        CALL <JMP.&msvcrt.printf>
00401326  .  B8 00000000        MOV EAX,0
```

**(a)** A normal breakpoint set at 0040131A

```
Address  | Hex dump                 | ASCII
0040131A  C7 04 24 00 30 40 00 E8   Ç□$.0@.è
00401322  7A 05 00 00 B8 00 00 00   z□..¸...
0040132A  00 C9 C3 90 90 90 55 89   .ÉÃ□□□U%
```

**(b)** Memory starting at 0040131A

**Fig. 2:** Alterations in the code caused by int 03h breakpoints are not showed by the disassembler

occurs, but if it is, the debugger will jump 1 byte, causing the following code to be incorrect. There are many other methods including stack segment manipulation (causes the debugger to execute instructions unwillingly) and detecting hardware breakpoints (implemented by Intel and are controlled by special registers).

*Antigoat:* Goat files are used by AV developers to better understand the infection, since the entire file's contents of the original file are known, making easy to separate the infection and the file. They typically contain NOPs instructions. The malware uses heuristics techniques to detect these files, e.g., by imposing some conditions like file size, if the file contains a large number of NOP instructions, etc.

## 4   Conclusions

This field is extensive and due the lack of space it is impossible to touch all the methods and to explain them in a more detailed manner. The number of techniques is overwhelming, and it was tried to balance to the maximum the coverage between passive and active protection methods. It is possible to observe how simple and efficient some methods are and how old methods like polymorphism and metamorphism are continuously and successfully applied to achieve evasion from AV software.

Can a system be safe with an up to date AV software? No. Since the offensive capabilities will always be ahead of the defensive capabilities, the malware writers will always be ahead from the AV developers. According to Bontchev, anything less than 100% detection is almost a complete failure, because it just takes 1 missed virus, in order to reinfect all the system again [9]. Due to the continue evolution of the techniques and tricks used by the malware writers and the development of new ones, the conclusion is straightforward: it is impossible to have a complete protection even with an up to date state of the art AV software, making prudent to always assume that a system is unprotected and that it is already infected.

# References

1. Honen, T.: A worm's evolution. http://www.virusbtn.com/conference/vb2004/abstracts/thonen.xml (2004)
2. Filiol, E.: Strong cryptography armoured computer viruses forbidding code analysis: the bradley virus. In: EICAR2005 annual conference 14, StJuliens/Valletta - Malta (2005)
3. Szor, P.: The Art of Computer Virus Research and Defense. Advances in Information Security , Vol.44. Addison-Wesley (2005)
4. Weaver, N., Paxson, V., Staniford, S., Cunningham, R.: A taxonomy of computer worms. In: WORM '03: Proceedings of the 2003 ACM workshop on Rapid malcode, New York, NY, USA, ACM (2003) 11–18
5. Staniford, S., Paxson, V., Weaver, N.: How to own the internet in your spare time. In: Proceedings of the 11th USENIX Security Symposium, Berkeley, CA, USA, USENIX Association (2002) 149–167
6. Gordon, S., Chess, D.: Where there's smoke, there's mirrors: The truth about trojan horses on the internet. http://www.research.ibm.com/antivirus/SciPapers/Smoke/smoke.html (1998)
7. Ford, R.: The wrong stuff? Volume 2., Los Alamitos, CA, USA, IEEE Computer Society (2004) 86–89
8. Szor, P., Ferrie, P.: Hunting for metamorphic. http://www.symantec.com/avcenter/reference/hunting.for.metamorphic.pdf (2003)
9. Watson, G.: Discussion of polymorphism. http://vx.netlux.org/ (1992)
10. Symantec Corporation: Understanding and managing polymorphic viruses. http://www.symantec.com/avcenter/reference/striker.pdf (November 1996)
11. Shevchenko, A.: The evolution of self-defense technologies in malware. http://www.hypponen.com/staff/hermanni/more/papers/retro.htm (2007)
12. Hypponen, M.: Retroviruses - how viruses fight back. http://www.hypponen.com/staff/hermanni/more/papers/retro.htm (1994)
13. Piper, F., Murphy, S.: Cryptography: A Very Short Introduction. Oxford University Press (1984)
14. Singh, A.: Identifying Malicious Code Through Reverse Engineering. Advances in Information Security , Vol.44. Springer (2009)

# Heuristics Study for the Traveling Salesman Problem

Joaquim Meireles[1]


[1]Faculdade de Engenharia da Universidade do Porto
Rua Dr. Roberto Frias, S/N 4200-465, Porto, Portugal
pro09011@fe.up.pt

**Abstract.** The Traveling Salesman Problem is very known problem with the aim to find the shortest route in a set of cities starting and ending in the same city. Running time of exact methods to solve the problem can be prohibitively so alternatives can be considered in situations where the optimal distance is not mandatory. There are some techniques like heuristics that can obtain an optimal or approximated solution to the problem. The application implemented will apply two known constructive and improvement heuristics, Nearest Neighbor and 2-Optimal, and a metaheuristic, Tabu-Search. It will be made a comparison of the obtained solutions by those algorithms defining their paper on the discovery for the global solution. It will also be demonstrated that metaheuristics are essential in order to achieve results closer to the optimal distance and can be considered for problems that don't need optimal solutions.

**Keywords:** Traveling Salesman Problem, Optimization, Heuristics, Metaheuristics.

## 1    Introduction

The Traveling Salesman Problem (TSP) is a known problem and is one of the most studied in the optimization scientific community. The problem is to find the shortest route of a traveling salesperson that starts at a home city, visits a prescribed set of other cities, each only once, and returns to the starting city[1]. For this study, and because there are many variations of the TSP, our work was based on the Symmetric TSP (STSP) without restrictions and, for simplicity purposes, it will be referred only by TSP.

Exact methods are used to solve TSP instances, but they are known to be time consuming, so they cannot be applied in more complex TSP problems. When instances become too large for exact methods, heuristics and in particular metaheuristics are often used as alternative. Heuristics algorithms have proved to obtain good solutions with limited computational effort and time. [2-3].

The main goal for this study is the development of a generic application to the TSP that shows the employment of known heuristics and their impact on the problem solution. Based on some known datasets supplied from TSPLIB[4], each one with their own characteristics, different heuristics (constructive, improvement and an extension of improvement heuristics: metaheuristics) will be applied and combined in order to find high quality solutions. The heuristics selected for this study are well

known and were chosen due to their simplicity and efficiency. It was chosen one heuristic of each type: Nearest Neighbor as a constructive heuristic and 2-Optimal as an improvement heuristic. A metaheuristic, Tabu-Search, was also applied over 2-Opt so better results could be achieved.

Despite those techniques may not retrieve optimal solutions, they can be implemented in TSP applications where optimal solutions are not essential to the problem resolution. Finally the quality of the solutions obtained will be analyzed and individual heuristics results will be compared. Final solutions are expected to be near the optimal value expecting results closer to the optimal after applying Tabu Search algorithm.

Section 2 gives a brief overview of state of the art, section 3 discusses the development of the application, focusing on the implementation of the constructive, heuristic and meta-heuristic algorithms. Section 4 will illustrate the results obtained from the application based on the TSPLIB sample data. Those results will be discussed focusing on the distance obtained for each heuristic, compare them to the library optimal value, and analyze the gain obtained with each one. Finally section 5 presents conclusions about this work showing the importance of metaheuristics to the resolution of the TSP and future work that could be done.

## 2     Background

In the TSP, a salesman needs to visit N cities with given distances $d_{ij}$ between two cities $i$ and $j$, returning to his city of origin. Each city must be visited only once, and the tour must be as short as possible. Mathematically speaking, given $n$ points and a cost matrix, *[c(i, j)]*, a tour is a permutation of the $n$ points. The points can be cities or arcs, and the permutation is the visit of each city exactly once, and then returning to the first city. The cost of a tour, *<i1, i2... in-1, in, i1>*, is the sum of its costs: *c(i1, i2) + c(i2, i3) + ... + c(in-1, in) + c(in, i1)*, where *(i1, i2, ..., in)* is a permutation of *{1,...,n}*[5-6]. Generally, the distance travelled in a tour depends on the order in which the cities are visited so the objective is to find the right sequence of the cities which will minimize the total cost.

When the cost matrix is symmetric, *d(i,j)=d(j,i)*, it is called a symmetric TSP, otherwise it's called asymmetric TSP (ATSP). Based on these fundaments, there are some TSP variations each of them with their specifications and objectives like the MAX TSP, Bottleneck TSP or TSP with Time Window (TSPTW)[1].Taking in consideration all the TSP variations, formulations of applications for this problem types are numerous, such as vehicle routing, computer wiring, cutting wallpaper and job sequencing[6-7].

In the theory of computational complexity, the decision version of TSP belongs to the class of NP-complete problems. Thus, it is assumed that there is no efficient algorithm for solving TSPs. It is likely that the worst case running time for any algorithm for TSP increases exponentially with the number of cities, so even some instances with only hundreds of cities will take many CPU years to solve it[1, 6]. There are two ways of finding the optimal solution: by exact or approximation algorithms. The purpose of exact algorithms is to find the optimal solution to a problem, so they will work reasonably fast for relatively small problem sizes. Finding

the exact solution to a TSP with *n* cities requires to check *"(n − 1)!"* possible tours and usually they "brute force search" must be done to check all possibilities, selecting the best one in the final.

Despite the fact that a significant progress has been made in the ability to solve TSP by exact algorithms, there are still several aspects of the problem to be solved to reach optimality, that are hard for exact algorithms. Moreover, even when an instance is solvable by an exact algorithm, the running time may become prohibitively large for certain applications. Therefore, a large number of heuristic algorithms start appearing , some of which normally produce solutions near optimal[1, 6].

TSP heuristics are generally classified in two types: tour construction heuristics and improvement heuristics, where there is also an extension of this last one: meta-heuristics. A construction heuristic is an algorithm that determines a tour according to some construction rules (incremental), but does not try to improve upon this tour. Some examples of constructive heuristics are Nearest Neighbor (NN) that is used in this work, Nearest Insertion, Furthest Insertion and Sweep heuristic.

Generally, constructive algorithms make initial solutions very quickly, but their quality is inferior by those produced by improvement algorithms. Improvement heuristics starts from an initial tour and search for a better one by iteratively moving from one solution to another, accordingly to adjacent relationships defined by a given neighborhood structure[1]. Some examples of these types of heuristics are the k-opt heuristic (2-Optimal and 3-Optimal) or the Lin-Kernighan heuristic.

The optimal solution that an improvement algorithm could find depends on the initial solution because the algorithm will only find the optimal solution if the search is made in the neighborhood of that optimal. In that context, an expansion of improvement heuristics appears to try to solve that problem: metaheuristics. Their main idea is to evade from a local optimal to a worst solution so that the algorithm could go into another space of solutions that might find the global optimal. In conclusion, the meta-heuristics algorithm extends the improvement algorithms allowing the passage through solutions that don't improve the solution in the perspective that will find a final better solution or even the global optimal[8]. Some examples of metaheuristics are Tabu-Search, Genetic Algorithms or Simulated Annealing.

## 3    Computational Work

For this computational work, it was used 30 datasets selected from the TSPLIB for testing purposes. Due to simplicity and computational speed reasons, datasets sizes were small (between 100 and 300 cities each). The total number of cities included in each library is presented in the name, for example, the library "*a280*" has 280 cities. Each line of the library represents a city with two coordinates *x* and *y* (Euclidean plan), so the formula to achieve the distance between two points ($d(i, j)$) is showed in Formula 1.

$$d_{ij} = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} \tag{1}$$

As referred in previous sections, the goal of this work is to acquire a final tour with the minimum distance to travel through the cities. To achieve this goal, it was developed an application that passes through several stages where different algorithms are computed in the following order: sequential algorithm (loading cities to an initial structure), Nearest Neighbor algorithm (constructive heuristic to create a good initial tour), 2-Optimal (improvement heuristic) and at last the metaheuristic Tabu-Search (implemented over 2-Opt heuristic whose purpose is to escape from the local optimal to achieve different solutions). The generic algorithm of the application can be seen in Figure 1.

| Step 1 | Get Cities Data |
|--------|-----------------|
| Step 2 | Generate Start Tour |
| Step 3 | Sequential Algorithm |
| Step 3.1 | Perform Constructive Heuristics |
| Step 4 | Chooses best Tour from Step 2 or Step 3 |
| Step 5 | Perform Improvement Heuristics |
| Step 5.1 | 2-Optimal |
| Step 6 | Perform Metaheuristic Algorithm |
| Step 6.1 | TABU Search |
| Step 7 | Show Results and Statistical Data (Interface) |

**Fig. 1.** General algorithm of the application

The main goal of this study wasn't the achievement of the optimal solution, but the application of different heuristics types and their impact in the problem resolution. With this idea in mind, it was chosen some simple heuristics that could represent the effort of each heuristic type to the problem resolution. Due to constructive heuristics nature, Nearest Neighbor algorithm was chosen because of the fast results it can achieve. Like with the constructive heuristic, a simple improvement heuristic was chosen from the "*k-opt family*", the 2-Optimal heuristic. As a metaheuristic, Tabu-Search was chosen due to its robustness, and also it was familiar to the author. This metaheuristic was applied over the 2-optimal heuristic

The application has a small interface where the libraries files can be specified and loaded into the application data structure with the possibility to define Tabu Search parameterization (defined in section 4, Table 1). This parameterization is made only for the metaheuristic algorithm with the goal of testing several parameters and their impact to the final solution. Finally, when the application has calculated its final solution, statically data is shown like, distances obtained for each heuristic and the final cities ordering.

A brief overview of each algorithm used in the application will be done focusing on their principal characteristics and generic algorithm.

### 3.1 Sequential Algorithm

The sequential algorithm is the most basic algorithm used. Its finality is the construction of an initial tour. This algorithm isn't a heuristic algorithm but just a way of loading the library cities data into the application data structure.

### 3.2 Nearest Neighbor Heuristic

The implemented constructive algorithm in the application is the Nearest Neighbor (NN) and its objective is the construction of an initial tour that will serve as a good starting tour for the improvement heuristics. The steps needed to perform this algorithm are explained in Figure 2.

| Step 1 | Choose start City as Current City (CC) |
| Step 2 | Mark CC as visited |
| Step 3 | While (there are cities not visited) Do |
| Step 3.1 | Get nearest city from CC that is not visited |
| Step 3.2 | Set CC as new chosen city |
| Step 3.3 | Mark CC as visited |
| | End while |
| Step 4 | Return to start city |

**Fig. 2.** Nearest Neighbor Algorithm

An extension to this algorithm is the repetition of every city as the starting point and then return only the best tour found[6]. This heuristic is called Repetitive Nearest Neighbor and was used in this scenario. The objective is to find the best solution possible for the NN because the minimum cost solution depends on its starting city, and a good solution is required in order for the improvement heuristic achieve better solutions. In this scenario NN algorithm will always retrieve the same solution for each library.

### 3.3 2-Opt Heuristic

Once a tour has been generated by the nearest neighbor heuristic, it was implemented an improvement heuristic over the solution obtained previously: 2-Opt[9]. This heuristic belongs to the k-Opt heuristics, where there is another famous heuristic named 3-Optimal. The reasons of choice for 2-Opt was due to its simplicity and because the dataset size was small, so a complex improvement heuristic is not necessary to achieve good results. The main idea of 2-Opt is the definition of a neighborhood structure on the set of all admissible tours. Its heuristic systematically looks at each pair of non-adjacent edges of the current tour and determines whether the tour length would be decreased by removing these two edges and adding the other possible pair of edges. In such a structure, the tour can iteratively be improved by always moving from one space to its best neighbor until no further improvement is possible [6, 10]. The generic algorithm can be seen in Figure 3.
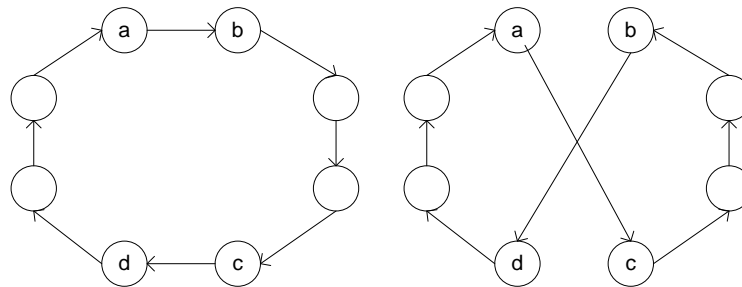
| Step 1 | Choose two arcs (a b) from Current Tour (CT) |
|---|---|
| Step 2 | While (there are pairs) Do |
| Step 2.1 | Choose two arcs (c, d) |
| Step 2.2 | Evaluate switch of arcs |
| Step 2.3 | If (evaluation<0) Then |
| Step 2.3.1 | Switch arcs |
| Step 2.3.2 | Invert Direction |
| | End If |
| Step 2.4 | Chose another two arcs |
| | End While |
| Step 3 | If (tour was improved) Then |
| Step 3.1 | Set CT equals to improved Tour |
| | End If |

**Fig. 3.** 2-Opt Algorithm

In practice, and shown in Figure 4, the 2-Opt movement consists in the replacement of two arcs (a, b) and (c, d) for another two arcs: (a, c) and (b, d). In the symmetric TSP case, the movement application makes an alteration in the solution cost that is given by Formula 2:

$$\Delta = d(a, c) + d(b, d) - d(a, b) - d(c, d). \tag{2}$$

If the formula's cost (delta) is negative, it means that the replacement of the arcs will produce a better solution. Considering a given orientation of the tour, it's also easy to see that upon the switch one of the paths, (b… c) or (d... a), a path will have to be inverted to maintain an admissible direction.



**Fig. 4.** Example 2-Optimal Movement

The final solution obtained with this heuristic will serve as starting solution for the metaheuristic algorithm applied over 2-Opt. In this scenario, since the 2-Opt algorithm depends on the starting solution provided by the NN algorithm, it will always retrieve the same solution for the next stage.

### 3.4    Tabu-Search

A common problem with neighborhood searches is that some of them can easily get trapped in a local optimum. This can be avoided by using a metaheuristic, in this case Tabu-Search (TS). The basic principle of TS is that when it encounters a local optimum it allows non-improving moves to the solution so it can search for solutions in a different neighborhood and, the possibility of better solutions. When this situation occurs, something needs to be done to prevent the search from tracing back its steps to where it came from. This is done by the use of memories, called tabu lists (TL), that record the recent history of the search. In its simplest form, a TL is a short-term memory which contains the solutions that have been visited in the recent past (less than $n$ iterations ago, where $n$ is the number of previous solutions to be stored) and are used to prevent cycling when moving away from local optimal through non-improving moves[11-12].

There are three ways of determining a TL size: fixed, a random value chosen from an interval or a dynamic TL size[13]. It's a known problem that TL size affects the performance of Tabu Search [14]. If it's too short it cannot escape from the local optimal, but if it's too long, it can limit and may force the heuristic to search values away from the optimal[15]. The size of the TL is also very important because depending on the type and size of the problem, different quality solutions are expected.

The most commonly used stopping criteria in TS are: after a fixed number of iterations (or a fixed amount of CPU time), after some number of iterations without an improvement in the solution (the criteria used in the majority of implementations) or when the objective reaches a pre-specified threshold value. In complex Tabu schemes, the search is usually stopped after completing a sequence of phases, being the duration of each phase determined by one of the above criteria[8].

For this scenario, TS algorithm had two termination criteria: based on a fixed number of iterations defined by the user, or after some iterations without an improvement to the solution.

As referred before, a metaheuristic isn't a new type of heuristic but an extension to an improvement heuristic, so it was applied over the 2-Opt heuristic. The TS algorithm was similar to the 2-Opt algorithm except it had two new concepts: a termination criteria (step 1 of Figure 5) and the TL implementation (step 3.3 and 3.4 of Figure 5). This algorithm will run until the number of iterations reaches the limit defined by the user, or the algorithm were running $n$ times without improvements on the solution. For each iteration, the arcs selected to switch are verified and if they are in the TL and if they not improve the solution they are ignored and another two arcs are chosen to be tested.

Despite the algorithm finds worst local solutions, it continues executing, since it can encounter better solutions in the next iterations (it allows to escape from a local optimum and search for other solutions in the neighborhood that in another case will never be considered). For each iteration new arcs are chosen and added to the TL, deleting the oldest entry if it's already full.

| Step 1 | **While (termination criterion not satisfied) Do** |
|---|---|
| Step 2 | Choose two arcs (a b) from Current Tour (CT) |
| Step 3 | While (there are pairs) Do |
| Step 3.1 | Choose two arcs (c, d) |
| Step 3.2 | Evaluate switch of arcs |
| Step 3.3 | If ((evaluation<0) or (**movement is not in Tabu List)**) Then |
| Step 3.3.1 | Switch arcs |
| Step 3.3.2 | Invert Direction |
| | End If |
| Step 3.4 | **Add Movement to Tabu List (delete oldest entry if necessary)** |
| Step 3.5 | Chose another two arcs |
| | End While |
| Step 4 | If (tour was improved) Then |
| Step 4.1 | Set CT equals to improved Tour |
| | End If |
| | End While |

**Fig. 5.** Tabu-Search applied to 2-Opt Algorithm

The TL size, number of executions and iterations were defined in running time.

## 4    Experimental Results

In order to achieve conclusions about the effects of heuristics on the final solution some experiments were made. The obtained measures were the total distance for each heuristic on each library. These values were compared for each heuristic type and to the optimal value of the library. The optimal value is the best value of the library known so far for the library. The main objective in this scenario is the comparison between the different algorithms relative to the optimal value. For that, it was used general specifications for all libraries which are demonstrated in Table 1.

**Table 1.** Application parametrization

| Heuristic | Parameter | Value |
|---|---|---|
| NN | Nr Executions | Total number of the test library's cities |
| Tabu-Search | Nr Executions | 300 |
| | Number Iterations | 100 |
| | Iterations without improvement | 5 |
| | Tabu list size (TL) | 5 and 7 |

Since NN algorithm is fast and the problem size is small, it was chosen the total number of cities on the test library as the total number of executions in order to achieve the shortest tour possible. This value is variable seeing as it depends on the total cities in the test libraries.

For Tabu-Search, two different Tabu List sizes of 5 and 7 were chosen.. A TL size between 5 and 12 are known for obtaining good results, so a value of 7 was chosen for testing purposes [14, 16]. Another TL of size 5 was tested in order to observe the

possible impact of different TL sizes on the same problem. The number of executions, iterations and iterations without improvements are experimental and were chosen due to good results obtained in first experiments.

The results obtained (in units) for all the libraries are shown in Table 2.

**Table 2.** Total Distance (in units) after applied Heuristics

| File | Optimal Distance | Initial Distance | NN | 2-Opt | Tabu-Search (TL: 7) | Tabu-Search (TL: 5) |
|------|------------------|------------------|------|-------|---------------------|---------------------|
| a280 | 2579 | 2819 | 3094 | 2765 | 2738 | 2724 |
| bier127 | 118282 | 393998 | 133971 | 130969 | 122169 | 123522 |
| ch130 | 6110 | 47801 | 7199 | 6715 | 6518 | 6456 |
| ch150 | 6528 | 52812 | 7078 | 6869 | 6577 | 6668 |
| d198 | 15780 | 22514 | 17810 | 16426 | 16230 | 16047 |
| eil101 | 629 | 2064 | 736 | 686 | 680 | 672 |
| gil262 | 2378 | 26296 | 2884 | 2696 | 2530 | 2552 |
| kroA100 | 21282 | 191394 | 24698 | 21749 | 21749 | 21457 |
| kroA150 | 26524 | 287850 | 31482 | 29250 | 28149 | 28409 |
| kroA200 | 29368 | 373943 | 34548 | 30820 | 29926 | 30315 |
| kroB100 | 22141 | 157185 | 25883 | 23442 | 22924 | 22627 |
| kroB150 | 26130 | 273236 | 31320 | 27669 | 26679 | 26932 |
| kroB200 | 29437 | 327452 | 35394 | 32106 | 31632 | 31741 |
| kroC100 | 20749 | 183465 | 23566 | 21971 | 21485 | 21551 |
| kroD100 | 21294 | 170991 | 26401 | 23305 | 22202 | 21382 |
| kroE100 | 22068 | 188350 | 24907 | 23201 | 22635 | 22390 |
| lin105 | 14379 | 36478 | 16939 | 15228 | 15173 | 14988 |
| pr107 | 44303 | 62757 | 46678 | 45338 | 44575 | 44575 |
| pr124 | 59030 | 98943 | 67057 | 60665 | 60304 | 60525 |
| pr136 | 96772 | 287026 | 114561 | 107354 | 104855 | 105320 |
| pr144 | 58537 | 93524 | 60963 | 60753 | 60753 | 60753 |
| pr152 | 73682 | 160980 | 79567 | 76108 | 75397 | 75189 |
| pr226 | 80369 | 110416 | 92553 | 83450 | 81964 | 82386 |
| pr264 | 49135 | 77976 | 54491 | 53889 | 52349 | 52401 |
| pr299 | 48191 | 83508 | 58288 | 51069 | 50002 | 49990 |
| rat195 | 2323 | 4038 | 2629 | 2436 | 2396 | 2405 |
| rd100 | 7910 | 50561 | 9427 | 8244 | 8225 | 8222 |
| ts225 | 126643 | 276532 | 140485 | 134761 | 128505 | 128691 |
| tsp225 | 3916 | 10300 | 4633 | 4187 | 4049 | 4011 |
| u159 | 42080 | 43376 | 48587 | 42976 | 42433 | 42433 |

The first and second columns of Table 2 are the name and the known optimal value of each library, in units, respectively. The next four columns represent the values obtained for the algorithms: Initial Distance shows the initial tour distance before any heuristic algorithm was applied, and, subsequently, the three heuristic algorithms used in the application: Nearest Neighbor (NN), 2-Optimal and Tabu-Search algorithm with a TL size of 7 and 5.

Table 3 shows the values obtained for each heuristic algorithm related to the optimal distance. This table allows seeing a direct relation of the algorithms impact on the optimal solution.

**Table 3.** Distances (in percentage) to Optimal after applied Heuristics

| File | Optimal Distance | Initial Distance | NN | 2-Opt | Tabu-Search (TL: 7) | Tabu-Search (TL: 5) |
|------|------------------|------------------|-----|-------|---------------------|---------------------|
| a280 | 2579 | 9,29% | 19,98% | 7,20% | 6,17% | 5,62% |
| bier127 | 118282 | 233,10% | 13,26% | 10,73% | 3,29% | 4,43% |
| ch130 | 6110 | 682,34% | 17,82% | 9,91% | 6,68% | 5,66% |
| ch150 | 6528 | 709,01% | 8,43% | 5,23% | 0,75% | 2,14% |
| d198 | 15780 | 42,68% | 12,86% | 4,10% | 2,85% | 1,69% |
| eil101 | 629 | 228,22% | 17,07% | 9,03% | 8,11% | 6,84% |
| gil262 | 2378 | 1005,79% | 21,27% | 13,39% | 6,39% | 7,32% |
| kroA100 | 21282 | 799,32% | 16,05% | 2,20% | 2,19% | 0,82% |
| kroA150 | 26524 | 985,24% | 18,69% | 10,28% | 6,13% | 7,11% |
| kroA200 | 29368 | 1173,30% | 17,64% | 4,94% | 1,90% | 3,22% |
| kroB100 | 22141 | 609,93% | 16,90% | 5,88% | 3,54% | 2,20% |
| kroB150 | 26130 | 945,68% | 19,86% | 5,89% | 2,10% | 3,07% |
| kroB200 | 29437 | 1012,38% | 20,24% | 9,07% | 7,46% | 7,83% |
| kroC100 | 20749 | 784,21% | 13,58% | 5,89% | 3,55% | 3,87% |
| kroD100 | 21294 | 703,00% | 23,98% | 9,44% | 4,26% | 0,41% |
| kroE100 | 22068 | 753,50% | 12,86% | 5,13% | 2,57% | 1,46% |
| lin105 | 14379 | 153,69% | 17,81% | 5,90% | 5,52% | 4,24% |
| pr107 | 44303 | 41,65% | 5,36% | 2,34% | 0,61% | 0,61% |
| pr124 | 59030 | 67,62% | 13,60% | 2,77% | 2,16% | 2,53% |
| pr136 | 96772 | 196,60% | 18,38% | 10,93% | 8,35% | 8,83% |
| pr144 | 58537 | 59,77% | 4,14% | 3,79% | 3,79% | 3,79% |
| pr152 | 73682 | 118,48% | 7,99% | 3,29% | 2,33% | 2,05% |
| pr226 | 80369 | 37,39% | 15,16% | 3,83% | 1,98% | 2,51% |
| pr264 | 49135 | 58,70% | 10,90% | 9,68% | 6,54% | 6,65% |
| pr299 | 48191 | 73,29% | 20,95% | 5,97% | 3,76% | 3,73% |
| rat195 | 2323 | 73,81% | 13,15% | 4,88% | 3,14% | 3,53% |
| rd100 | 7910 | 539,20% | 19,18% | 4,22% | 3,98% | 3,94% |
| ts225 | 126643 | 118,36% | 10,93% | 6,41% | 1,47% | 1,62% |
| tsp225 | 3916 | 163,02% | 18,31% | 6,93% | 3,40% | 2,43% |
| u159 | 42080 | 3,08% | 15,46% | 2,13% | 0,84% | 0,84% |

The format of Table 3 is similar to Table 2, but the obtained results are in percentage and related to the libraries optimal value.

Table 2 and 3 shown that the sequential algorithm, which is the simplest of the algorithms used, is the less efficient except in two cases, A280 and U159, where the variation rate related to optimal value was very small. This happened because the cities were already well ordered in the library. The NN algorithm improved, in general, the initial solution except in some special cases referred previously. The results obtained by this algorithm are not very efficient, but it provides a good initial solution for the 2-opt algorithm. As expected, the 2-Opt algorithm made a boost on the solutions quality, improving all the solutions, in average, 9.01% points and 6.38% points worst than the optimal value. In this scenario, and due to the initial solution is always the same, 2-Opt solutions always retrieved the same solution for each library, so it's independent of the number of executions. TS solutions were even better than 2-opt improving the solution in all the test libraries.

As tables 2 and 3 illustrate, the results obtained are smaller every time an algorithm is applied to the tour. For a better understanding of the algorithms

performance, Figure 6 shows a graphic with the average for each algorithm related to the library optimal value.
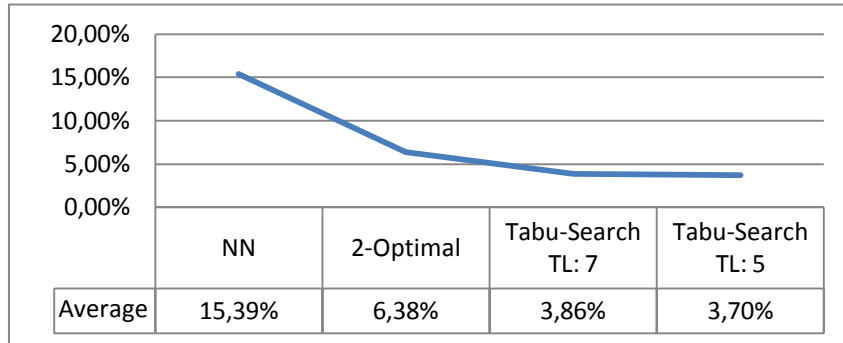


| | NN | 2-Optimal | Tabu-Search TL: 7 | Tabu-Search TL: 5 |
|---|---|---|---|---|
| Average | 15,39% | 6,38% | 3,86% | 3,70% |

**Fig. 6.** Heuristics Average Variation to Optimal Distance

As shown in figure 6, the TS algorithm was the best algorithm, but even with the same algorithm, there were some differences caused by the parameterizations. We can see that, in average, the best results were obtained with a TL of size 5, but there were some libraries where a TL of 7 worked better. In some instances, the differences are significant like in *eil101* and *kroA100*. Different TL sizes offer different solutions so different parameters correctly applied on each library, may result in better solutions. The heuristics approaches have proved to guarantee good approximations to the optimal solution, especially with the implementation of a metaheuristic.

## 5      Conclusions and Future Work

In this paper, was demonstrated the use of different heuristic types and compared them for the Traveling Salesman Problem. Those heuristics were applied in 30 common problems from TSPLIB with a range of 100 to 300 cities. It was exposed that heuristics can be considered when optimal solutions are not required. Combination of different types of heuristics should be used resulting in better solutions. In all the studied problems, metaheuristics have improved the solution so they are fundamental to obtain the best possible solution.

It's also evident that algorithms configurations should be chosen carefully (ex: TL size) because bad configurations can retrieve worst solutions than those expected, as we saw in subsection 3.4. Quality solutions and execution time should be considered in the problem resolution: obtain high quality solutions in the shortest time possible. It's was not the study objective the achievement of optimal solutions, but the techniques it can be used to accomplish it.

Even with few tests, it was demonstrated that heuristics approaches can accomplish good results and are real alternatives to exact methods when solving TSP problems that don't required the achievement of optimal solutions.

In the future, different initial solutions (ex: NN different solutions) should be applied to observe the algorithms comportment. Also, different parameterizations, especially on the Tabu Search, like number of iterations or Tabu List sizes (adapted to

the problems) should be considered. This could be done by doing more experiments and observe the obtained results.

Test libraries with more cities should be considered to observe the behavior of the application on them and make some adaptations, perhaps dynamic. Also, different heuristics should be applied to the problem, so more data should be collected and compared. Although execution times were not discussed, information about it is needed to do more valid conclusions.

# References

1. Gutin, G., Punnen, A.P.: The Traveling Salesman Problem and Its Variations. Kluwer Academic Publishers (2007)
2. Jourdan, L., Basseura, M., Talbia, E.-G.: Hybridizing exact methods and metaheuristics: A taxonomy Eur J Oper Res 199, 620-629 (2009)
3. Ponnambalam, S.G., Aravindan, P., Chandrasekaran, S.: Constructive and improvement flow shop scheduling heuristics: an extensive evaluation. Production Planning & Control; 12, 335-344 (2001)
4. TSPLIB-Symmetric traveling salesman problem (TSP), http://elib.zib.de/pub/mp-testdata/tsp/tsplib/tsplib.html
5. INFORMS Computer Society, http://glossary.computing.society.informs.org/second.php?page=T.html
6. The Traveling Salesman Computational Solutions for TSP Applications. Springer Berlin / Heidelberg (1994)
7. Lenstra, J., Kan, A.R.: Some simple applications of the travelling salesman problem. Operational Research Quarterly (1975)
8. Gendreau, M.: An Introduction to Tabu Search. Handbook of Metaheuristics, pp. 18. Kluwer Academic Publishers (2002)
9. Hansen, P., Mladenović, N.: First vs. best improvement: An empirical study. Discrete Appl Math 154, 802-817 (2006)
10. Rosenkrantz, D.J., Stearns, R.E., II, P.M.L.: An analysis of several heuristics for the traveling salesman problem. In: Netherlands, S. (ed.) Fundamental Problems in Computing, pp. 45-69 (2009)
11. Glover, F., Marti, R.: Tabu Search. Metaheuristic Procedures for Training Neutral Networks, vol. 36, pp. 53-69. Springer US (2006)
12. Gendreau, M.: An introduction to Tabu Search. Handbook of metaheuristics, vol. 57, pp. 37-54. Springer New York (2003)
13. Salhi, S.: Defining tabu list size and aspiration criterion within tabu search methods. Comput Oper Res 29, 67-86 (2002)
14. Glover, F.: Tabu Search-- Part I. ORSA journal on computing 1, 190 (1989)
15. Tsubakitani, S., Evans, J.R.: Optimizing tabu list size for the traveling salesman problem. Comput Oper Res 25, 91-97 (1998)
16. Glover, F.: Future Paths for Integer Programming and Links to Artificial-Intelligence. Comput Oper Res 13, 533-549 (1986)

# Methods, Difficulties and Practical Solutions in Implementing CMMI® High Maturity Levels

Isabel Lopes Margarido*[1], Marco Vieira[2], Raul Vidal[1]

[1]Faculty of Engineering, University of Porto, Rua Dr. Roberto Frias, 4200-465 Porto, Portugal
{isabel.margarido, rmvidal}@fe.up.pt

[2]Faculty of Sciences and Technology, University of Coimbra, Pólo II, Pinhal de Marrocos,
3030-290 Coimbra, Portugal
mvieira@dei.uc.pt

**Abstract.** Many companies face problems when implementing CMMI® (Capability Maturity Model Integration) High Maturity Levels since it is necessary to implement complex practices that include collecting metrics, do quantitative management and have effective performance models that allow predicting the course that controlled processes are going to take. A key question is whether CMMI® implementers are using the scientific knowledge available or taking other kind of sustained guidance to achieve CMMI® levels 4 and 5. In this study we analyse the scientific recommendations that are applicable to the CMMI model and what is being done in Critical Software, S.A., a Portuguese company recently appraised CMMI® level 5, where we performed a case study. With our study we verified that the company did not apply scientific knowledge in the CMMI High Maturity Levels implementation and used tools such as Balanced Score Card, Gold Question Metric Indicator, Ishikawa tools and Six Sigma.

**Keywords:** Software Engineering, CMMI High Maturity Levels, Project Management Metrics, Software Lifecycle Metrics.

## 1    Introduction

Costumers in healthcare and defence, require CMMI Level 5 for important contracts and recognize that it gives high predictability, better engineered products for scalability, maintainability, adaptability, and reliability [1]. CMM (Capability Maturity Model) was created by the Carnegie Mellon University (CMU) and the SEI (Software Engineering Institute) in 1991. An extended version, CMMI was created in 2002 and in 2006 CMMI for Development version 1.2 was published [1]. CMMI is a process improvement maturity model for the development of products and services [2]. Being only a model it does not provide strict guidelines on how its practices and goals may be implemented.

Many companies face problems when implementing CMMI® High Maturity Levels (HML) that arise from complex practices such as measurement and quantitative management or the use of effective performance models for predicting

the future course of controlled processes. In fact, part of the difficulties found in the processes evolution and new Process Areas (PA) implementation are related to the need of moving towards a statistical thinking and quantitative management [3]. That problem has already been highlighted by the SEI, when it concluded that some companies did not understand the statistical nature of CMMI level 4 and certified CMMI HML companies did not had a consensus on the necessary characteristics of level 4 [4].

Certified companies generally do not share the implementation knowledge to maintain the competitive advantage, therefore new companies aspiring to become certified on the HML lack proper guidance.

Nowadays the scientific knowledge is approaching CMMI implementation associated with development, management and improvement methodologies. There is also literature from the 70's to our days that introduces measurement methodologies for project management and quality assurance [5]. A key question is to understand whether companies are using this knowledge base to support CMMI implementation.

Critical Software, S.A. (Critical) is a SME (Small Medium Enterprise) Software House (around 400 employees) working on a wide market that includes, among other types of customers, highly demanding clients such as space, industry, defence and government. Critical achieved CMMI level 3 in 2006 [6]. In 2007 a SCAMPI$^{SM}$-C[1] (Standard CMMI Appraisal Method for Process Improvement) for level 5 allowed further analysis of the company's Quality Management System (QMS) improvement opportunities, the company was growing and the QMS should support employees on their daily tasks in order to develop high quality products [7]. On the 18th of December of 2009 Critical was appraised CMMI level 5, the first Portuguese company to achieve High Maturity Level Certification.

Our investigation intents to prove the following hypothesis: *h1- many companies adopt CMMI level 5 motivated by client demands*[1]. Benefits for adopting CMM/CMMI such as quality, cost/schedule and customer satisfaction improvement, add predictability and rework reduction were demonstrated [8]. A different motivation from the one indicated by Sutherland [1] can drive our future work objectives, by indicating which metrics are important to collect in order to monitor those goals*; h2- companies follow practitioners and SEI guidance to implement CMMI level 5*. Practitioners follow proven effective methods that other practitioners have already used [3], ignoring scientific knowledge. We identify scientific knowledge adequate for CMMI HML implementation and in our future work we will test it in practice*; h3- dissemination of new practices is not 100% effective.* In the year of 2000 a survey was conducted to understand what CMMI level 4 and 5 companies used in the CMMI implementation and it shown that practices were not clearly institutionalized [9]. We intend to identify barriers to dissemination, to find ways of improving it, in our future work. To answer our questions we performed a case study in a SME Software House.

When performing a case study in a HML company it is necessary to take under consideration what CMMI HML implies: 1.the organization collects metrics for controlling, managing and improving processes; 2. projects are quantitatively

---

[1] SCAMPI is a service mark of the Carnegie Mellon University.

managed; 3. process improvements are statistically proven and controlled, and mapped with organization's goals; 4. process, procedures and projects are monitored and when problems occur root causes are identified and resolutions are implemented. Our method was composed by four phases: 1. *Bibliographic Review*, finding scientific knowledge to support CMMI HML implementation; 2. *Relate Scientific Knowledge with the Model*; 3. *Survey a CMMI 5 SME*, in order to understand how the company put HML into practice; 4. *Analysis and Conclusions*.

The company's solutions came out from SEI additional documentation, books and other implementers' suggestions. From our scientific research and our case study we gathered tools to implement HML that are used in practice, tolls from the scientific knowledge and context information that will be useful in our future research work.

This section introduces the problem under investigation, hypothesis, applied method and results. In Section 2. *Background and Related Work*, we present the CMMI levels 4 and 5 process areas and tools provided by science that may be used to implement them. Section 3. *The Case Study Methods and Goals*, presents the methodology used to perform this investigation. Section 4. *Findings*, explains how the organization moved from level 3 to 5. Section 6 concludes this paper eliciting resources to implement CMMI HML and indicates future research opportunities.

## 2    Background and Related Work

Table 1 presents the CMMI levels 4 and 5 Process Areas (PA) purpose and some of their Specific Practices (SP).

**Table 1.** CMMI HML PAs, purpose and some of their SPs [2].

| ML 4 PAs | Details |
|---|---|
| Organizational Process Performance (OPP) | **Purpose:** "establish and maintain a quantitative understanding of the performance of organization's set of standard processes in support of quality and process-performance objectives, and to provide the process-performance data, baselines, and models to quantitatively manage the organization's processes." <br> **SPs:** standard processes; lifecycle model descriptions – Process Performance Models (PPM); tailoring criteria and guidelines; organisation's measurement repository; organisation's process asset library; work environment standards. |
| Quantitative Project Management (QPM) | **Purpose:** "quantitatively manage project's defined process to achieve the project's established quality and performance-objectives." <br> **SPs:** establishing and maintain the project's quality and process performance objectives; based on stability and capability data, compose the project's defined process; select the sub-processes for statistical management; manage the project performance and statistically manage the sub-processes performance. |
| **ML 5 PAs** | **Details** |
| Organizational Innovation and Development | **Purpose:** "deploy incremental and innovative improvements that measurably improve the organization's processes and technologies. The improvements support the organization's quality and process performance objectives, which derive from the organization's business objectives." |

| | **SPs:** collecting/identifying and analysing improvement proposals and innovations; conduct pilot projects to select which ones are worth for implementation and selecting improvements for a planned and managed deployment. Improvement effects need to be measured |
|---|---|
| Causal Analysis and Resolution | **Purpose:** "identify causes of defects and other problems and take action to prevent them from occurring in the future." <br> **SPs:** selection of defects and analysis of their causes; implementing solutions; evaluating the changes' effect and recording data. |

Measurement Analysis (MA) is a Maturity Level (ML) 2 PA, crucial for HML, it establishes measurement objectives; specifies measures, data collection storage and procedures; analysis procedures. This PA SPs also include the collection and analysis of measurement results.

The scientific community analyses measurement in the Software Engineering and Management areas. To implement CMMI HML it is necessary to have management and software engineering metrics and tools support.

The Balanced Score Card (BSC) created by Kaplan and Norton in 1992 was evolved by the authors along the years becoming in 2001 "a strategy implementation tool to facilitate and control performance measurement and management." [10] BSC may be used to map processes which performance is relevant for the organization business goals and allows it to control processes, as required by OPP.

The Goal Question Metric (GQM) tool is identified as a tool that structures data collection that affects organization's goals and its success in diagnosing the efficiency and effectiveness of processes is recognised [11].

Investigators in the software engineering area propose the usage of Six Sigma to implement HML, namely in the improvement of processes. Gonçalves *et al* apply Six Sigma tools combined with CMMI process areas in the identification, classification and prioritization of improvement. Six Sigma tools are applied according with the complexity degree of innovations [12].

Xiasong *et al* propose a model that integrates DMAIC (Define Measure Analyse Improve Control) tool with CMMI process areas OPP, QPM, MA and CAR. The authors do not explicitly refer OID [13].

Schalken and Van Vliet propose a framework that combines quantitative and qualitative data to analyse software engineering processes, applicable to post-mortem projects [11]; we consider that such a tool is adaptable to implement OID, the iterative process is useful to implement CAR in problems identified on closed projects.

Cohan and Glazer refer agile development teams improvements to achieve ML 5, however their success has not yet been proven [14]. Agile is a process performance improvement in companies that already achieve ML 5 [1], [15].


## 3    The Case Study Methods and Goals

In order to understand how practitioners implement the HML and which bibliographic support they are using we performed a case study in Critical Software, S.A., a company appraised CMMI Level 5.Our research method followed a set of procedures organized in 4 phases.

*Phase 1: Bibliography Review* – research the existent scientific publications and other bibliography about the CMMI HML implementation, measurement (and related) tools, and improvement tools.

Phase 1 allowed us to know the directions given by the scientific community for the implementation of CMMI levels 4 and 5.

*Phase 2: Relate Scientific Knowledge with the Model* – establish a connection between our findings and the CMMI for Development model, in particular to the HML Process Areas.

*Phase 3: Survey a CMMI 5 SME* – understand what SMEs are doing to achieve ML 5.

Phase 2 gave us the necessary information to begin phase 3. We had the practices related to the CMMI model, so we could understand what practitioners were doing and establish a relationship with the outcomes of our bibliographic review.

In order to understand difficulties in implementing level 5, we performed a survey in Critical Software, S.A., a company evolving from level 3 to the level 5. We needed to understand the process that SMEs go through when implementing HML, including: knowledge building; evolution of existent processes; creation of non-documented processes, characteristic of the Process Areas of Levels 4 and 5; support tools; team building and adaptation to the new mindset; deployment process; and adherence to the new practices. Our investigation methods included: *SCAMPI-A outcomes* analysis; *Workshops*; *Data collection and analysis tools study*; *Interviews* with the CMMI appraisal team, the project sponsor, and employees; Consulting *process* and *performance models*.

This phase included ongoing interviewees' identification, derived from documentation. From interviews answers new relevant stakeholders were identified.

*Phase 4: Analysis and conclusions* – gather recommendations for Software Houses in the SMEs category and advance with the tendency evolution for the CMMI for Development model.

To perform the case study the following field methods were applied:

- *Exploratory interviews* – a method that would allow us to access information that in direct questions are blocked;
- *Documentation analysis* – analysis of the Quality Management System (QMS), documentation of the CMMI program, requirements of Information Systems (IS), projects documentation;
- *Tools analysis* – verifying changed/acquired tools functionalities that support levels 4 and 5. Taking into consideration the high maturity companies expected characteristics we assume that IS exist to support metrics collection, analysis, treatment and interpretation; therefore we needed to understand which tools were deployed or changed. We assumed the following tools existed and would probably have been updated: Project management; Defect tracking; Improvements management; Data analysis tool.
- *Individual interviews* – directed interviews with specific questions derived from the previous methods, knowledge of the company, bibliographic research. In the individual interviews, questions derived from other interviewees inputs and the own interviewee were posed. To follow IS alterations we needed to interview the IT manager and the changes implementers.

# 4  Findings

After performing all the interviews and consulting artefacts (program and project documents) and software tools, the relevant topics were categorized. A transcription of the interviews was made placing the evidences under the identified categories.

This chapter presents the results of the field study in subchapters related to each category.

## 4.1  Program History

The CMMI program evolution was marked by some difficulties and paradigms. A project that had a planned duration of 6 months (the program manager informed that there are records of companies who achieved CMMI level 5 in that time) reached its goal 2,5 years later. CMMI 5 project, as it was known in Critical started in June of 2007. According to the program sponsor there was no clear vision or guidelines to achieve levels 4 and 5 and the company had little perception of what should be done. He also considers SEI's guidance was insufficient: did not have a vision on the questions of the HML Implementation, did not provide guidance on how to do it.

Unanswered questions were: *How to put CMMI levels 4 and 5 into practice? How to collect metrics in practice? Are there any tools to support metrics collection?* Those questions remained unanswered until the middle of 2008. After a year of work without results, the SEI indicated to the company the Six Sigma tool, as being appropriate to implement levels 4 and 5. The program sponsor revealed that Six Sigma provided a set of tools for control and management and with those the theory became systematized. According to the interviewee Six Sigma identifies the processes usage and readjustment of measures in order to improve the model, helps companies to percept the models and their behaviour – whether they are performing as expected.

## 4.2  Motivations

The Critical CMMI program sponsor indicated what motivated the CMMI 5 project:
1.  Respond to business strategy, the company is still growing and moving to markets with highly demanding clients, with higher maturity and more complex necessities.
2.  Continuous improvement, the company recognized the value of achieving level 3, the processes were already defined but there was a lack of efficiency and improvement, there was still a path to walk.

Considering *h1- many companies adopt CMMI level 5 motivated by client demands* it is partially proven because the maturity improvement was result of business goals. The program sponsor adds "In the USA, clients normally demand to work with companies that have at least CMMI level 3." The organization never felt the ML as a limitation; level 3 or level 5, the ML in business is considered as an add-on that brings value and therefore companies take profit of it.

## 4.3   Program Management

The program management itself faced the difficulties of poor information on how to implement the HML, allied with lack of human resources and limited predictability on the project's real duration.

**Human Resources.** The human resources involved in the project and role are listed in table 2. The selection criteria were to involve people with experience, who had proven knowledge of their area that could bring added value to the CMMI5 project. People were involved as appraisal team members, areas consultants and experts.

**Table 2.** Team members directly involved in the CMMI5 project.

| Element | Role |
| --- | --- |
| Sponsor | Involved in progress and evaluation meetings. Support the program management resolving conflicts in the organization: human resources, internal software tools development competition. |
| Quality Department (QD) Organization's areas members | QD had to be focused on CMMI5 project and avoid reticence to evolve the QMS. Having areas members support, dissemination on those areas would be guaranteed (Human Resources, Business Development, Financial, Engineering, etc.). Appraisers. |
| Directors | Guarantee their support on practices and goals as well as their commitment with the program. |
| IT Manager Tools responsible | Evolve IS to sustain HML practices and 'Statistical Thinking'. Relevant stakeholders and key users were involved in the tools projects (ex: QD member to validate requirements, Software Product Assurance Manager to validate process and quality charts). |
| Top Management Area Managers | Define goals, relevant processes and indicators in each process. |
| Process owners, QD Manager, pilot projects | Processes updates and processes creation. Some projects involved key-users and technical experts. |
| Scientific Community | Involved in some processes definition (ex: reuse and estimation). |

**Project Plan.** There were several different plans created after concluding that milestones could not be met. The major rupture of the project plan occurred when the company finally adopted Six Sigma tool to break the how to implement HML paradigm. Before this achievement the effort to implement new PAs based on statistical thinking was unsuccessful. After that it was necessary to implement Six Sigma following the CMMI model and introducing changes in software tools and QMS processes. The program sponsor referred that an additional layer regarding metrics collection, control and management had to be added, which required certain projects and software tools adaptation, both becoming more detailed.

The program manager elucidated us on the several plans and their failures causes which are mapped in table 3.

**Table 3.** Program plans and reasons for SCAMPI A re-schedule.

| Plan | Details |
|------|---------|
| 1st (June07 – May08) | SCAMPI C for ML 5: August 2007; SCAMPI A 6 months later. **Problem:** insufficient time to address identified problems. |
| 2nd (December08) | **Reason:** Readiness review in Jan08 shown that project was slipping. GAPs projects were not finished. |
| 3rd (April09) | **Reason:** Program was running late. Six Sigma was identified as HML implementation tool and its deploy would take time. |
| 4th (October09) | **Reason:** Several trainings (Six Sigma, CMMI High Maturity Practices). New way of work was in place and pilot projects were in progress. |
| 5th (Oct09 and December09) | **Reason:** Delay on software tools implementation; key team members moved to other critical projects. Partitioned SCAMPI A scheduled levels 2, 3 for October and HML for December. |

Median duration of CMMI5 projects is approximately 2,5 years (19 months to go from ML 3 to ML 4 and 13 months to go from ML4 to ML5 [16]), few cases are known of companies that were appraised such level in 6 months and we believe that those cases are of companies already mature, with measurement systems already in place, therefore their effort only implies processes performance monitorization and *mis en place* of QPM, OPP, OID and CAR.

**Implementation Steps.** The HML implementation comprehended two phases:

- *Phase 1- Architecture Definition*, responsibility of the Program Manager, included steps as understanding organization's objectives; monitoring the processes performance to reach those objectives; studying the metrics collection tools. During this phase changes on the QMS and IS occurred.
- *Phase 2- Design*, responsibility of the IT human resources, included steps as discussing tools needs the IT Manager; specify tools requirements; project management tool improvement and integration with other tools when reporting effort: defect tracking, UML tool (requirements and tests). During this phase Six Sigma Green Belt training, quantitatively managed pilot projects and on job training were taking place.

## 4.4 Implementation

**Theoretical Support.** The theoretical support for the program included trainings, books and SEI articles and is resumed on table 4.
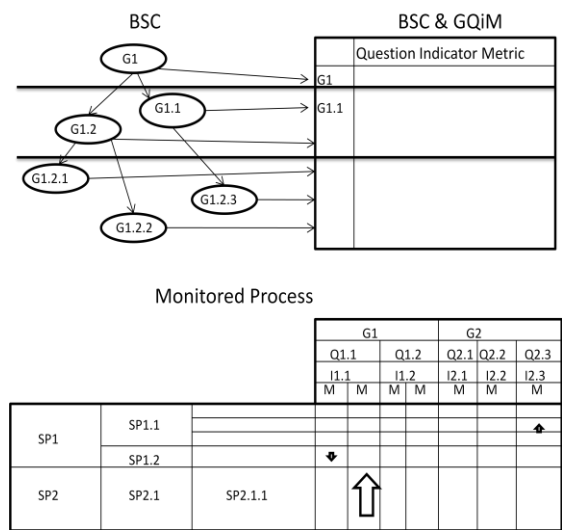
**Table 4.** Theoretical support used by practitioners.

| Artefacts | Details |
|-----------|---------|
| Trainings | SCAMPI appraisal team training; Six Sigma training involving people with different roles on adequate training sessions; Green Belt training; Six Sigma consultant supporting pilot projects kick-off; SEI training: Understanding High Maturity Practices. |
| Books | *Six Sigma Pocket Guidebook*; Donald Wheeler: *Understanding Variation: The Key to Managing Chaos and Understanding Statistical Process* |

| | |
|---|---|
| | *Control*; Grady: *Practical Software Metrics for Project Management and Process Improvement*; Carlton and Florac: *Measuring the Software Process*; Grady and Caswell *Software Metrics: Establishing a Company-Wide Program*; Kan: *Metrics and Models in Software Quality Engineering*; Siviy, Penn and Stoddard: *CMMI and Six Sigma*; Augustine A. Stagliano: *Six Sigma Advanced Tools Pocket Guide*; Michael Brassard and Diane Ritter: *The Memory Jogger*; Michael L. George, David Rowlands, Mark Price and John Maxey: *The Lean Six Sigma Pocket Tool Book*; John H. Baumert and Mark S. McWhinney: *Software Measures and the Capability Maturity Model* and [2]. |
| Other | Internet researches – on processes and GQiM (Gold Question indicator Metric); IBM taxonomy for defects classification; SEI's documentation. |

The referred bibliography is composed by practitioner's books and SEI documentation. Such evidence corrobates hypothesis *h2- companies* follow *practitioners and SEI guidance to implement CMMI level 5*.

**Metrics Establishment and Collection.** Metrics were established by applying the BSC, which mapped organization goals that were drilled down to lower level goals (organization – areas – operation). Metrics were derived by using the Goal Question indicator Metric (GQiM). The most relevant goals/sub-goals for business strategy were elicited and realistic objective values established for those goals indicators. The process steps are illustrated in figure 1.



**Fig.1.** Mapping *BSC* into *GQiM*, into processes and *sub-processes (SP)*. Monitored processes are being followed. Realistic *metric (M) goals (G)* are established (may be to decrease a metric value, such as number of defects; or increase a metric value, as % of code being reviewed).

Indicators for goals are calculated through metrics which need to be collected. Each indicator may be guaranteed by more than one process and sub-processes. In order to understand which Process Performance Models are necessary each process is studied separately for each GQiM. Each process of the QMS has a set of activities

(sub-processes) that on its details describes the steps to be followed. Some of them do not have meaning for the metrics but others influence the metric behaviour, imply usage of adopted tools and may be derived in performance goals.

Taking into account that Earned Value (EVA) metric fulfils CMMI level 3 requirements but is insufficient for level 5, because it does not concern with the quality of produced work, only with task closure, a triangle of significant metrics was established: effort (measured in hours), volume (measured in thousand lines of code (KLOC)) and defects (number). These measures allow understanding the execution speed, volume of work produced and its quality. The metrics are being collected at the operational level. Table 5 maps metrics with the tools that collect or use them.

**Table 5.** Each tool collects data on certain metrics which is integrated in the data warehouse tool, allowing data management and analysis and generating cubes of information. The project management tool provides monitoring information.

| Tool | Metrics Information |
|---|---|
| Project Management | Collects effort per phase. Effort spent on fixing defects is mapped with correspondent defects; effort spent on tests (specification or execution) is mapped with tests. Provides graphical information for project management: *S-Curve*, *Burndown*, *SPI* and *CPI*. Process metrics: *defect detection rate*, *defect fix cost*, each performance process is graphically represented for testing (unit, system, acceptance) and code reviews; *code review rate* (KLOC reviewed per hour). Provides graphics with organization information from all projects: *estimation error, SPI, CPI* and *defects by reporter* (internal or client). |
| Defect Tracking | Identifies defects reporter: client, internal. Defects are imported by the project management tool for effort measurement purposes. |
| UML | Requirements and tests are imported by the project management tool for effort measurement. *Requirements stability* is read by the data warehouse tool. |
| Data Warehouse | Is the core of the IS and integrates all metrics. It provides cubes for metrics analysis, graphics construction and integrity checks on data. |

**Dissemination.** The introduced changes dissemination was performed by several means: *1. Documentation* – processes, procedures, guidebooks were deployed on the QMS; PMO guidelines support project managers on new practices; *2. Intranet* – videos about the practices, PPB (Process Performance Baselines), PPM, BSC and announcements; *3.Training* – open workshops, practices presentations, brainstorming sessions; *4.On job training (pilot projects)* – projects are accompanied by trained people, *patronos*, and a SPAE (Software Product Assurance Engineer); Quality Manager participates in progress and data integrity check validation meetings and clarifies doubts; *5. Other* – a resume tool is available in each workstation with compilation of changes made on tools and explanation of concepts.

Even though people realize communication has improved there are still difficulties in the operational layer due to dissemination deficiencies. There were no trainings on software tools usage and associated concepts changes. Some tools help is outdated and in others only 'tool tips' and brief explanations are provided. People using cubes

generated by the data warehouse tool, started their work unaware of how to build graphics to understand information, only after seeing how it was done they learnt.

SPAEs and Project Managers (PM) had Quantitative Management training, although PMs did not participate in the practical part of the workshop. However, dissemination agents, such as SPAEs, had not a complete understanding on some of the metrics and what to monitor. They only guaranteed processes baselines were being respected and performed integrity checks. Actions on metrics at operational level are still very few. Therefore hypothesis *h3- dissemination of new operation practices is not 100% effective* has been proved.

### 4.5    Adopted Tools

BSC is used in goals elicitation and identification of performance objectives. The company defined the BSC perspective based on its values: Financial, Customer, Learning and Internal Processes. Some perspectives are aligned with Braams four critical perspectives: financial, customer, learning, and growth [10].

The company applied the Six Sigma tool not only for processes improvement but to answer measurement related needs. Six Sigma provided better understanding on processes, sustained quantitative projects management and supported quantitative quality assurance. Each project establishes performance objectives at the kick-off, controls processes application with the organization's baselines, performs integrity checks on data and, whenever problems occur, invokes the CAR process to determine the root cause of the problem and determine its resolution.

CAR is performed by applying the fishbone diagram, one of the Ishiwaka tools. Other Ishikawa tools, *Pareto chart*, *run chart*, *histogram*, *scatter diagram* and *control chart*[17] are also applied on metrics analysis, PPMs construction and PPBs control.

The tools used for HML implementations are similar to the ones used by other practitioners [3], [9]. To implement HML, practitioners use methods that have proved to be effective in practice and their choices are driven by the SEI and consultants.

## 5    Conclusions and Future Research

The following tools may be used in practice by SME Software Houses to sustain CMMI levels 4 and 5: *BSC*, *GQiM* [18]. *Ishiwaka* tools are applicable in Quality Assurance activities, from Quantitative Management to Causal Analysis and Resolution. Quantitative Project Management is being performed applying Six Sigma [19]. The Six Sigma tool is also being used to develop Process Performance Models.

Regarding the hypotheses of this study *h2* and *h3* were confirmed. Scientific knowledge needs to be tested in the industry and dissemination process needs to be improved. *h1* was partially confirmed, companies also intend to improve performance and quality which indicates what kind of metrics will need to be monitored.

For future research we consider that in a constantly changing world the scientific direction of CMMI is pushing towards agile management and development techniques, which may rapidly adapt to changes. CMMI is a more dynamic model than CMM however it still needs to evolve to keep up with the highly dynamic

projects and maintain room for the more stable ones. We have reached a point where innovation becomes one of the major concerns of companies that will need to be constantly adapting and trying to survive to a highly competitive market.

In our research we found industry testimonials of agile implementation and some articles about the usage of CMMI and agile, although for High Maturity agile has proven to be a effective improvement but not a mean of CMMI 5 implementation [1]. We believe CMMI next evolution step will be towards knowledge, innovation and agile (management and development).

## References

1. Sutherland, J., Jakobsen, Johnson, C.R. K.: Scrum and CMMI Level 5: The Magic Potion for Code Warriors. In: Hawaii International Conference on System Sciences, Proceedings of the 41st Annual, pp. 466--475 (2008)
2. Chrissis, M.B., Konrad, M. Shrum, S.: CMMI® Guidelines for Process Integration and Product Improvement. Software Engineering Institute. Addison-Wesley, Massachusetts (2006)
3. Takara, A., Bettin, A.X., Toledo, C.M.T.: Problems and Pitfalls in a CMMI level 3 to level 4 Migration Process. In: Sixth International Conference on the Quality of Information and Communications Technology, pp. 91--99 (2007)
4. Hollenbach, C. Smith, D.: A portrait of a CMMI$^{SM}$ level 4 effort. Systems Engineering, pp. 52--61 (2002)
5. Jones, C.: Applied Software Management: Assuring Productivity and Quality. Software Engineering Series. McGraw-Hill, Inc., New York (1991)
6. Critical Software, S.A.: Annual Report 2006. Critical Software, S.A. (2006)
7. Critical Software, S.A.: Annual Report 2007. Critical Software, S.A. (2007)
8. Goldenson, D.R., Gibson, D.L., Ferguson, R.W.: Why Make the Switch? Evidence about the Benefits of CMMI. In: SEPG. Carnegie Mellon Software Engineering Institute (2004)
9. Radice, R.: Statistical Process Control in Level 4 and Level 5 Software Organizations Worldwide. In Software Technology Conference (2000)
10. Braam, G.J.M.: Performance effects of using the Balanced Scorecard: A note on the Dutch experience. Long Range Planning, vol. 37, issue 4, pp. 335--349 (2004)
11. Schalken, J., Vliet, H. van: Measuring where it matters: Determining starting points for metrics collection. The Journal of Systems and Software, vol. 81, pp. 603--615 (2008)
12. Gonçalves, M.G.S., et al.: A Strategy for Identifying, Classifying and Prioritizing Improvement and Innovation Actions: A CMMI Level 5 and Six Sigma Approach. In: 19th Australian Conference on Software Engineering, pp. 104--111 (2008)
13. Xiaosong, Z., et al.: Process Integration of Six Sigma and CMMI. In: Industrial Informatics, 2008. INDIN 2008. 6th IEEE International Conference, pp. 1650--1653 (2008)
14. Cohan, S., Glazer, H.: An Agile Development Team's Quest for CMMI® Maturity Level 5. In: 2009 Agile Conference, pp. 201--206. IEEE Computer Society (2009)
15. Jakobsen, C.R., Sutherland, J.: Scrum and CMMI - Going from Good to Great. In: 2009 Agile Conference, pp. 333--337. IEEE Computer Society (2009)
16. SEI, Process Maturity Profile - CMMI® For Development SCAMPI$^{SM}$ Class A Appraisal Results 2009 Mid-Year Update. Carnegie Mellon University (2009)
17. Kan, S.H., Metrics and Models in Software Quality Engineering. Addison-Wesley (1995)
18. Park, R.E., Goethert, W.B., Florac, W.A.: Goal-Driven Software Measurement - A Guidebook. Software Engineering Institute (1996)
19. Raisinghani, M.S., et al.: Six Sigma: Concepts, Tools, and Applications. Industrial Management & Data Systems, vol. 105, issue 4, pp. 491--505 (2005)

# Multi-Agent Systems and Robotics

# Evaluating Agent-based Methodologies with Focus on Transportation Systems

Lúcio Sanchez Passos

Department of Informatics Engineering
Artificial Intelligence and Computer Science Laboratory (LIACC)
Faculty of Engineering, University of Porto (FEUP)
Rua Dr. Roberto Frias, S/N • 4200-465 Porto • Portugal
pro09026@fe.up.pt

**Abstract.** To develop a Multi-agent System (MAS) that achieves all requirements of the system, the used agent-based methodology has to address all issues related to it. So, to apply these methodologies in transportation system we have to evaluate then in its context. Take advantage of methodology's features make easier to design the optimal MAS. In this work, we used a MAS Analysis and Design Framework to evaluate three main agent methodologies: Gaia, Tropos and MaSE. A comparative table and methodology combination proposal are the achieved results.

**Keywords:** agent-based methodologies, transportation system

## 1      Introduction

The computer systems in recent decades become facilitators of many tasks performed in various society's fields, and its architecture was originally simple and possessed a small number of capabilities. In addition, with the opening of new alternatives for applicability, as well as its complexity increment, were created an automation function for some parts of the system. Thus, over the years we are searching to overtake higher levels of computational complexity that aims, in an utopico view, to mimic some of the human being capabilities.

In this context, the last challenges for researchers in computing are related to the arrival of distribution in enterprise systems. The prerogative of these systems is that information alone can be combined and used as a new platform of knowledge. New problems, such as systems heterogeneity and communication standards, led the scientific community to seek accommodate these new requirements, bringing a vision of computer system as an organization and so create agents' concepts. However, do not exists a complete definition of agent, but in a few words an agent is persistent computational entity that can perceive, reason, act and communicate [1].

Considering the characteristics of the current transportation system, we can fit it as a distributed and decentralized system, suggesting an enormous application potential of Multi-Agent Systems (MAS) [2]. Further, MAS has extensively applied in

transport systems, since of traffic control [3], vehicle to vehicle Communication [5], and autonomous driving [4]. However, to design these systems several methods are being used, and in some cases the achieved results are different from those expected due to the choice of an inappropriate methodology.

## 1.1 Related work

The ability to choose between different agent methodologies is the first step towards a MAS which integrates all the paradigm's advantages. Evaluate agent-based methodologies received much attention for the variety of existing ones, although most of these approaches were aimed to specifics application. In these approaches, the authors propose delimitation in the subject, which varies from scenarios, system orientation, to a specific methodology.

In [7], the author focuses on analyse different concepts of roles and its encapsulation. This perspective is justifiable because roles can be used as base for MAS, however, these are others interesting aspects to observe. The reference [8], presents a series of issues that need to be checked in a methodology, and it encouraged other authors to propose their own evaluation framework [6], [9]. Instead, reference [10] reports a methodologies analysis with a case study. These are studies evaluating general methodologies as Gaia, Tropos and MaSE and reference [11] is a very complete one, but do not centre in a specific scenario where some features can be more valuable. Roughly, was not found a study that evaluates agent-based methodologies oriented in transportation system despite its large use.

## 1.2 Contributions of this paper

In this work, we will first describe the three main agent-based methodologies: Gaia, Tropos and MaSE. After that, we basically analyze and identify the potential of three methodologies for designing MAS focused on creating systems for the transport scenario. Instead of creating a new framework for evaluation of agent methodologies, we prefer to adapt from Multi-agent System Analysis and Design Framework (MASADF), extracting those characteristics to create a comparative table with important key-points in transportation environment. Many interesting question marks rise in this context, stimulating further discussions and research, some of which are also discussed in this paper.

Following this contextualization of the subject, the remaining of this paper is organised as follows. In the next section will be described an evaluation framework, citing the main points and adaptation. Section 3 is going to present an overview of analysed methodologies. In Section 4, a comparison will be made, followed by a discussion. The conclusion is the last part of this paper.

# 2     Evaluation Framework

Evaluate the strengths and weakness of agent methodologies is an important role to determine which of them are in a mature or described in sufficient details to design real systems. Thus, to compare the methodologies presented in this paper we choose a framework described in [6], because covers the needed areas and it is simple to understand and apply. The framework extends MASADF and is divided in six major areas: concepts, notation, process, pragmatics, support for software engineering and marketability. An overview of above areas is given in sequence.

As the name says, the concepts area criticizes how well the methodology achieves it. This criterion is divided in two parts; first the general concepts include the agent's abilities to operate without supervision, to react to the environment, to pursue new goals, and to collaborate with others agents. Second, the agent's lateral concepts, such as beliefs, desires and intentions, are taken into account. Besides, observe communication, operation and socialization characteristics.

Next, modelling techniques used to represent system's concepts have to be evaluated. To do so, notation criteria are simplicity to understand and use it. Another one is the capability to express and deal with different levels of details in various development phases. Also, preciseness present in the symbols, syntax and semantic definitions, and traceability and modularity abilities.

In process phase, we compare and evaluate the life cycles standards for each methodology, reminding it could be totally different. Firstly, life cycles are studied, then manageability, i.e., project, configuration, verification and quality plans. Specify the development context which the methodology will be used.

Pragmatics criteria deals with practical aspects of deploying and using a methodology. An important part is extant resources, such as development tools and available information. Required resources and project adaptability are treated as evaluation of methodology's usability for users and its versatility in other context.

The software engineering support considers the methodology's ability to execute it. It is divided in operational and maintenance (reusability, testability, extensibility, modifiability, maintainability), and the integration of new terms. Finally, the final criterion, marketability, refers the satisfaction of all people involved in the system's development.

In this paper, we changed only the metric used in the original framework, a crescent scale of 0 to 7. The new qualification levels were proposed to simplify some evaluation aspects, because this approach can bring objectivity to the final results. Instead, we propose scale below, reminding the focus on transportation scenario:

- Low: Indicates that the methodology poorly address the criteria.
- Medium: Indicates that the methodology address the criteria, yet some major issues are lacking
- High: Indicates that the methodology fully or almost fully address the criteria

# 3 Analysed Agents Methodologies

In this section we give an overview of three well-known methodologies: Gaia, Tropos and MaSE. We first introduces Gaia and after that its improvements and extension. In Tropos each stage is presented a short description, such as MaSE. For the methodologies we cited their known advantages and limitation, as seen below.

The current transportation system uses computers to improve efficiency of traffic control and coordination, as well as transport planning. Mainly, centralised decision can be a good approach to traffic coordination; however their performance depends on system scale. Another issue is real-time constrains and the presence of heterogeneous participating entities. Finally, to address a new perspective MAS can be used to give a distributed view of the transportation system [19].

In an ideal application of MAS technology, by Parunak [18], it has to be: modular (entity is well defined), decentralized (application can be decomposed in stand-alone processes), changeable (can modify quickly and frequently), ill-structured (all information is not available), and complex. So, we can see that transportation system is a great application for MAS, because fulfil all requirements above.

## 3.1 Gaia

The first complete agent methodology was Gaia, described in [12], which sees MAS as a computational organization where each role is interpreted by an agent and cooperate with others ones to achieve an application's main goal. Despite, Gaia was quite used, it presents limitations: could only work with closed MAS (agents must always cooperative), and its representation do not apply software engineering standards. Therefore, were proposed two improvements, the Gaia V.2 [14] and ROADMAP [15], besides these are matched with AUML [16] representation to overcome Gaia's poor notation. In this subsection, we will explain all phases of Gaia and also each improvement cited above.

In applying Gaia, we start with abstract concepts and, in each step, models are constructed until concrete concepts that could be implemented to satisfy requirements. The methodology consists in two parts: analyse and design (authors explain that requirements capture is treated as an independent phase). The analysis phase aims to model the system's structure, describing each organisation part to better understand it. Additionally, the design has to transform the results of analysis in a low level abstraction to clarify manners the agents may be implemented.

Entering deeper in the analysis stage in Gaia, the first step is indentify roles in the system, i.e., determine functionalities in a system. Finally, we have a prototypical roles model, delimiting each role including their permission, responsibilities, activities and protocols. Then an interaction model is constructed to capture the communication that occurs in the system between roles. Thus final output from analysis stage is a fully elaborated roles model, which documents all important aspects, including protocols and activities. In Gaia v.2, was added an environmental model very useful for dynamic system, which environment has great influence on roles. Another step was the organizational rules to better define an organizational

structure in the design stage. Instead, a useful improvement seen in ROADMAP was the knowledge model associating roles and knowledge handle by them.

The design process derives from roles, interaction models, agent's types and instances that will be used in the system. Also a services model indentifies the required services to achieve the agent's role, and for last the acquaintance model focuses on agent's communication. Nevertheless, in Gaia v.2 the design stage was quite modified introducing an organizational structure and patterns, which influences aims a more complete role and interaction models. Further, these three new steps contribute for agent and services models. The differences between Gaia v.2 and ROADMAP design stage are the modes to supply Gaia's deficiencies, but are generally similar.

Another improvement in Gaia was related with notation and the reason was explained above. To supply this gap, Gaia's versions use the Agent UML (AUML) to improve its expressiveness to represent MAS. This extension of Gaia was quite accepted by the scientific community because AUML brings the system closer to an implementation level, and is more compact and efficient express communication protocols.

Although Gaia's improvements try to fulfil some lacks, there still others facets that has to be observed. The first limitation is capturing properly the dynamic aspects of the environment, which are essential for complex MAS systems. Also, the organization structures are modelled using a simple notation what difficult construction for large systems. The final limitation refers to requirement engineering, which is not treated in Gaia such as in other methodologies.

### 3.2 Tropos

To supply the need covering all the phases of software engineering was created Tropos [13]. The methodology uses agent and related concepts applied in stages of system development to better fulfil the specification. Besides, the early requirements analysis is a crucial Tropos part to capture all goals. Tropos adopts Eric Yu's i* model [17] that provides notions of actors, goals and actors dependencies to model an application. An overview of Tropos' four development phases is presented below.

The early requirements analysis is a goal-oriented task and to do it we have to model stakeholders' intentions as goals. From these goals, we can guide the analysis into system requirements. The output of this step is a strategic dependency model for describing a relationship network among actors, i.e., a graph constructed to explain what kind of dependencies an actors have, called dependum. The depedum can be goal (fulfil a precisely goal), softgoal (fulfil a subjective goal), task (perform a given activity), and resource (provide a resource). So, these elements produce an organizational model.

In the next phase, late requirements analysis, a complete operational description of the system-to-be is desired. In other words, the system is added in an organizational model as a centre piece, and expresses its internal and external relations. With this model, can be visualised mains tasks and goals that a system has to achieve, however

cannot be given a complete system's characterization because is missing architecture description, made in the next stage.

In sequence, an architectural design is defined to express system's architecture, which presents how system's components work together. Tropos has available a group of predefined organizational architectural styles, based on research in organization management field, to help match MAS in organization context. To execute this stage, we have to search the best architectural style and adapt it to the system organization wanted using analytical procedures.

Finally, the detailed design can be made to introduce complementary information in each architectural component. In particular, this methodology phase is concerned with how each assigned goal is achieved, and an important manner is represents social patterns between agents. The process to identify social patterns began focusing in an actor and decomposes its depedums in new relationships, adding involved agents, subgoals and subtasks. A specification of agent communication and behaviour is made in detail design and some extensions are proposed by the author, e.g., the use AUML to notation and FIPA-ACL as language communication.

Tropos is a goal-oriented methodology and uses actor and its goals in all development steps. A strong Tropos' feature is the requirements analysis which helps to group important information about system needs. In other hand, Tropos model the environment but not gives methods to design an actor who could be dynamic interactive with it. In addition, BDI applications are Tropos' mainly focus, limitating its use to others agents' architectures.

### 3.3 MaSE

The Multiagent System Engineering (MaSE) [20] is a full-lifecycle methodology that tries to suppress all needs in develop a heterogeneous MAS. It derives from object-oriented paradigm and treat agent as a specialized objects with capabilities to interact to achieve an individual or collective goal. MaSE explicit covers analysis and design, as others methodologies, however specifies a logical order to construct each model. Further, we give an overview of MaSE explaining the steps to design a MAS.

Starting with analysis phase in MaSE, a software requirement specification should be done pointing system's behaviours that have to be fulfilled, but MaSE do not proposes any guide for this stage. Next, in Capturing Goals, we must be identified from extracted scenarios, presented in the initial specification, the goals. These goals are structured in a Goal Hierarchy Diagram (GHD), which is a graph with non-defined hierarchy, where nodes represent goals and arcs define a sub-goal relationship. In GHD, goal can have types to easily manage and understand it.

The second analysis step is Applying Use case that aims translate goals in roles and associated tasks, so we can define desired system actions. Sequences Diagrams aims present an events' flow between roles, including agents' communication. Thus, from the results of last two steps (GHD and Sequences Diagrams), we can explicit the refining roles. This stage form the foundation of design phase, because depict system's base, the roles. Summarizing, refining roles is a great role model and define role behaviour using concurrent task model.

Designing a system with MaSE has to pass through four steps: creating agent classes, constructing conversations, assembling agent classes and system design. Firstly, to create agent classes each role is assigned to a class and can be reallocated to reach the best design. Others issues, such as various levels of cohesion, are taken into account. Identify conversations between agents classes are also treated in this step, which represent the system using an agent class diagram. Entering deeper in communication, next step is constructing conversations that were identified in last phase. A conversation is composed by an initiator, a responder and a coordination protocol. So, the communication class diagram for each agent in the conversation, represent state diagram with conditions.

In assembling agents, internal features of agent class are modelled, divided in agents and architecture's component, both flexible to the designer will. The agent architectural diagram shows components, their visibility and external resources. Finally, the last and simplest step of MaSE is the system design where deployment diagrams are used to present every agent instance.

Therefore, MaSE is an interactive methodology that enables the user modifies any system's part and easily perpetuates this new feature over all stages. On the other hand, there are no references to environment model, i.e., in MaSE the ambience is considered static which is not true for almost MAS. Additionally, like Gaia, MaSE do not suggest any path to guide the user in system's requirements.

## 4  Comparative Evaluation and Discussion

For comparison, we use the framework presented in section 2. In order to better understand the evaluation, the above methodologies will be discussed in terms of its features for each criterion. After, we will comment all frameworks, and a comparative table is going clarify the final grades. At the end, we generally discuss the obtained results with focus on transportation system. Before begin the evaluation has to be made an observation that here will be considered Gaia v.2 and all methodologies are extended with AUML.

### 4.1  Concepts

The agent's general concepts are presented in Gaia and the agent can be fully specified with autonomous, proactiveness, however integration model not cover all social connections. Another limitation is reactiveness, for static environments the methods are enough, but Gaia cannot model actions for dynamic scenarios. In lateral concepts, Gaia covers all communication, operation and socialization features, but not defines any mental state for the agent.

Tropos also has limitation on reactiveness, because we can only design MAS for static scenarios. Also, proactiveness problem is related with the hierarchical relationship. Instead, in agent's lateral concepts Tropos cover all features. Thus, MaSE presents limitations in the same parts for the same reason and mental attitudes can be seen as defined goals.

## 4.2 Notation

The Gaia notation is simple, so is easy to be understand and used, but is incomplete to define and deals with more complex relations and level of description. Nevertheless, with AUML extension a larger numbers of notations were available, but still there is dependency verification.

Tropos and MaSE have originals notation more complete than Gaia's and were improved with AUML use. These two methodologies have advantages in traceability and modularity compared to Gaia, making possible an easier implementation.

## 4.3 Process

The life cycle could be mapped with Gaia, but only two phases are referenced: analysis and design. None of those plans cited in section 2 are provided by Gaia. Instead, it was a large development context, i.e., can be applied to create, reengineering systems.

MaSE reference only three software engineering phases: analysis, design and implementation. It also provides large development context, plus is an interactive design process. In other hand, Tropos is most complete methodology in terms of lifecycle, missing only test stage and share the MaSE development context features.

## 4.4 Pragmatics

In terms of pragmatics, Gaia is the worst because it has not extensive documentation, and requires a solid background. The adaptability is low and could only be applied a certain environment. One thing that can be highlighted is its non-limitation in respect of platform language.

Tropos was a lot of publication and tools for animation and model checking. It is a naturally understandable, but is BDI-based, limitating the domain application and language suitability. MaSE requires background knowledge but was a large community and documentation. It shares openness in implementation with Gaia.

## 4.5 Support for Software Engineering

All methodologies are very similar in this criteria, they permit some reusability, extensibility and modifiability. Nevertheless, MaSE and Tropos was the advantage to possess the testability.

## 4.6 Marketability

This item is hard to be evaluated, because is a subject statement. From the end-user's viewpoint any methodologies has a total safety, because agent paradigm is new. In developer side, it satisfaction is deeply related to facility the design and find information. In this direction, Tropos is better because of pragmatics issues, followed

by MaSE. Another role in this manager satisfaction, but none of them provides any management plans.

Thus, the comparative table is present below, as Table 1

**Table 1.** Comparation table between the analised methodologies

| Evaluating Criteria | Gaia | Tropos | MaSE |
| --- | --- | --- | --- |
| Concepts | High | High | High |
| Notation | High | High | High |
| Process | Medium | High | Medium |
| Pragmatics | Low | Medium | Medium |
| Support for Software Engineering | Low | Medium | Medium |
| Marketability | Low | Medium | Low |

Focusing in the transportation scenario, there some observation that have to be done to obtain better results. Any of theses methodologies fit perfectly for transpostation environment, because of the high dynamism. Thus, there is a lack in notation of organization and its relations. From this evaluating, we can see that exists enormous lack in agent methodologies, if they are observed in a "big picture" view.

Working with three methodologies presented, we can propose a solution based in the complexity of the desired system. For small systems, MaSE is the best methodologies because is quite flexible, covers almost all software design and implementation, if it is used within agentTool. Nevertheless, for more complex system an unique solution could not be achieve, so we bring a solution used in other works, such as [10], which is to use early requiments phase to characterise the system's needs and to next phases Gaia v.2 extended with AUML to provide flexibity in choose agent's architecture. Nevertheless, this perspective could be improved by the knowlegde model described in MaSE, because in transportation, agents has to deals with many kinds of information.

## 5    Conclusion and Future Work

In this paper was presented the three well-known agent-oriented methodologies where every specification phase was describable. An overview of the evaluation framework was done. Finally, a complete evaluation and comparison between Gaia, Tropos and MaSE was made. Additionally, the work performed here provides a starter point to researcher in the field to guides their efforts.

In future works, we point out some directions. First, a deeper study has to be done to concrete results described here, because agent-oriented systems are nowadays subject and new perspectives are always emerging. Study specific methodologies can be great to bring new features, which can be grouped in a base methodology, e.g., ADELFE propose the use of AMAS theory to solve environment's dynamic issue.

Nevertheless, the most hard-working path is to create a whole new methodology that covers all issues related with transportation system.

# References

1. Huhns, M. N., and Singh, M. P.: Agents and Multiagent Systems: Themes, Approaches, and Challenges. In: Reading in Agents, pp. 1-23 (1998)
2. Schleiffer, R.: Intelligent agents in traffic and transportation. In: Transportation Research 10C, no.5-6 (2002)
3. Oliveira, Denise de and Bazzan, Ana L. C. and Silva, Bruno C. da and Basso, E. W. and Nunes, L. and Rossetti, R. J. F. and Oliveira, E. C. and Silva, R. and Lamb, L. C.: Reinforcement learning based control of traffic lights in non-stationary environments: a case study in a microscopic simulator. In: 4th European Workshop on Multi-Agent Systems, pp. 31-42 (2006)
4. Vasirani, M and Ossowski, S.: Evaluating Policies for Reservation-Based Intersection Control. In: 14th Portuguese Conference on Artificial Intelligence, (2009)
5. Gonçalves, J.F.B., Esteves, E.F., Rossetti, R.J.F., Oliveira, E.C.: Simulating Communication in a Service-Oriented Architecture for V2V Networks. In: 14th Portuguese Conference on Artificial Intelligence, (2009)
6. Akbari, Z.O., Farrahi, A.: Evaluation Framework for Agent-Oriented Methodologies. In: Journal of World Academy of Science, Engineering and Technology, no. 45, pp. 418-423 (2008)
7. Puviani, M., Cabri, G., Leonardi, L.: Agent roles: From methodologies to infrastructures. In: International Symposium on Collaborative Technologies and Systems, pp. 416-423 (2008)
8. Yu, E., Cysneiros, L.M.: Agent-Oriented Methodologies – Towards A Challenge Exemplar. In: 4th International Bi-Conference Workshop on Agent-Oriented Information System, pp.1-13 (2002)
9. Abdelaziz, T., Elammari, M., Unland, R.: A Framework for the Evaluation of Agent-Oriented Methodologies. In: 4th International Conference on Innovations in Information Technology, pp. 491-495 (2007)
10. Castro, A. and Oliveira, E.: The rationale behind the development of an airline operations control centre using Gaia-based methodology. In: Int. J. Agent-Oriented Software Engineering, vol. 2, no. 3, pp.350–377 (2008)
11. Bergenti, F., Gleizes, M.P. and Zambonelli, F.: Methodologies and Software Engineering for Agent Systems: The Agent-Oriented Software Engineering Handbook, Kluwer Academic Press, ISBN: 1402080573.
12. Wooldridge, M. J., Jennings, N. R., and Kinny, D.: The Gaia Methodology for Agent-Oriented Analysis and Design. In: International Journal of Autonomous Agents and Multi Agent Systems, pp. 285-312 (2000)
13. Castro, J., Kolp, M., and Mylopoulos, J.: Towards Requirements-Driven Information Systems Engineering: The Tropos Project. In: Information Systems. Elsevier (2002)
14. Wooldridge, M. J., Jennings, N. R., and Zambonelli, F.: Developing Multiagent Systems: The Gaia Methodology. In: ACM Transactions on Software Engineering and Methodology (TOSEM), vol. 12, no. 3, pg 317-370 (2003)
15. Juan, T., Pearce, A., Sterling, L.: ROADMAP: Extending the Gaia Methodology for Complex Open Systems. In: Autonomous Agents and Multi-Agent Systems (2002)
16. Bauer, B., Müller, P.J., Odell, J.: Agent UML: A Formalism for Specifying Multiagent Interaction. In: Springer-Verlag, Berlin, pp.91-103 (2001)

17. Yu, E.: Social Modelling and i*. In: Conceptual Modelling: Foundations and Applications, LNCS vol. 5600, pp. 99-121 (2009)
18. Parunak, H.V.D.: Industrial and Practical Applications of DAI. In: Multiagent System (AGENTS'01), MIT Press (1999)
19. Passos, L.S.; Rossetti, R.: Intelligent Transportation Systems: a ubiquitous perspective. 14th Portuguese Conference on Artificial Intelligence, (2009)
20. DeLoach, S.A.: Multiagent Systems Engineering: A Methodology and Language for Designing Agent Systems. In: Agent-Oriented Information System 99 (AOIS'99), pp. 45-57 (1999)

# Interaction Protocols for Electronic Institutions in B2B inter-organizational business

Pedro Brandão Neto

Faculdade de Engenharia da Universidade de Porto – DEI-LIACC
Rua Dr. Roberto Frias, 420-465 Porto, Portugal
pro09012@fe.up.pt

**Abstract.** Electronic Institutions came up because of the increasing need of business. The Electronic Institutions, defined through the multi-agent paradigm, constitutes the virtual marketplace. The agents exchange messages to reach their goals over negotiation protocols. In order to improve the interactions between participants in B2B inter-organizational business are proposed in this paper different versions of interaction protocols for Electronic Institution, based on FIPA-Contract-Net Specification. The negotiation protocols based in this specification improve the participant agents negotiation of the proposed Electronic Institutions. To analise the interactions between the Electronic Institution agents using the developed protocols is defined a multi-agent system, where agents got to fulfill its tasks satisfactorily in the business satisfactorily, using communicative acts.

**Keywords:** multi-agent; contract net protocol; electronic institutions.

## 1 Introduction

Platforms for Electronic Institutions (EI) emerged due to increasing need to business [1] [2]. One of the essential steps in the B2B business activity is the selection of business partners that is naturally accomplished through the negotiation. Nowadays, there are EIs that standardize and define the steps that are considered fundamentals in the negotiation process.

The EIs constitute Virtual Marketplaces (VM) that institute norms and rules of market operations as well make available useful services, and this abstraction is defined in this work through the Multi-Agent System (MAS) paradigm. The participant agents of a MAS in the negotiation process need to share a common language and to establish communication protocols in order to provide a conversation structure.

Therefore, the agents in a MAS exchange messages to reach their goals. Negotiation protocols can be used in this step. Several situations between agent conversation are led by adoption of Interaction Protocols (IP) standards [3] [4] [5]. The standards define the screenplay in the conversation between agents in the negotiation process. These standards are specified by FIPA (Foundations of Intelligent Physical Agents) [6] and implemented by the JADE platform [7], reason why will be this platform chosen for accomplishment of this work. And the

programmers can redefine these protocols including the needed logic according to the specific application domain.

In order to improvement regarding interaction standardization between participant entities in a VM, we will define an Electronic platform of negotiation based on FIPA-Contract-Net specification [6] intending to optimize the process of B2B inter-organizational negotiation.

However, there is not a global consensus at the moment to adoption of a protocol in the conversation between agents in the negotiation process in EI, and several researchers [1] [3] [4] [8] [2] exploit this ambit aiming to delineate new models of conversation. This paper investigates the FIPA Contract-Net specification in order to develop means more efficient in negotiation process between agents.

In the JADE environment when is utilized the FIPA-Contract-Net protocol in the conversation between agents, is distinguished the Initiator role (agent that initiates the conversation) and the Responder role (agent inserted in the conversation after be contacted by some other agent). Therefore, based on the FIPA-Contract-Net (ContractNetInitiator and ContractNetResponder classes that belong to jade.proto package) protocol we propose three protocols, in both interaction roles, that are: open protocol, proposal protocol and closing protocol.

In this work is proposed a MAS in implementation of an Electronic Platform of Negotiation. This MAS is applied in VM that groups interested enterprises (i.e. agents) in offer and them that want to buy. In this VM the defined protocols (open protocol, proposal protocol and closing protocol) are used by Buyer and Supplier agents. Once established, that Electronic Platform can be used for future researches.

This study is the basis for the implementation of a negotiation algorithm more powerful, the Q-negotiation [1] algorithm, which introduces advanced features in electronic markets negotiation, such as multiple-attribute negotiation, learning in negotiation and distributed dependencies resolution.

The problematic presented in this work excites questions such the proposed agents will deliberate, reason and decision making in a dynamic environment, in order to finalize a B2B transaction in a virtual marketplace

Thus, the intention of this paper is to answer the following hypothesis: do the negotiation protocols based on FIPA-Contract-Net specification improve negotiation of participant agents of the Electronic Institution?

## 2 Literature Revision on Electronic Institution

Due to internet explosion, the electronic transactions became themselves an unavoidable trend. Electronic commerce represents a significant piece of the global market and it merits the investments made in researches in that new technology. It is an area that requires attention and efforts of several researchers. Several Artificial Intelligence (AI) techniques are applied in the area of electronic commerce, of among which stands out: Agent technology [9] and machine Learning Techniques [10] [11] [12].

Recently, the growth of Information Technology and Communication (ICT) is transforming the way as the traditional commerce is done. With the advance and

expansion of these new technologies is eliminated time and space restrictions. In this scenario appears the Electronic Commerce, which is a new commerce modality based on ICT and comprehends the following fields: business-to-consumer (B2C) and business-to-business (B2B) electronic commerce. In the B2C electronic commerce, business participants are sellers and final buyers. Already in the B2B electronic commerce, as opposed to B2C, the objective of the commercial transaction is not a final product and, generally, business participants are enterprises that need to include in their own products processes that are outside of their expertise domain or resources that do not own [13].

The development of new systems, through new paradigms, contributed significantly in the areas of business in the B2B modality, as appointed in [14] [15] [16]. So, it is possible to reduce negotiation costs, marketing and contract formation, among others. Like this, the enterprises always aim to improve and increase their capacity to interact in a more efficient way with other organizations minimizing their costs and maximizing their profits.

In that new era of the digital economy, the commercial relationships between business partners became themselves more flexible and the businesses tend to be created whenever that appears business opportunity. The emergent need of new products and services with increase in quality, shorten time to market and the instability in demand of the product. This new reality brought new technological, social, economical and ethical challenges [17].

In [18] emphasizes the fact that the development of electronic commerce resulted in a considerable increase to competitiveness causing to need of new organizational models. Already in [13] refers himself that the competitiveness of electronic commerce appears due to market openness and of their dynamism, permitting the emerging of new organizational structures, as is the case of Virtual Enterprises.

The Virtual Enterprise (VE) can answer effectively to the market requirement in new demand, because it combines the essential competencies of the independent and heterogeneous enterprises that collaborate in a temporary and loosely linked network, thereby presenting high flexibility and agility. The enterprises try to answer these new market requirements by engaging themselves in temporary corporations presenting a flexible structure that changes dynamically according to current market situations. All the enterprises collaborate for a global goal with their competencies, knowledge and resources ([13]).

Recently, came up Electronic Institutions (EI) concept, proposed in [8] [1] [2] [14], as virtual counterpart of real world institutions. According to [13], EI is a platform that enables, through a communication network, the automatic transactions between parties involve in a business activity, establishing sets of institutional norms and rules of behavior. In this way, the EI ensures the needed trust in any electronic transaction.

The EI has two main purposes: (1) to support interaction between representative artificial agents of business activity participants (enterprises) as a coordination framework, making the business agreement establishment more efficient, and (2) provide a trust level through a normative environment ([19]). These environments are created to establish some kind of social order [20] that allows successful well interactions between the heterogeneous and autonomous entities [21]. Therefore, researches as [22] [23] [24] [2] [1] [25] [26] demonstrated the relevance of EI and application of several technologies, which aim improvement in the EIs.

Some approaches include Multi-Agent System (MAS) paradigm that aims to automate the creation and operation process of VE [13]. The MASs intend to provide the best mean to find out solutions to problems which are not easily resolved by a single agent [27]. There is an emphasis towards to develop mechanisms that regulate MAS in distributed environments. Thus, whilst agent theory describes agents as autonomous entities of self interests, rather that interact in openness environments, the MAS application in real-world scenarios raises an important question: how to ensure a cooperative behavior in scenarios populated with heterogeneous agents and self interests ([18])?

Efforts in interaction infrastructure for MAS have been extensively studied among the participant agents in an EI because them require and involve specific interactions during a determined period of time.

Being the entity agents of heterogeneous and autonomous software, there is the need to define negotiation protocols between the agents in order to select participants that will be able to make the ideal business according to their own objectives ([13]).

Such as happen in the real market, the agents of virtual market engages themselves in an iterative negotiation process composed of multiple rounds of proposals and counterproposals. That iterative process of negotiation enables the learning during the negotiation process. This learning is realized through the Q-negotiation algorithm that includes observation (learns with actions already known) capabilities and exploitation (learns with actions don't unknown) [1].

In this work, we start up of the study of FIPA-Contract-Net protocol for the development of negotiation protocols, which are applied in the proposed Electronic Platform.

The FIPA-contract-net specification came up with the original Contract-Net interaction protocol pattern that it adds rejection and confirmation communicative acts [28]. In this interaction protocol, the agents of a MAS can be divided into categories according with their roles. In the contract net IP, one agent (the Initiator) takes the role of manager which wishes to have some task performed by one or more other agents (the Participants) and further wishes to optimize the function that characterizes the task.

Some works [1][13][29] have been used to provide means more elaborated in the negotiation effectuation in partner selection of the VE. Nevertheless, even with the progress obtained with application of those techniques from Artificial Intelligence (AI), specifically the software agent application, the negotiation process, still needs of improvement.

This paper will utilize FIPA specification protocols available in the JADE (Java Agent Development Framework) platform.

## 3  WORK PROPOSAL

This section presents a Multi-Agent system that was developed to represent an Electronic Platform, which will operate in the textile virtual market where buyer and supplier agents convene themselves constituting a virtual marketplace. The Buyer and Supplier agents represent the enterprises that wish to buy and sell, respectively. That

agents meet themselves in the VM in the expectancy of fulfill business, such as happen in traditional commerce, where humans meet themselves  personally and exchange proposals and counterproposals attempting to reach agreement in an business electronic transaction. Like this, the Buyer and Supplier agents engage themselves in an iterative negotiation process composed of multiple rounds of proposals and counterproposals.

In an Electronic Platform, the participant agents to communicate themselves needs to share a common language and to establish communication protocols to provide a common structure to the negotiation process. Each agent will can, also, retains a decision making mechanism that it will enables to determine the criterions for an agreement. In this VM, three protocols were defined, which are: open protocol, proposal protocol and closing protocol. Therefore, is important to perceive that these three protocols were defined on the basis of FIPA-Contract-Net Protocol.

In synthesis, the developed Electronic Platform is composed of two agent types: Suppliers and Buyers; the "Buyers" interested in buying a given product make their announcement; the "Suppliers" await for proposal requisition of the Buyer agents and it send proposal to each Buyer agent that it requested; the Buyer Agent selects the best proposal and then divulges the winner, and the negotiation is finished.

In the following subsections will be described as Buyer and Supplier agents are organized in the Electronic Platform and how they accomplish your tasks; as are defined the three protocols to support the negotiation between the agents and, by last, a scenario is detailed.


### 3.1 Multi-Agent System (MAS)

The proposed Electronic Platform, in this paper, offers the needed infrastructure for the multi-agent system's agents implemented on the JADE (Java Agent Development Framework) [7] platform fulfill business in a VM. The two developed main agents are: Supplier and Buyer. The other agent in development phase is the proxy agent.

The Buyer agent, which represents the enterprises that wish to buy, has the following tasks:
- Subscribes in a VM;
- Finds the Supplier agents that met themselves in the VM, only those that able to satisfy your needs.
- Requests proposals from all Supplier agents;
- Analises the received proposals from Supplier agents;
- Chooses the best proposal;
- And, to send a message toward supplier agent that offered the best proposal;

The Supplier agent represents the enterprises that wishes to sell, your tasks are:
- Subscribes itself in a VM;
- Registers your services in the yellow pages;
- Has a data structure to ensures your product catalogue;
- Waits by proposal requisition from Buyer agents;
- Formulates proposals for a Buyer agent;

- Announces and to officialize the Buyer agent, in toward to accept the conditions and compromises itself to realize the agreement;

The proxy agent will make easy the communication process between MAS and the external interface. The role of this agent is to offer an online environment for enterprises represented by Buyer and Supplier agents. Therefore, the need of this agent succeeded to interconnect the real enterprises with MAS. This will enable transparent communication process between enterprises and MAS.

The negotiation between these agents is done using the available behaviors in the JADE platform and the Interaction Protocols according to the FIPA (Foundation for Intelligent Physical Agents) specification, in case FIPA-Contract-Net. TTable 1 below shows the Buyer and Supplier agent's behaviors associated with yours respective tasks.

**Table 1.** Behaviours list associated to the Buyer and Supplier agents.

| Agents | Behaviors | Tasks |
|--------|-----------|-------|
| Buyer | WakeBehaviour | Start up the agent execution and, also, it queries DF (Directory Facilitador) agent. |
| | OneShotBehaviour | Send a message toward Supplier agents. |
| | Generic Behaviour | Stay waiting by responses of Supplier agents. |
| | Generic Behaviour | Determine the winner, the agent that offers the best proposal. |
| Supplier | OneShotBehaviour | Start up the agent execution and register itself in the yellow pages services (the DF agent). |
| | CyclicBehaviour | Analise the CFP (Call For Proposals) from Buyer agent. |
| | CyclicBehaviour | Verify himself whether won the dispute and compromises itself with agreement. |

### 3.2 Negotiation protocol

In SMA, the agents exchange messages during the necessary interaction and cooperation to achieve their goals. The negotiation protocol can be used. Several common situations can be conducted adopting interaction protocol patterns. Such protocols are specified by FIPA and they are available in JADE. The programmers can redefine these protocols including the necessary logic according to specific domain.

In the JADE environment, when using the FIPA-Contract-Net protocol in the conversation between agents, is distinguished the Initiator role (agent that initiates the conversation) and the Responder role (agent inserted in the conversation after had been contacted by some other agent). Therefore, based on the FIPA-Contract-Net (ContractNetInitiator and ContractNetResponder classes that belong to jade.proto package) protocol we propose three protocols, in both interaction roles, that are: open protocol, proposal protocol and closing protocol.

Open protocol:

The Buyer agent makes the announcement for the Supplier agents that are both in the VM. The Buyer agent specifies the task and also, any conditions that it, the

initiator, places on the task. The performative act utilized is the *CFP (Call for Proposal)* that starts the negotiation between the Buyer and Supplier Agents. Moreover, it can establishe a deadline to determine a period to receive proposals.

Proposal protocol:

1) the Supplier (the participant)  agent may accept or reject the Buyer Agent condition. If the condition is accepted, the Supplier agent formulates his bid establishing their conditions for fulfillment of the task, the performative act to be utilized is the *proposal*, otherwise, it is the *refuse*, 2) the Buyer agent receives all the proposals from Supplier agents with the *proposal* performative act. The Buyer Agent can accept or reject the proposal, with an *accept-proposal* or *refuse-proposal* performative act.

Closing Protocol:

Once that Buyer agent determined the best proposal, the Supplier agent compromises itself with the fulfillment of the task and it indicates the *informe-done* performative act.

The protocols explained above base himself on version of the FIPA-Contract-Net. This interaction protocol is shown in Figure 1.



**Fig. 1.** FIPA-Contract-Net Interaction Protocol [6]

### 3.3 Scenario

A Virtual Marketplace (VM) starts with clusters of the Supplier Agents. Each Supplier agent in the beginning performs his *OneShotBehaviour* behavior. This behavior is responsible to register the agent in the services of yellow pages. Also, this behavior creates two objects utilized by this agent, which are: Fabric and Catalog. The "Fabric" object is a stock template of the Supplier agent and it describes the product characteristics. Already the "catalog" is utilized by the agent at time of execution for load the values of their stock. The next behavior to be performed by the supplier agent is the *CyclicBehaviour*. This behavior remains to wait messages from the Buyer agent and while to arrive a message with the *CFP performative*, the Supplier agent takes an object of the message, which specifies the negotiation conditions. In the case that this agent can meet the demands of the Buyer agent using a *proposal* will be done using the *proposal performative*. In the other way, it will refuse the requisition (the CFP) of the Buyer sending them a message with the *refuse performative*. The last behavior that is performed of this agent is the *CyclicBehaviour*, which is fulfilled when one of the Supplier agents wins the dispute. That is, this behavior is performed when the agent receives a message from the Buyer agent with the *accept-proposal* or *refuse-proposal performative*. If the *performative* is of the *accept-proposal* type the supplier agent updates its stock and compromises himself to fulfill the task with the *inform-done performative*.

Whenever, a Buyer Agent can participate of the VM since that it has available Supplier agents. It starts their execution through the *WakerBehaviour* behavior that realizes a supplier agents search when it sends a message to DF agent (agent that provides the services of yellow pages) and creates an object of the "TargeProduct" type. In this object is described the initial desire of the agent. Following, the Buyer agent starts the negotiation using *OneShotBehaviour* behavior that sends a message to all Supplier agents of the VM. This message is of the *CFP (Call for Proposal)* type and in the same goes an object of the "TargeProduct", which describes the desire of the agent. The other behavior to be executed is a *generic behavior* that waits for responses, as messages, from agents suppliers. The message can be of two types represented by *proposal or refuses* performative acts. By last, the Buyer agent performs other *generic behavior* that fulfills two tasks. The first determines the best proposal received from the Buyer agents and the last sends a message to winning agent with the *accept-proposal performative*, the other agents receive *refuse-proposal*.

Thus, to have a detailed scenery of the functioning of proposed Electronic Platform in this work will be described to follow the negotiation between 3 Supplier agents (s1, s2 and s3) and 4 Buyer Agents (b1, b2, b3 and b4). Every interaction, in the negotiation, of these agents is through the open, proposal and closing protocols included in the FIPA-Contract-Net negotiation. Note that these protocols are encapsulated through behaviors of the table 1 above and are utilized the performative acts.

The Supplier agents assume the manufacturer role of textiles. All Supplier agents get your configuration through a set of values that are attributed randomly in creation time. Then, each new supplier agent will have available 4 fabric types presented in

Table 2 (1st column) that are associated randomly with the quantities and prices in the table (2nd and 3rd column).

**Table 2.** Type of fabrics, quantities and prices that can be transacted in the virtual marketplace

| Type of Fabric | Quantities (meters) | Price (100 meters) |
|---|---|---|
| Chiffon | 500,600,680,700,750,780,800,820, 900 | 1630,1633,1638,1640,1642,1648,1655,1559, 1675,1680,1684,1685,1700,1720, 1780 |
| Cotton | 1450,1490,1500,1580,1620,1648, 1698,1720,1840 | 1080,1190,1205,1215,1280,1298,1300,1320, 1397,1415,1438,1447,1452,1498, 1500 |
| Voile | 2510,2580,2595,2699,2730,2860, 2900,3115, 3120 | 1005,1020,1080,1088,1096,1100,1120,1136, 1158,1170,1178,1180,1500,1503, 1555 |
| Nylon | 900,908,915,930,980,999,1102, 1108,1250 | 1680,1982,1983,2010,2012,2017,2222,2230, 2238,2242,2246,2250,2352,2469, 2680 |

The Supplier agents, s1, s2 and s3, provide fabrics and the Buyer agents, b1, b2, b3 and b4 are the enterprises that want to buy the following fabrics: cotton, chiffon, voile and nylon respectively. To demonstrate how the agents behave themselves in the application operation, we will describe this steps:

1. The agents s1, s2 and s3 meet themselves to constitute the VM;
   - Each agent offers its products (see table 2 above) with the respective available quantity, that was randomly attributed;
   - The agents keep updated the product's quantities;
2. The Buyer agents make a CFP for Supplier agents through the *open protocol* aiming to receive proposals;
3. The Supplier agents elaborate proposals through the *proposal protocol* and send them to the Buyer agent;
4. The Buyer agent selects the best proposal received taking into account several situations, not only the best price using the *proposal protocol*;
5. Following the Supplier agent, the contracted, firms an agreement with the Buyer through the *closing protocol*;
6. And, the transaction is finished;

The Figure 2 below shows the exchanges of messages between the supplier agents (s1, s2, s3 and s4) and buyer agents (b1, b2, b3 b4) in the negotiation.

**Fig. 2.** Exchange of messages between the Supplier and Buyer agents in the negotiation.

# 4 Conclusion

Looking for improvements in the interaction protocols between the participating agents in the inter-organizational B2B negotiation process we defined a SMA to ensure a Platform Electronic. This paper investigates the protocol FIPA-Contract-Net in order to develop more efficient negotiation process. Therefore, we developed a Platform Electronic, the textile industry, for the negotiation between the participating agents. These agents are in the MV to make business getting involved in an iterative process of negotiation through the three protocols specified (opening and closing proposal) that constitute the three stages of the negotiation protocol FIPA-Contract-Net. The two main agents are developed: Suppliers and Buyers. The other agent that is under development is the proxy agent.

Therefore, a detailed overview of operation of how the agents perform their tasks was described. To this end, we used 3 agents Suppliers and 4 agents Buyers. All interaction between the agents was done through open protocols, of propose and closures, which were developed in this work. These protocols have been encapsulated in them as behaviors and were used performative acts established by FIPA.

The purpose of this article resulted in the following contributions: 1) the three protocols of negotiation for B2B transactions, which was based specification from the FIPA-Contract-Net. 2) The definition of a proxy agent. This agent let the process transparent communication between real companies with company officials. Thus,

companies will have a very friendly interface which will facilitate communication with the multi-agent system.

As future work will be implemented the proxy agent from the specification proposed in this paper. This study is the basis for the implementation of a trading algorithm more powerful, the trading algorithm-Q [1], introducing advanced features in trading on electronic markets, such as negotiation of multiple attributes, learning in negotiation and resolution dependencies in a distributed way.

# References

1. Rocha, A. P.: Metodologias de Negociação em Sistemas Multi-Agentes para Empresas Virtuais. Engenharia Electrotécnica e de Computadores, Faculdade de Engenharia, Universidade do Porto, Porto, (2001)
2. Sierra, C.: Agent-Mediated Electronic Commerce. In: Autonomous Agents and Multi-Agent Sytems, 9, pp. 285--301, Spain (2004)
3. Karacapilidis, N., Moraïtis, P.: Decision Support Systems archive. Decision-making and E-commerce systems. 32, 53--69 (2001)
4. Chen, Y. C., Huan, Q., Wang, S. S.: Multi-agent Pursuit-Evasion Algorithm Based on Contract Net Interaction Protocol. Advances in Natural Computation. 3612, 482--489 (2005)
5. Fan, G., Huang, H., Jin, S.: An Extended the Contract Net Protocol Based on the Personal Assistant. In: ISECS International Colloquium on Computing, Communication, Control, and Management, vol. 2, pp.603-607 (2008)
6. FIPA, http://www.fipa.org/specs/
7. JADE, http://jade.tilab.com/doc/tutorials/
8. Rocha, A. P., Oliveira, E.: Electronic Institutions as a framework for Agents' Negotiation and mutual Commitment. In: Artificial Intelligence: Knowledge Extraction, Multi-agent Systems, Logic Programming, and Constraint Solving, LNAI 2258, Springer, pp. 232-245, Brazil, (2001)
9. Sierra, J. A., Aguilar, R., Vigil, P. B. N., Arcos-Rosell, J.L., M.V.: Engineering Multi-Agent Systems as Electronic Institutions. A World of Agents. Vol. V, Nº.4, august (2004)
10. Jin, Y.: User Heterogeneity and its impact on Electronic Auction Market Design: An Empirical Exploration. MIS Quarterly. 28, 21--43 (2004)
11. Gao, S., Xu, D.: Conceptual modeling and development of an intelligent agent-assisted decision support system for anti-money laundering. Expert Systems with Applications. 36, n.2, 1493--1504 (2009)
12. Lee, H. J., Lee, J. K.: An effective customization procedure with configurable standard models: Decision Support Systems. 41 n.1, 262--278 (2005)
13. Rocha, A. P., Cardoso, H. L., Oliveira, E.: Contributions to an Electronic Institution supporting Virtual Enterprises life cycle. In: Virtual Enterprise Integration: Technological and Organizational Perspectives, Idea Group Inc., ISBN 1-59140-406-1, Ch. XI, pp. 229-246, (2005)
14. Esteva, M., Cruz, D., Sierra, C.: ISLANDER: An Electronic Institutions editor. In: Proceedings of the first international joint conference on Autonomous agents and multiagent systems, Italy (2002)
15. Grieger, M.: Electronic marketplaces: A literature review and a call for supply chain management research. 144, 280--294 (2003)
16. Minghua He, Nicholas R. Jennings, and Ho-Fung Leung on Agent-Mediated Electronic Commerce. In: IEEE Transactions on Knowledge and Data Engineering, 15, 4 (2003)

17. Urbano, J., Rocha, A. P., Oliveira, E.: Trust Evaluation for Reliable Electronic Transactions between Business Partners. In: AAMAS'09 Workshop on Agent-based Technologies and applications for enterprise interOPerability, pp. 85-96, Hungary, (2009)
18. Cardoso, H. L., Oliveira, E.: Virtual Enterprise Normative Framework within Electronic Institutions. In: Engineering Societies in the Agents World V, LNAI 3451, Springer, ISBN 3-540-27330-1, pp.14-32 (2005)
19. Cardoso, H. L., Oliveira, E.: Electronic Institutions for B2B: Dynamic Normative Environments. In: Artificial Intelligence & Law (special issue on Agents, Institutions and Legal Theory), Springer, Vol. 16, 1, pp. 107-128, ISSN 0924-8463, (2008)
20. Oliveira, E., Rocha, A. P.: Agents Advanced Features for Negotiation in Electronic Commerce and Virtual Organisations Formation Process. In: Agent Mediated Electronic Commerce: The European AgentLink Perspective, LNAI 1991, Springer, pp. 78-97, (2002)
21. Cardoso, H. L., Oliveira, E.: A Contract Model for Electronic Institutions. In: Coordination, Organizations, Institutions, and Norms in Agent Systems III, LNAI 4870, Springer, ISBN: 978-3-540-79002-0, pp. 27-40, (2008)
22. Cardoso, H. L., Oliveira, E.: A Context-based Institutional Normative Environment. In: Coordination, Organizations, Institutions, and Norms in Agent Systems IV, LNAI 5428, Springer, pp. 140-155, (2009)
23. Urbano, J., Rocha, A. P., Oliveira, E.: Computing Confidence Values: Does Trust Dynamics Matter?. In: Artificial Intelligence – 14th Portuguese Conference on Artificial Intelligence, Springer, ISBN 978-3-642-04685-8, pp. 520-531, Portugal, (2009)
24. Cardoso, H. L., Malucelli, A., Rocha, A. P., Oliveira, E.: Institutional Services for Dynamic Virtual Organizations. In: Collaborative Networks and Their Breeding Environments - Sixth IFIP Working Conference on Virtual Enterprises, Springer, pp. 521-528, Spain, (2005.)
25. Cardoso, H. L., Oliveira, E.: Risk Tolerance and Social Awareness: Adapting Deterrence Sanctions to Agent Populations. In: Progress in Artificial Intelligence – 14th Portuguese Conference on Artificial Intelligence, Springer, ISBN 978-3-642-04685-8, pp. 560-571, Portugal, (2009)
26. Cardoso, H. L., Oliveira, E.: Monitoring Cooperative Business Contracts in an Institutional Environment. In: 11th International Conference on Enterprise Information Systems, Italy, (2009)
27. Wooldridge M. An introduction to multi-agent systems. Chichester, England: Wiley (2002)
28. Smith, R. G.: The Contract Net Protocol: High-Level Communication and Control in a Distributed Problem Solver. In: IEEE Trans. Computers, 1104--1113 (1980)
29. Vytelingum, P., Cliff, D., Jennings, N.: Analysing Buyers and Sellers Strategic Interactions in Marketplaces: An Evolutionary Game Theoretic Approach. In: 9th International Workshop on Agent-Mediated Electronic Commerce, pp. 14-15, Hawaii, (2007)

# Real-Time Obstacle Avoidance for Intelligent Wheelchairs

Marcelo R. Petry[1,2], Rodrigo Braga[1,3], Luis Paulo Reis[1,3], António Paulo Moreira[1,2]

[1]FEUP - Faculty of Engineering of the University of Porto,
Rua Dr. Roberto Frias s/n, 4200-465, Porto, Portugal
[2]INESC-P Institute for Systems and Computer Engineering of Porto,
Rua Dr. Roberto Frias 378, 4200-465, Porto, Portugal
[3]LIACC - Artificial Intelligence and Computer Science Lab.
Rua Dr. Roberto Frias s/n, 4200-465, Porto, Portugal
{marcelo.petry, rodrigo.braga, lpreis, amoreira}@fe.up.pt

**Abstract.** Intelligent wheelchairs operating in dynamic environments need to sense its neighbourhood and adapt the control signal, in real-time, to avoid collisions and protect the user. In this paper we propose a robust, real-time obstacle avoidance extension of the classic potential field methodology. Our algorithm is specially adapted to share the wheelchair's control with the user avoiding risky situations. This method relies on the idea of virtual forces, generated by the user command (attractive force) and by the objects detected on each ultrasonic sensor (repulsive forces), acting on the wheelchair. The resultant wheelchair's behaviour is obtained by the sum of the attractive force and all the repulsive forces at a given position. Experimental results from drive tests in a cluttered office environment provided statistical evidence that the proposed algorithm is effective to reduce the number of collisions and still improve the user's safety perception.

**Keywords:** Intelligent wheelchair, obstacle avoidance, potential field methods

## 1   Introduction

Motivated to answer to numerous mobility problems, many intelligent wheelchair related projects have been created in the last years [1]. In fact, intelligent wheelchairs (IWs) are a very good solution to assist handicapped people who are unable to operate classic electric wheelchairs by themselves in their daily activities. While some initiatives have improved the autonomous function of the mobility aid [2], [3], others focused their work in sharing the wheelchair's control with the user [4], [5]. Shared control initiatives take advantage of the user's intelligence and assist the driver in the navigation process when dangerous situations are detected, extending and complementing user capabilities.

In such a way, techniques as obstacle avoidance developed in fields of robotics have the potential to improve user's safety and reduce the navigation complexity. These methodologies consist basically on shaping the robot's path to overcome

unexpected obstacles. A number of algorithms were develop to overcome obstacles and differ basically in the sensorial data used and control strategies. However, not all techniques are suitable to be implemented in a shared control paradigm. Some of the desired properties of shared control algorithms are:

- Avoid obstacles in real-time. Since wheelchairs operate in dynamic environments, it is not feasible to implement popular time-consuming global path planners. Instead, such application is more suitable to approaches based on fast response like reactive/reflexive controls.
- Low computational consuming. Low memory and processing consuming algorithms are more likely to achieve a real-time reflexive behaviour in embedded systems.
- Increase user safety and user safety perception. Beyond a quantitative reduction in the number of collisions, shared control initiatives may consider qualitative evaluations of the wheelchair's overall behaviour. In spite of imposing the control to the wheelchair, the algorithm may adapt the control signal to reduce the discomfort caused in driving tasks.

Furthermore, once intelligent wheelchairs are designed to carry people with disabilities, they should have the same durability, functionality and ergonomics concern of the standard powered wheelchairs. It not only constrains the number of sensors, but their type and position on the wheelchair. Therefore, the shared control algorithm may be robust enough to ensure the user safety even with non-optimal amount of information.

Following the above referred features, this paper proposes and implements an extension of a classic obstacle avoidance technique known as potential field. Special attention was given to user the autonomy, assisting the wheelchair's control just when dangerous situations were detected.

The rest of the paper is subdivided as following: a brief description of some obstacle avoidance methodologies is presented in Section 2. Section 3 presents a Potential Field extension specially designed to assist users in intelligent wheelchair's navigation. Section 4 presents experimental results obtained with both simulated and real drive tests and section 5 presents discussions, final conclusions and some future research topics.

## 2 Related Work

This section presents a concise summary of some classic obstacle avoidance methodologies developed in robotics. In particular, shortcomings and characteristics of each technique are discussed with focusing its implementation on intelligent wheelchairs with shared control paradigm.

### 2.1 Edge-detection

Considered one of the first methodologies proposed to avoid obstacles during robot navigation, edge-detection became popular in late eighties. Through ultrasonic-sensor

readings, the algorithm tries to map the position of the vertical edges of every obstacle in the robot surroundings. Then, once new edges are found, a temporary map is updated and an optimum path planning algorithm is applied to plan the subsequent path [6]. Whenever the robot moves, a scanning mode starts taking alternate sonar's samples. Once measures are under a certain safety distance, the robot stops, take a panoramic scan, apply the edge-detection methodology and restart the cycle all over again [7], [8].

Edge-detection methodology is not itself an obstacle avoidance technique. Actually, it can be better described as an approach to represent the environment based on geometrical primitive line segments. Therefore, off-line path planners are still needed in order to yield obstacle-free paths, limiting its implementation in low-resource embedded systems.

## 2.2    Certainty grids

Certainty grid (CG) method is a probabilistic representation of obstacles in a grid based world model. This world model has been developed for mobile robots in Stanford and CMU for more than ten years, and was originally designed to handle sonar's inaccuracies shortcomings [9].

In this method, the robot's work area is modelled as a 2-D array of square elements, called cells. Each cell of the grid contains a likelihood estimate (certainty value) that indicates confidence that an obstacle is placed within the corresponding region of space. Once readings are more likely to detect objects closer to the acoustic axis of the sonar, a probabilistic function updates more the certainty value in this region  than in the other areas enclosed by the sensor [9], [10].

In spite of some improvements presented by CG methodology, some drawbacks can compromise its implementation in real-time applications. Firstly, the accuracy provided is too much dependent of the cell size. Secondly, as the robot moves over large areas, lots of memory and processing power are required, restricting the application of CG especially in some embedded systems. Finally, the subsequent robot's path shall be computed off-line, by a global path-planning.

## 2.3    Vector Field Histogram

Introduced by Borenstein and Korem [6], the Vector Field Histogram (VFH) uses a polar histogram instead of a 2-D Cartesian grid to avoid collisions and steer the mobile robot to the target. This method employs a two-stage data reduction process in order to compute the control command to the robot.

In the highest level of data, VFH stores a detailed 2-D histogram grid map of the robot's neighbourhood. As just only one cell in the histogram is updated for each range reading, it takes just a small computational overhead. Thus, a probabilistic distribution is obtained by continuously and quickly sampling each sensor while the robot moves [11]. At the second level, data is mapped onto a one-dimensional polar

histogram, that comprises $n$  angular sections each with width $\alpha$ . Each sector in the polar histogram contains a value representing the polar obstacle density in that direction. Finally, based on the obstacle polar density (1-D histogram), VFH selects

the best steering direction for the robot and computes the reference values for driving the robot (third level of data representation) [6].

As can be observed, the VFH overcomes some issues shown by the other methods described above. In fact, the influence of low accuracy distance measures is minimized through the histogram representation. In addition, the world representation is restricted to the robot's surrounding trough a bi-dimensional sliding window, reducing the computational overhead. On the other hand, local minima problems are still not solved by the algorithm itself, which has to invoke a global path planner when these situations are flagged. Finally, like edge-detection and CG methodologies, VFH depends not only from the data gathered by the sonars, but from an accurate localization system. Otherwise, inaccurate robot's position can introduce more errors and disturb the object mapping.

### 2.4    Potential Field Methods

First suggested by Andrews and Hogan [12] and Khatib [13], the Potential Fields methodologies (PF) relies on a simple and powerful principle, the artificial potential field concept. In this method, the robot is considered immersed in a potential field generated by the target and by obstacles. In this field, obstacles generate imaginary repulsive forces, while the target generates an attractive force to the robot. The resultant robot behaviour is obtained by the sum of all attractive and repulsive forces at a robot's given position.

After the original work, a number of improvements and extensions have been published. Krogh [14] has computed forces not only to steer the robot around objects, but to set its speed as well and Seiki [15] has introduced the consideration about the nonholonomic motion constrains and the robots shape into the PF. Khatib and Chatila [16], considered, besides distance, the robot's relative orientation to the obstacle in order to compute forces. Bicho [17] implemented a dynamic approach using low level sensory information, in which each sensor generates a repulsive force that drives the direction and the speed of the robot.

In its original version the PF methodology exhibit many shortcomings, in particular the sensitivity to local minima that arises mostly due to the symmetry of the environment. Furthermore, it tends to be very susceptible to misreading (since it takes into account just one set of data) and to the sonar most common issues. Some versions still assume a known and prescribed world model to evaluate off-line the potential field. Finally, some implementations present significant problems related to oscillations in narrow passages and in the presence of obstacles [6], [18].

### 3    Local Obstacle Avoidance

The potential field concept was chosen as base for our implementation given its simplicity. Especially due to the possibility to easily adapt the algorithm to cover the specific requirements of shared control paradigms and to run it on the limited computational capability of our prototype's embedded system. However, our work

differs from the original PF because it does not try to build a world map of the environment. Instead, our it is closer to the implementation described by Bicho [17], where each ultrasonic range reading is treated as a repulsive force.

Once an object is detected by a sensor $S_i$, a virtual repulsive force $F_i$ towards the robot is computed. The direction of each repulsive force is determined by the direction of $\sigma_i$, from the object point $O_i$ to the Robot Center Point $C$ , Fig.1. Notice that since sonar sensors return radial measures of the environment, it is not possible to determine precisely the angular location of the object. However, it is much more likely that the detected object is closer to the acoustic axis of the ultrasonic transceiver then in the periphery of the conical field of view [13]. Thus, the position of obstacle $O_i$ is computed as the measured distance $D_i$ under the acoustic axis of the sensor.

$$\sigma_i = tan^{-1}\frac{O_{iy}}{O_{ix}} \tag{1}$$

Where

$O_{ix}, O_{iy}$         Relative position of obstacle detected by the sensor $S_i$ .

$\sigma_i$         Direction from the object $O_i$ to the wheelchair's center point $C$ .



**Fig. 1.** Repulsive forces          **Fig. 2.** Resultant force that steers the whelchair

In order to keep user autonomy at the utmost, control signals are only adapted in situations which the user faces an eminent risk of collision. Therefore, repulsive forces start acting just when a safety range is reached. Due to inertia, the distance needed to completely stop the wheelchair increases with its speed $S$ . Thus, the collision risk is considered as a bi-dimensional variable, both distance and speed dependent, Fig. 3 and Fig. 4.

Such safety range is designed not just to avoid obstacles in the wheelchair's neighbourhood, but also to avoid oscillations that non-critical far objects could cause in the control's behaviour. The magnitude of the repulsive forces grow exponentially accordingly the pair ($D_i$, $S$):

$$[\square|F|_i] = \alpha * exp(-\beta * D_i + \omega * S) * [\square|F|_\alpha] \tag{2}$$

where

$\alpha$, $\beta$, $\omega$      Positive constants deduced from the desired safety range.

$F_\alpha$      Attractive force

$D_i$      Distance measured by the sensor Si



**Fig. 3.** Safety distance range in function of the wheelchair's speed and distance      **Fig. 4.** Speed component of the safety range

Once all repulsive forces are computed, they are added up to yield a resultant repulsive force $F_r$.

$$F_r = \sum_{i=0}^{n} F_i \tag{3}$$

Next, the virtual attractive force $F_a$ induced by the target is updated. In the wheelchair implementation the force $F_a$ is directly proportional to the current user input, which can be either the standard wheelchair's joystick or a special user interface which is based on the user's head position. Summing both the resultant repulsive force $F_r$ and current attractive force $F_a$ it is possible to derive the final force $F_t$ that steers the wheelchair, Fig. 2.

$$F_t = F_\alpha + F_r \tag{4}$$

# 4    Experiments

In order evaluate the efficiency of the proposed algorithm six volunteers performed each one set of four drive tests. Sets were composed four laps: two laps in a simulated environment (with and without the assistance of the algorithm) and two laps in a real environment (with and without the assistance of the shared wheelchair control). All of the six recruited participants were aged between 26 and 39 years old, and have spent around 40 minutes running the experiments and answering a post-test questionnaire.

Based on the work proposed by Parikh [19], a well-defined protocol to conduct the test was designed. It aims to ensure that data were collected accurately and in the same way across the tests, and will be better explained in the section 4.1.

Participants were instructed about the objective of the task and about the closed circuit they should drive, Fig. 6. It was reinforced that their main goal was to drive safely and then to finish each lap as fast as they could. Time was just mentioned as a secondary objective to prevent volunteers from navigating too slowly, and was not used on the evaluation process.

Real tests were run using IntellWheels intelligent wheelchair prototype, Fig 5. It is equipped with a ring of eight ultrasonic sensors distributed accordingly Fig. 1, 2 encoders and one embedded ATmega1280 microcontroller board to run the algorithm. Due to some constrains related to the update rate of the ultrasonic sensors, each algorithm cycle spent 80 ms to be computed. Simulation tests were run under the IntellWheels Simulator, emulating the same characteristics of the real prototype and the real environment. Further details about the prototype and the simulator can be found in [20], [21], [22].



**Fig. 5.** Intelligent wheelchair prototype used during the tests in the real environment



**Fig. 6.** Closed circuit which the experiments were conducted

During these trials, some conditions faced by handicapped individuals have been simulated. To accomplish that, all participants were asked not to drive the wheelchair using its standard hand driven joystick. Instead, volunteers were requested to perform all four laps using a special human-machine interface based on the user's head position [21], [23].

### 4.1 Experiments Protocol

This experiment protocol has been defined to standardize the results of both tests, and consists basically of seven steps:

- Step 1: volunteers have been instructed about test procedure and about their objectives during the four drive tests.
- Step 2: it was given to the participant a 10 minutes driving trial in a simulated environment. Thus, the user could experiment the wheelchair and make the necessary adjustments to the special human-machine interface.
- Step 3: once prepared, the participant was asked to drive the wheelchair (1 lap) through the circuit in the simulated environment with the manual control paradigm.
- Step 4: after the first test, it was asked to the volunteer to drive the wheelchair (1 lap) through the same circuit in the simulated environment, but with the assistance of the shared control.
- Step 5: accomplished both tests in the simulator, the user were asked to drive the wheelchair (1 lap) in the real environment with the manual control.
- Step 6: in the last test the user had to drive the wheelchair (1 lap) in the real environment with the shared control paradigm.
- Step 7: to evaluate the shared control paradigm, the user safety perception and to conclude the set of experiments, a pot-task questionnaire was applied.

### 4.2 Results

From the set of experiments described above, both quantitative and qualitative data have been generated. All analysis were performed within subjects, which allowed us to estimate if providing assistance actually helped each individual, rather than testing the performance of individuals against each other. Based on the number of collisions of each trial, the shared algorithm performance could be evaluated in the simulation environment, Fig. 7, and in the real environment, Fig. 8.



**Fig. 7.** Number of collisions per volunteer on the simulated environment

**Fig. 8.** Number of collisions per volunteer on the real environment

On the other hand, it is important to evaluate the algorithm from the user's perspective. Related projects concluded that in spite of their algorithm were

responsible in reducing the number of collisions, users did not felt safer and preferred driving with no assistance. In order to measure the user feeling, the questionnaire applied was composed of five statements for each control paradigm, in which respondents were invited to specify their level of agreement on a five-point Likert scale (1 = Strongly disagree, 2 = Disagree, 3 = Neither agree nor disagree, 4 = Agree, 5 = Strongly agree):

1. I feel comfortable when driving the wheelchair.
2. I feel that I have the control of the wheelchair behaviour.
3. It is easy to drive the wheelchair in cluttered spaces.
4. Driving the wheelchair requires little attention.
5. The wheelchair has the same behaviour either in the simulated and the real environments.
6. I believe that the shared control helped me during the navigation task.

In our analysis, the user safety perception was treated as an indirect variable measured through the sum of the points of the first four statements, Fig. 9.



**Fig. 9.** User's safety perception with and without the assistance of the shared control.

Another inference can be done regarding to the behaviour of the wheelchair in the simulator. Through the fifth statement, we tried to measure how close to the reality the simulated behaviour of the wheelchair is. A threshold value of 3 was used to compare results, Fig. 10. Remember that in Likert scale a value of 3 means that respondents neither agree nor disagree with the statement, thus a value greater than 3 means that simulated wheelchairs react just like the real ones in the user's perspective.



**Fig. 10.** How realistic is the simulator on manual and shared control paradigms.

**Fig. 11.** User's help perception with the assistance of the shared control.

Finally, one last result of the questionnaire is the user's perception of the help provided by the wheelchair. In this case it was evaluated through the statement 6, only present in the shared control section of the questionnaire, and compared with a threshold value of 3, Fig. 11. Similar to what was mentioned before, a value greater than 3 mean that the user felt helped by the algorithm.

## 5    Conclusions

This paper presented a new approach to improve safety for wheelchair's users. To assist patients in their navigation tasks, an obstacle avoidance methodology has been adopted. Based on the dynamic approach of the classic field of forces concept, this work extends and complements the potential field methodologies from a shared control perspective. To reduce the computational cost and run the algorithm in real-time, each ultrasonic range reading is treated as a repulsive force. Thus, it is not necessary to build a map of the environment and compute thousands of parameters. Furthermore, as localization is not required, dead reckoning errors are not introduced when computing the distance to obstacles.

On the other hand, such proposal is still very sensor-depended, and does not overcome by itself the intrinsic sonar shortcomings. Further improvements should include some sensor filtering to increase robustness and reduce measures oscillations. Another problem detected during the experiments regards to the prototype and not to the algorithm itself. It was noticed the presence of two small blind spots for object closer than 25 cm, one at each side of the wheelchair. Both of them cannot be reached during a forward/rewind displacement, but only when turning in sharp corners. In our tests it was the main reason of those collisions during the tests with the shared control paradigm, and could be eliminated with the addition of two more sonars.

As depicted in Figures 7 and 8, the number of collisions in the simulated environment is clearly much greater, and can be explained through the volunteer's comments and their behaviour during the tests. First, the 3D environment of the simulator could not provide an accurately perception of depth and distance to objects, causing collisions in the cluttered test circuit. However, the second reason is related to the user way of driving. In the simulated environment volunteers tended to relax and to reduce the attention to the circuit mostly because they were not going to suffer the physical damages of a collision.

The most interesting results came out after a statistical analysis of the experimental data. Thus, based on a paired two-tailed t-student test with $p<0.05$, it is possible to conclude that:

1. In the simulated environment, the results indicate that there is significant difference between the number of collisions with and without the shared control paradigm (t= 6.028, p= 0.002), providing statistical evidence that the shared control is effective in reducing the number of collisions.
2. In the real environment, the results indicate that there is significant difference between the number of collisions with and without the shared control paradigm (t= 2.582, p= 0.049), providing statistical evidence that the shared control is effective in reducing the number of collisions.

3. A significant difference between the user safety perception with and without the shared control paradigm (M=-1.166, SD = 0.983, N= 6) was significantly greater than zero (t= -2.907, p= 0.034), providing statistical evidence that the shared control is effective to improve user's safety perception.

4. The mean of the user's help perception variable (M=3.67, SD= 0.516, N= 6) was significantly greater than the test value 3 (t= 3.162, p = 0.025), providing evidence that volunteers indeed felt that the shared control paradigm helped then to drive the wheelchair.

5. However, through the t-student test it was not possible to state with a confidence level of 95% that the wheelchair has the same behaviour in the real and simulated environments for both manual and shared control paradigms. In spite of only one value under the threshold, the low number of samples could be considered the main cause of this result.

Our research has demonstrated that it is possible to increase safety with low cost sensors and improve the acceptance of intelligent wheelchairs.

Future work will address the wheelchair simulator by increasing its realism in order to achieve very similar real and simulated behaviours. Future work will also be concerned with increasing the number of volunteers in the drive tests, and testing the algorithm on impaired people with different diseases.

### Acknowledgments

## 6 References

1. Simpson, R., LoPresti, E., Hayashi, S., Nourbakhsh, I., Miller, D.: The Smart Wheelchair Component System. Journal of Rehabilitation Research and Development 41, pp. 429--442 (2004)
2. Hamagami, T., Hirata, H.: Development of Intelligent Wheelchair Acquiring Autonomous, Cooperative, and Collaborative Behavior. In: IEEE International Conference on Systems, Man & Cybernetics, pp. 3525--3530 (2004)
3. Bourhis, G., Agostini, Y.: The VAHM Robotized Wheelchair: System Architecture and Human-Machine Interaction. Journal of Intelligent & Robotic Systems 22, pp. 39--50 (1998)
4. Bell, D.A., Borenstein, J., Levine, S.P., Koren, Y., Jaros, L.: An Assistive Navigation System for Wheelchairs Based Upon Mobile Robot Obstacle Avoidance. In: IEEE International Conference on Robotics and Automation, 2018-2022 (1994)
5. Lankenau, A., Rofer, T.: A Versatile and Safe Mobility Assistant. IEEE Robotics & Automation Magazine 8, pp. 29--37 (2001)
6. Borenstein, J., Koren, Y.: The Vector Field Histogram - Fast Obstacle Avoidance for Mobile Robots. IEEE Transactions on Robotics and Automation 7, pp. 278--288 (1991)
7. Borenstein, J., Koren, Y.: Obstacle Avoidance with Ultrasonic Sensors. IEEE Journal of Robotics and Automation 4, pp. 213--218 (1988)

8.  Elfes, A.: Sonar-Based Real-World Mapping and Navigation. IEEE Journal of Robotics and Automation 3, pp. 249--265 (1987)
9.  Wang, T.Z., Yang, J.: Certainty Grids Method in Robot Perception and Navigation. In: IEEE International Symposium on Intelligent Control, pp. 539--544 (1995)
10. Moravec, H.P.: Sensor Fusion in Certainty Grids for Mobile Robots. Ai Magazine 9, pp. 61--74 (1988)
11. Borenstein, J., Koren, Y.: Histogramic In-Motion Mapping for Mobile Robot Obstacle Avoidance. IEEE Transactions on Robotics and Automation 7, pp. 535--539 (1991)
12. Andrews, J.R., Hogan, N.: Impedance control as a framework for implementing obstacle avoidance in a manipulator. In: ASME Winter Conference, pp. 243--251 (1983)
13. Khatib, O.: Real-Time Obstacle Avoidance for Manipulators and Mobile Robots. International Journal of Robotics Research 5, pp. 90--98 (1986)
14. Krogh, B.H.: A Generalized Potential Field Approach to Obstacle Avoidance Control. In: International Robotics Research Conference (1984)
15. Seki, H., Shibayama, S., Kamiya, Y., Hikizu, M.: Practical Obstacle Avoidance Using Potential Field for a Nonholonmic Mobile Robot with Rectangular Body. In: IEEE International Conference on Emerging Technologies and Factory Automation, pp. 326--332 (2008)
16. Khatib, M., Chatila, R.: An Extended Potential Field Approach for Mobile Robot Sensor-Based Motions. In: International Conference on Intelligent Autonomous Systems, pp. 490--496 (1995)
17. Bicho, E., Mallet, P., Schoner, G.: Using Attractor Dynamics to Control Autonomous Vehicle Motion. In: 24th Annual Conference of the IEE Industrial Electronics Society, pp. 1176--1181 (1998)
18. Koren, Y., Borenstein, J.: Potential-Field Methods and Their Inherent Limitations for Mobile Robot Navigation. In: IEEE International Conference on Robotics and Automation, pp. 1398--1404 (1991)
19. Parikh, S.P., Grassi, V., Kumar, V., Okamoto, J.: Usability Study of a Control Framework for an Intelligent Wheelchair. In: IEEE International Conference on Robotics and Automation pp. 4745--4750 (2005)
20. Braga, R.A.M., Petry, M., Moreira, A.P., Reis, L.P.: IntellWheels - A Development Platform for Intelligent Wheelchairs for Disabled People. In: 5th International Conference on Informatics in Control, Automation and Robotics, pp. 115--121 (2008)
21. Braga, R.A.M., Petry, M., Moreira, A.P., Reis, L.P.: Concept and Design of the Intellwheels Development Platform for Intelligent Wheelchairs. In: Cetto, J.A., Ferrier, J.L., Filipe, J. (eds.). LNEE, vol. 37, pp. 191--203. Springer-Verlag, Heidelberg (2009)
22. Braga, R.A.M., Malheiro, P., Reis, L.P.: Development of a Realistic Simulator for Robotic Intelligent Wheelchairs in a Hospital Environment. In: Baltes, M.Lagoudakis, T.Naruse, S.Shiry (Eds.) Robocup 2009. LNAI, vol. 5949, pp.23--34. Springer, Heidelberg (2009)
23. Reis, L.P., Braga, R.A.M., Sousa, M., Moreira, A.P.: Intellwheels MMI: A Flexible Interface for an Intelligent Wheelchair. In: Baltes, M.Lagoudakis, T.Naruse, S.Shiry (Eds.) Robocup 2009. LNAI, vol. 5949, pp.296--307. Springer, Heidelberg (2009)

# Biped Walking using Coronal and Sagittal Movements based on Truncated Fourier Series

Nima Shafii[1], Luís Paulo Reis[1,2], Nuno Lau[3,4]

[1]Faculty of Engineering of the University of Porto, Rua Dr. Roberto Frias s/n, Porto, Portugal
[2]Artificial Intelligence and Computer Science Lab., Rua Dr. Roberto Frias s/n, Porto, Portugal
[3]University of Aveiro, Campus Universitário de Santiago, Aveiro, Portugal
[4]IEETA – Institute of Electronics and Telematics Engineering of Aveiro, Portugal
nima.shafii@gmail.com, lpreis@fe.up.pt, nunolau@ua.pt

**Abstract.** Biped walking by using all joint movements and DOFs in both directions (sagittal plane and coronal plane) is one of the most complicated research topics in robotics. In this paper, angular trajectories of a stable biped walking are generated by a Truncated Fourier Series (TFS) approach. The movement of legs and hand in sagittal plane are implemented by an optimized gait generator and for the first time a model is presented that can produce the movement of legs in coronal plane based on Truncated Fourier Series. Particle Swarm Optimization (PSO) is used to find the best angular trajectories and optimize TFS. Experimental results show that training of the robot can be successfully performed by our method, thus allowing the biped robot to walk faster by using all joints movement in sagittal plane and coronal plane.

**Keywords:** Bipedal Locomotion. Gait Generation. Particle Swarm Optimization.

## 1. Introduction

Bipedal walking has become more popular in the last few years, since robots can do their work in a hard and complex environment by means of legged locomotion. But bipedal walking still needs more time to achieve this goal. Various methods for bipedal locomotion have been presented until now in the literature. We can classify these methods in two large groups: model-based and model free. In model-based approaches first the physical model of the robot is designed and after that a controller is built for it. "Zero Moment Point" (ZMP) [1] and "Inverted Pendulum" [2] are two methods that belong to this approach.

In model–free approaches the sensory information is used for making motions. In fact, there isn't any consideration of physical model in this method and skills are implemented in an easier way. Passive Dynamic Walking (PDW) [3], Central Pattern Generator (CPG) [4] and Ballistic Walking [5] are the most known methods of model-free approach. In the PDW approach, the robot does not have any actuators' model and it walks just by utilizing the gravity force. The Ballistic Walking is originated from the simple human walking model based on the observation of human walking in

which the muscles of the swing leg are activated only at the beginning and the end of the swing phase. In the CPG approach, special neural circuits take the role of rhythmic walking controller using the non-linear equations to model the neural activities. Researchers usually focus on complex mathematical models like Hopf [6] or Matsuoka [7] to model these neural activities and generate rhythmic walk patterns (Gait). In 2007 Truncated Fourier Series Formulation method was used as a gait generator in bipedal locomotion [8]. In this article Truncated Fourier Series (TFS) together with a ZMP stability indicator was used to prove that TFS could generate suitable angular trajectories for controlling bipedal locomotion but it was not implemented on real robot [8]. In 2009, an optimized gait generator based on TFS was implemented in a simulated humanoid robot and TFS parameters were also reduced by 2 dimensions (down to 6 dimensions) [9]. Shafii extended the basic TFS which is capable to produce hand angular trajectories with emphasizing the role of hands in smooth and robust walking and used a new method to refine signals for reducing the role of inertia to improve the speed and robustness of the robot [10].

In this paper the results of the two previous papers are used to produce walking movements on sagittal plane. The method was tested on a simulated NAO humanoid and the experiments were performed using Rcssserver3d [11], a generic three-dimensional simulator which is based on Spark and Open Dynamics Engine (ODE). The robot model has 22 DOF with a height of about 57cm, and a mass of 4.5kg.

The paper structure is as follows. First, optimized TFS gait generator is introduced to generate walking movement on the sagittal plane. Hand angular trajectories generator and the method for reducing the role of inertia are also explained.

Then a new model of hip angular trajectory generator based on TFS is introduced which can produce the leg's movement on the coronal plane (Y direction). Particle Swarm Optimization (PSO) is used to optimize the produced signals, to overcome inherent noise of the simulator, Resampling algorithm is implied which could lead to robustness in nondeterministic environments. At the end of the article, results of this approach are presented and the efficiency of the method on producing trajectories to walk a robot in the forward direction is shown.

## 2. Movements in Sagittal Plane

There are three DOFs in each leg move in sagittal plane; one in the hip, one in the ankle and one at the knee. In this work, similar to [12], foot in sagittal plane was kept parallel to the ground by using ankle joint. This is done in order to avoid collision. Therefore ankle trajectory can be calculated by hip and knee trajectories and ankle DOF parameters are eliminated. Trunk sagittal and coronal plane motion is fairly repeatable [13] therefore Fourier series can be used.

### 2.1 An optimized gait generator for leg's movement

In this model, legs joint angular trajectories in sagittal plane are divided in two parts; the upper portion and the lower portion. $C_h$ is offset of hip trajectory and $C_k$ is offset

of knee trajectory. From $t_1$ to $t_2$ the left leg is considered as supporting leg and the variation of its knee angle is so little that it can be assumed fixed. This duration of walking is named knee lock phase. In addition, the shift phase of the two leg trajectories signal is as half of the period of each signal so by producing the trajectory of one leg the other leg's trajectory can be calculated. The trajectories for both legs are identical in shape but are shifted in time relative to each other by half of the walking period. The TFS for generating each portion of hip and knee trajectories are formulated in below.

$$
\theta_h^+ = \sum_{i=1}^{n} A_i . \sin\left(iw_h t\right) + c_h, w_h = \frac{2\pi}{T_h}
$$

$$
\theta_h^- = \sum_{i=1}^{n} B_i . \sin\left(iw_h t\right) + c_h, w_h = \frac{2\pi}{T_h}
$$

$$
\theta_k^+ = \sum_{i=1}^{n} C_i . \sin\left(iw_k t\right) + c_k, w_k = w_h
$$

$$
\theta_k^- = c_k \geq 0
$$

**(1)**

In these equations, the plus (+) sign represents the upper portion of walking trajectory and the minus (-) shows the lower portion. $i=1$ and $A_i$, $B_i$, $C_i$ are constant coefficients for generating signals. The $h$ and $k$ subscripts stands for hip and knee respectively. $C_h$, $C_k$ are signal offsets and $T_h$ is assumed as the period of hip trajectory. Considering the fact that all joints in walking motion have equal movement frequency and stride rates is statistically equal, the equation $w_k = w_h = \frac{2\pi}{T_h}$ can be concluded.

Parameter $t_1$ is the start time of lock phase for knee joint and parameter $t_2$ represents the end time of lock phase and in this duration of time $\theta_k^- = c_k \geq 0$.

According to [9], by specifying the start and end time of the lock phase, two parameters of $t_1$, $t_2$ could be eliminated. Therefore the number of variable for optimization to produce legs movement in sagittal plane decreased to 6. This omission has many advantages such as; reducing the search space of optimization problem and increasing the convergence speed of PSO.

## 2.2 Modeling of arm motion in sagittal plane

In sagittal plane, during human walking, the arms normally swing in opposite manner to legs, which helps to balance the angular momentum generated in the lower body [14]. Humans swing their arms close to 180º out of phase with their respective legs during walking [15]. Fig. 1 shows the trajectory of legs and arm swings and the relation between them in a stable straight walking [14]. It is shown that Trajectory of arms is similar to sinusoidal signal with same frequency of legs.
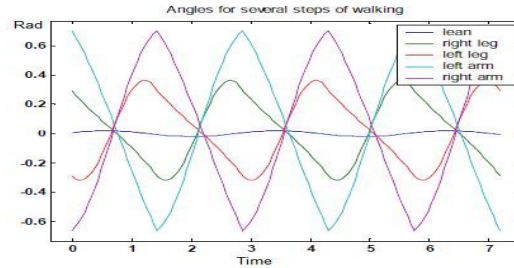
**Fig. 1.** Trajectories of legs and arms

It is found that Speed of walking has strong effect on arm swinging during gait. By increasing gait's speed, the arms may swing higher and faster to reduce the effects of longer, quicker steps by the legs [16].

It can be expected that the utilization of arm swing provides good performance to yaw moment stability, and recovery from stumbling. The effectiveness of this method is confirmed with an improvement of the accuracy of straight walking at different speeds. As has been shown, the trajectory of arms is a sinusoidal signal; therefore, to produce proper angular trajectories to arms swing, it is enough to obtain proper parameters for the following equation (2).

$$f(t) = A \sin(\omega_{arm} t), \omega_{arm} = \omega_{arm} \qquad \textbf{(2)}$$

In the above equation, $A$ and $\omega$ are assumed as the amplitude and frequency of the signal, respectively. In addition the shift phase of the two arm trajectories signals is half of the period of each signal, so by producing the trajectory of one arm the other arm's trajectory can be calculated. Since legs and arms have the same frequency, $\omega arms$ can be considered equal to $\omega legs$. According to the fact that the angle of arms is zero at the start of walking, shift phase factor is assumed as 0.

According to mentioned equation only the proper value of parameter $A$ must be obtained. To reach this aim we consider this parameter together with legs trajectory parameters in the optimization problem.

## 3. Increasing Speed at the Start Time of Walking in TFS

According to [10], increasing the walking speed and amplitude from standing to running the robot can walk more stable and faster. We implemented a model for the robot to walk from smaller gait with lower amplitude to bigger gait with higher speed and acceleration. In this model a linear equation is used to lead the robot to increase the amplitude of trajectory linearly form zero (stop state) to desired angular trajectories. $T$ is assumed as a parameter to determine how much time is needed for this increment algorithm to reach these desired trajectories. All angular trajectories such as arms and legs will be multiplied by the product of the following equation.

$$K = time / T , time < T \qquad \qquad \textbf{(3)}$$
$$K = 1, time >= T$$

## 4. Movements in Coronal plane

The range of motion in the coronal and transverse plane is less than that is seen in the sagittal plane [17] but it has an important role to keep the balance of walking and reach the highest speed of walking. The range of its motion depends on the speed of walking at higher speeds this range is smaller. Coronal plane movements are periodic motions [13]. Abduction and adduction are terms for movements of limbs relative to the coronal plane.

To produce legs' motion in coronal plane and also considering keeping the balance of robot, we proposed a walking sequences and scenario (Fig. 2). It illustrates the walking sequences in a walking period. $\theta$ is assumed as the maximum of legs movement. Like in the previous section in coronal plane, feet were kept parallel to the ground in order to avoid collision and considering of the fact that just one of each hip's joint moves each time, the angle of ankle is equal to the hip's angle of the opposite leg.



**Fig. 2.** Coronal plane view of proposed walking Sequence



**Fig. 3.** Left leg and right leg hips angular trajectories

Considering the walking sequence, Fig. 3 can be assumed as the hip angular trajectory in one period of walking. It is a sinusoidal signal that has a lock phase at

zero degree. Therefore in order to produce proper angular trajectories to move the hips in coronal plane, proper parameters for the following equation must be obtained(4).

$$f(t) = H \sin(wt), t < T_h / 2 \qquad \qquad \textbf{(4)}$$
$$f(t) = 0, t > T_h / 2$$

In the above equation, $H$ and $\omega$ are assumed as amplitude and frequency of signal respectively and $T_h$ is assumed a period of hip. As mentioned before, ankle trajectories can be calculated from hip trajectories and as it is shown, left and right hip angular trajectories are the same but with a phase shift of -pi. The period of walking in sagittal plane and coronal plane is equal. Therefore $T_h$ and $\omega$ is eliminated in this method and for producing the proper abduction and adduction, the proper value of $H$ parameter must be found.

According to the fact that movements in coronal plane decrease at higher speeds of walking, so by increasing the speed of walking mentioned in previous section, the movement of legs in coronal plane must be reduced. Therefore all angular trajectories generated by above model will be multiplied by the reverse product of the equation 3.


# 5. Implementation

Bipedal walking is known as a complicated motion since many factors affect walking style and stability such as robot's kinematics and dynamics, collision between feet and the ground. In such a complex motion, relation between Gait trajectory and walking characteristic is nonlinear. In this approach the best parameters to generate angular trajectories for bipedal locomotion must be found. According to [18], for this kind of optimization problem, Particle Swarm Optimization can achieve better results. Therefore PSO seems to be an appropriate solution.


## 5.1 PSO algorithm

The PSO algorithm consists of three steps; generating primitive particle's positions and velocities, velocity update and position update [19]. These parts will be described in sections 5.2, 5.3 and 5.4 respectively.


## 5.2 Initializing particles' positions and velocities

Equations (5) and (6) are used to initialize particles in which $\Delta t$ is the constant time increment. Using upper and lower bounds on the design variables values, $X_{min}$ and $X_{max}$, the positions, $X_k^i$ and velocities, $V_k^i$ of the initial swarm of particles can be first generated randomly. The swarm size will be denoted by $N$. The positions and velocities are given in a vector format where the superscript and subscript denote the $i^{th}$ particle at time $k$.

$$X_0^i = X_{min} + Rand(X_{max} - X_{min}) \tag{5}$$

$$V_0^i = \frac{X_{min} + Rand(X_{max} - X_{min})}{\Delta t} = \frac{Position}{time} \tag{6}$$

## 5.3 Updating Velocities

The fitness function value of a particle is used to determine which particle has the best global value in the current swarm ($P_k^g$), and to determine the best position of each particle over time ($P^i$).

The three values that affect the new search direction, namely, current motion, particle own memory, and swarm influence, are incorporated via a summation approach as shown in Equation below (7) with three weight factors, namely, inertia factor, $w$, self confidence factor, $C_1$, and swarm confidence factor, $C_2$, respectively.

$$\underset{\substack{Velocity\ of\ Particle \\ i\ at\ time\ k+1}}{V_{k+1}^i} = \overset{[0.4,1.4]}{w}\ \underset{\substack{Current \\ Motion}}{V_k^i} + \overset{[1,2]}{C_1}\ Rand\ \underset{Particle\ Memory\ Influence}{\underbrace{\frac{(P^i - X_k^i)}{\Delta t}}} + \overset{[1,2]}{C_2}\ Rand\ \underset{Swarm\ Influence}{\underbrace{\frac{(P_k^g - X_k^i)}{\Delta t}}} \tag{7}$$

The inertia weight $w$ controls how much of the previous velocity should be retained from the previous step. A larger inertia weight facilitates a global search, while a smaller inertia weight facilitates a local search [20]. Introducing a nonlinear decreasing inertia weight as a dynamic inertia weight into the original PSO significantly improves its performance through the parameter study of inertia weight [20]. This nonlinear distribution of inertia weight is expressed as follow:

$$w = w_{init} * U^{-k} \tag{8}$$

Where $w_{init}$ is the initial inertia weight value selected in the range [0, 1] and $U$ is a constant value in the range [1.0001, 1.005], and $k$ is the iteration number.

## 5.4 Updating the Position

Position update is the final step of each iteration and it is done by using the current particle position and its own updated velocity vector shown in the Equation below.

$$X_{K+1}^i = X_K^i + V_{K+1}^i \Delta t \tag{9}$$

In summary the PSO algorithm is:

```
Initialize Position (X₀) and Velocity of N particles
according to equation (4 and 5)
P¹=X₀
DO
    k=1
    FOR i = 1 to N particles
            IF f(Xᵢ) < f(Pⁱ) THEN Pⁱ= Xᵢ
            Pₖᵍ = min (P)
            Calculate new velocity of the particle according
            to equation (6, 7 and 8)
            Calculate new position of particle according to
            equation (9).
    ENDFOR
    k=k+1
UNTIL a sufficient good criterion (usually a desirable
fitness or a maximum number of iterations) is met.
```

## 5.5 PSO implementation

In PSO, the parameters of the problems are coded into a finite length of string as a particle. According to above sections, for producing movements in sagittal plane, TFS has 6 parameters to generate legs joints angular trajectories and 2 parameters are assumed to swing arms and to increase the speed of walking. There is also one parameter to produce proper legs' movement in coronal plane. Therefore there is a 9-dimension search space for the PSO to find the optimum solution.

Angular trajectories produced by each particle are applied to a simulated robot to make it walk. To use angular trajectory for walking, all individual robot's joints attempt to drive towards their target angles using proportional derivative (PD) controllers. To equip the robot with a fast walking skill a fitness function based on robot's straight movement in limited action time is considered. First the robot is initialized in x=y=0 (0, 0) and it walks for 15 seconds then the fitness function is calculated whenever the robot falls or time duration for walking is over. Fitness function formulation is simply expressed as the distance travelled by the robot along the x axis.

Due to the fact that there is noise in simulated robot's actuators and sensors, training walking task in this approach can be viewed as an optimization problem in a noisy and nondeterministic environment. Resampling is one of the techniques to improve the performance of evolutionary algorithms (EAs) in noisy environment [22]. In Resampling, the individual set of parameters (particle) $y_i$, the fitness $F(y_i)$ is measured $m$ times and averaged yielding fitness. According to (10) the noise strength of $\overline{F}$ is reduced by a factor $\sqrt{m}$.

$$\overline{F\left(y_i\right)} = \frac{1}{m}\sum_{k=1}^{m} F\left(y_i\right), y\left(i\right) = const. \Rightarrow \overline{\sigma_e} = \sqrt{Var\left[\overline{F\left(y_i\right)}\right]} = \frac{\sigma_e}{\sqrt{m}} \tag{10}$$

Since particles may not satisfy some constraints after updating position procedure,

constraint handling is a vital step in PSO algorithm. There are many constraints on parameters in this study (i.e time parameters in TFS must be positive). Therefore Pareto [23] with multi-objective modeling is used for handling constraints.

In Pareto, a solution, *x(2)*, is dominated by solution, *x(1)*, if *x(1)* is not worse than *x(2)* in all objectives, and for at least one of the objectives, *x(1)* is strictly better than *x(2)* .Without loss of generality, these conditions can be expressed as follows for the case where all of the objective functions are to be minimized:

$$fm\left(x\left(1\right)\right) \leq fm\left(x\left(2\right)\right) for \forall m = 1, 2, ..., M \quad \text{and}$$

$$fm\left(x\left(1\right)\right) \prec fm\left(x\left(2\right)\right) \quad \text{for some } m.$$

Each constraint is assumed as an object in which parameters must be satisfied .So according to Pareto method, a particle can be considered to find Pi, Pgk when it satisfies objects and constraints. Therefore calculating fitness for particles that cannot satisfy constraint is not necessary.

We considered various values for each parameter of the algorithm and tried all possible combinations. Finally we chose the best combination of the parameters regarding the dynamic inertia weight and test results that $C_1$ and $C_2$ are assumed as 1, 1.5, $w_{init}$ as 0.8, $U$ as 1.0002 and $\Delta t$ as 1, respectively. We have also implied a swarm consisted of 100 particles (N = 100) and maximum iteration of 10 and Resampling factor *m* is assumed as 3.

# 6. Results

To compare the presented method with Basic TFS method which uses the joints movement just in sagittal plane, both of them were tested using the same system and the same specification. We also optimized both methods by utilizing PSO with equal specifications and with the same fitness function.

Using basic TFS with 8 parameters and running PSO algorithm on a Pentium IV 3 GHz Core 2 Duo machine with 2 GB of physical memory, 3000 trials were performed in 4 hours. Finally the robot could walk 8.6m in 15s with average body speed of 0.57m/s and period of 0.41s for each step. Fig. 4 shows the average and best fitness values during these 10 generations.



**Fig. 4.** PSO convergence for previous TFS

Using the new approach presented in this paper, after 3300 trails and 5 hours from starting PSO in a machine with the same specifications, the robot could walk 11.5 m in 15 s. Average and best fitness are shown in fig. 5.
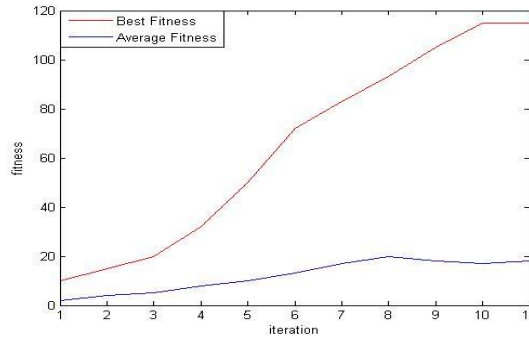


**Fig. 5.** PSO convergence for new approach and TFS model

Robot could walk with average body speed of 0.77 m/s by using the TFS with arm swing and increasing speed technique. The learned trajectories of left hip and knee after robot started to walk are shown in Fig. 6.A. It is determined that robot increased its speed in 0.20m/s. The learned trajectory of left arm is also shown in Fig. 6.B. and finally, the learned trajectory of left hip in coronal plane for abduction and adduction movements is also shown in Fig. 6.C.



**Fig 6. A)** Left Hip and Knee trajectories in the sagittal plane; **B)** Left Arm trajectory; **C)** Left hip trajectory in coronal plane (hipY)

After the learning procedure, the robot could walk with the average speed of 0.77 m/s. The best results of walking behavior for the teams that participated in RoboCup 2008 were chosen for the comparison [24]. Table 1, presents the comparison of the

best results of RoboCup teams, compared with the proposed approach. Analysing the table it is clear that forward walking achieved by our approach, may outperform the same skill of all teams analysed except SEU.

**Table 1.** Comparsion the average speed for forward walking in different teams (m/s)

|  | Proposed approach | FCPortugal | SEU | Wright Eagle | Bats |
|---|---|---|---|---|---|
| Forward Walking | 0.77 | 0.51 | 1.20 | 0.67 | 0.43 |

## 7. Conclusions

This paper presented a model with 9 parameters for producing all walking angular trajectories that uses all joints' movements in coronal and sagittal plane. An optimized Truncated Fourier Series is used to produce leg movements in sagittal plane and a model for swinging arms. A new model was also used to generate legs walking movement in coronal plane. We are able to increase the speed and stability of the robot's walking when compared to previously TFS model by using this model and the method mentioned in sec. 4 which was used both in sagittal and coronal motion.

According to the fact that this approach is model free and based on robot learning, it is capable of being used on all kinds of humanoid robots with different specifications. In future works we would like to expand this model to produced turn motion and improve the approach so that the robot can walk in any direction.

## References

1. Vukobratovic, M., Borovac, B., Surdilovic. D.: Zero-moment point proper interpretation and new applications. In: 2nd IEEE-RAS International Conference on Humanoid Robots, pp. 237-244 (2001)
2. Kajita, S., Kanehiro, F., Kaneko, K., Yokoi, K., Hirukawa. H.: The 3D linear inverted pendulum mode A simple modeling for a biped walking pattern generation. In: Proceedings of the 2001 IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 239–246 (2001)
3. McGeer, T. :Passive dynamic walking. International Journal of Robotics Research, vol. 9(2), pp. 62--82 (1990)
4. Pinto C. ,Golubitsky, M.: Central Pattern Generator for Bipedal locomption. J. Math. Biol. 53, pp. 474--489 (2006)
5. Mochon, S., McMahon, T.A.: Ballistic walking. J. Biomech. 13, pp. 49--57 (1980)
6. Buchli, J., Iida, F., Ijspeert, A.J.:Finding Resonance: Adaptive Frequency scillators for Dynamic Legged Locomotion. In: Proceedings of the 2006 IEEE/RSJ International Conference on Intelligent Robots and Systems. pp. 3903--3910 (2006)
7. Matsuoka, K.: Sustained oscillations generated by mutually inhibiting neurons with adaptation. Biol. Cybern. 52, pp. 367--376 (1985)

8. Yang, L., Chew, C.M., Poo, A.N.:Adjustable Bipedal Gait Generation using Genetic Algourithm Optimized fourier Series Furmulation. In: Proceedings of the 2006 IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 4435--4440 (2006)
9. Shafii, N., Javadi, M.H., Kimiaghalam B.: A Truncated Fourier Series with Genetic Algorithm for the control of Biped Locomotion. In: Proceeding of the 2009 IEEE/ASME International Conference on advanced intelligent Mechatronics, pp. 1781--1785, (2009)
10. Shafii, N., Khorsandian, A., Abdolmaleki, A., Jozi, B.: An Optimized Gait Generator Based on Fourier Series Towards Fast and Robust Biped Locomotion Involving Arms Swing. In: proceeding of the 2009 IEEE International conference on Automation and Logistics, pp. 2018 -- 2023 (2009)
11. Boedecker, J.: Humanoid Robot Simulation and Walking Behaviour Development in the Spark Simulator Framework. Technical report, Artificial Intelligence Research University of Koblenz (2005)
12. Kagami, S., Mochimaru, M., Ehara, Y., Miyata, N., Nishiwaki, K., Kanade, T., Inoue, H.: Measurement and comparison of humanoid H7 walking with human being. Robotics and Autonomous Sys, vol. 48, pp. 177--187 (2003)
13. Konz, R.J., Fatone, S., Stine, R.L., Ganju, A., Gard, S.A.,Ondra, S. L.: A Kinematic Model to Assess Spinal Motion During Walking, Journal of Biomechanics, vol. 31, pp. E898-E906, (2006)
14. Elftman, H.: The function of the arms in walking. J. Hum. Biol. 11, pp. 529--535 (1939)
15. Collins, S. H., Wisse, M., Ruina, A.: A 3-D passive dynamic walking robot with two legs and knees. Int. J. Robot. Res. 20, pp. 607--615 (2001)
16. Murray, M.P., Sepic, S.B., Barnard, E.J.: Patterns of sagittal rotation of the upper limbs in walking. In: *Physical Therapy* , vol. 47, pp. 272--284 (1967)
17. Thewlis, D., Richards, J., Hobbs, S.: The Appropriateness Of Methods Used To Calculate Joint Kinematics. Journal of Biomechanics, vol. 41, pp. S320--S320 (2008)
18. Shafii N., Aslani, S., Nezami, O.M., Shiry ,S.: Evolution of Biped Walking Using Truncated Fourier Series and Particle Swarm Optimization. In: CD proceeding of Robocop Symposium (2009)
19. Shi Y., Eberhart R.C.: Parameter selection in particle swarm optimization. In: Evolutionary programming VII proceedings of the seventh annual conference on evolutionary programming, pp. 591--600 (1998)
20. Jiao, B., Lian, Z., Gu X.: A dynamic inertia weight particle swarm optimization algorithm. j Chaos, vol. 37, pp. 698--705 (2008)
21. Beyer, H.G.: Evolutionary algorithms in noisy environments: theoretical issues and guidelines for practice. Comput. Methods Appl. Mech. Engrg., vol. 186, pp. 239--267 (2000)
22. Coello, C. A., Pulido, G. T.,Lechuga, M.S. :Handling Multiple Objectives With Particle Swarm Optimization. IEEE Transaction on Evolutionary Computation, vol. 8 , pp. 256--279, (2004)
23. Salman, A., Ahmad, I., Al-Madani, S.,: Particle Swarm Optimization For Task Assignment Problem, Microprocessorsand Microsystems, vol. 26, pp. 363--371 (2002)
24. Picado, H., Gestal, M. , Lau, N., Reis, L.P., Tomé, A.M., Automatic Generation of Biped Walk Behavior Using Genetic Algorithms. In  J.Cabestany et al. (Eds.), 10th International Work-Conference on Artificial Neural Networks, IWANN 2009, Springer, LNCS Vol. 5517, pp. 805--812, Salamanca, Spain, June 10-12, (2009)

# Protocols and Services

# Improving the Scalability of IEEE 802.11s Networks with a DHT-based Routing Protocol

Silvio Sampaio[1] * , Francisco Vasques[1], Pedro Souto[2], and Marcos Pinheiro[3]

[1] IDMEC, FEUP, University of Porto, Rua Dr. Roberto Frias,
4200-465, Porto, Portugal
{silviocs, vasques}@fe.up.pt

[2] ISR, FEUP, University of Porto, Rua Dr. Roberto Frias,
4200-465, Porto, Portugal
pfs@fe.up.pt

[3] DIMAP, Federal University of Rio Grande do Norte, Campus Universitario Lagoa Nova,
59078-970, Natal/RN, Brazil
marcos@dimap.ufrn.br

**Abstract.** IEEE 802.11s mesh network is gaining popularity among researchers and wireless equipment vendors. However, it still needs to improve its mechanisms to deal with scalability issues. In this paper we addressed the scalability problem in IEEE 802.11s mesh networks under two perspectives: (i.) capacity to handle a large amount of STAs information with a low overhead; (ii.) small number of communication hops. As proposed solution, we presented the DHT-based Cluster Routing Protocol (DCRP), a routing protocol based on DHTs, clustering of nodes and use of proxies. DCRP allows to archive (i.) by keeping STAs information in a DHT and avoiding flooding it over the network. Additionally, its modified forwarding mechanism that jointly use the data stored on DHTs and the cluster approach allow to archive (ii.). A preliminary performance evaluation, comparing the number of TC messages generated by DCRP vs. OLSR, is presented.

**Key words:** Distributed Hash Table, Routing, Mesh Networks, IEEE 802.11s

## 1 Introduction

Wireless Mesh Networks (WMNs) has got much attention for diverse applications such as office, campus/public access, residential automation, public safety, military, and industrial. WMN is a generic term to define a communication network build of

---

wireless nodes in a mesh topology. In a mesh topology each node act as an independent router, forwarding packets on behalf of other nodes. This type of wireless network is characterized by dynamic self-organization, self-configuration and self-healing. These properties enable, among other, fast deployment, low installation cost, and reliable communication. WMNs reduce the time and work required for creating or updating an existing wireless network, as the mesh nodes can dynamically cooperate to update/re-arrange the network.



**Fig. 1.** Classic 802.11 WLAN.

Classic WLANs, shown in figure 1, present a "wireless paradox" as each wireless *Access Point (AP)* need to be connected to a wired network in order to allow communication over the Extended Service Set (ESS). A mesh AP requires only to be connected to a power line. Another important difference is the communication among access points. In classic WLANs this communication is made through the wired network, as shown in figure 1. In a WMN each AP become a mesh point and, in a mesh fashion, it can establish mesh links with any other visible mesh point, as shown in figure 2. It enables multiples communication paths for data transmission in an ad-hoc fashion.

This work is focused in a particular type of Mesh Networks: the IEEE 802.11s Wireless Mesh Networks, which are the developing IEEE standard for *IEEE 802.11 Wireless Local Area Networks (WLANs)* based mesh networking. The IEEE 802.11s standard is detailed in session 2.1.

**Fig. 2.** Mesh 802.11 WLAN.

Typical deployed IEEE 802.11 WLANs consist of a series of wired APs that rely on a wired infrastructure to extend its connectivity. As result, the dimension of the network is largely restricted by the wired infrastructure. Conversely, in a IEEE 802.11s Mesh Network, APs can be interconnected in a multi-hop fashion. Thus, the IEEE 802.11s approach allows the setup of large sized networks (with a larger number of nodes) covered by a large number of APs. However, as the network size increases, the mesh network may face severe scalability problems. In essence, the available throughput will decay as the network gets bigger. One reason is the increasingly large amount of STAs association information that should be handled by the mesh APs. In a IEEE 802.11s mesh network, legacy stations or just STAs (802.11 devices), which do not support a mesh path selection mechanism, are associated to a mesh AP in order to have access to the mesh network. Then, the mesh AP have to handle the association information of STAs. Another reason for this behavior is the increase of the number of hops in the multi-hop network. The longer hop distance in a path will lower down the available throughput over the relay links.

The main target of this paper is to propose a novel scheme for path selection and message forwarding in IEEE 802.11s networks, that (i.) enables to handle a large amount of STAs information with a low overhead and (ii.) reduce the number of communication hops in order to increase the scalability of IEEE 802.11s networks. Preliminary results of this work already allow to state that (i.) can be achieved by the joint use of (Distributed Hash Tables) DHTs to store and quickly find information about the nodes. We also strongly believe that as our approach applies a clustering

technique and a novel forwarding mechanism that makes use of the data stored at the DHTs, it will reduce the number of hops on the multi-hop network. Thus (ii.) will be achieve.

The remainder of this paper is organized as follows. In Section 2 we introduced the most relevant concepts for the understanding of the proposed approach. An overview of the proposed scheme is discussed in Section 3. In Section 4, the proposed scheme is compared to state-of-art solutions that can be found on the literature. Preliminary results from a simulation assessment are presented in Section 5, and finally some concluding remarks are given in Section 6.

## 2 Background

In this section we briefly outline the main concepts used in our approach. First, the IEEE 802.11s standard is discussed. Then, the Distributed Hash Tables mechanism is summarized. After that, the path selection protocol *Radio Aware Optimized Links State Routing* (RA-OLSR) and its mechanisms is discussed. To finish, the clustering technique applied to create the clusters is presented.

### 2.1 IEEE 802.11s standard

The scope of the IEEE 802.11s Task Group is to extend the IEEE 802.11 architecture and protocol for providing the functionality of an Extended Service Set Mesh, i.e., access points able of establishing wireless links among each other that enable an automatic topology learning and dynamic path configuration.

Although the 802.11s standardization is still in progress, the most recent draft is already quite stable. In addition to generic IEEE 802.11 mechanisms, the IEEE 802.11s draft also addresses issues related to the MAC protocol, security and routing, where the latter is the subject of greater relevance [1, 2].

In a IEEE 802.11s Mesh Network is possible to find tree types of nodes: *Mesh Points (MPs)*, *Mesh Access Points (MAPs)*, and *Mesh Portal Points (MPPs)*. A MP is an IEEE 802.11s device that participates in the mesh routing process and can forward frames on behalf of other MPs in an ad-hoc way. In the IEEE 802.11s draft standard some Mesh Points that have additional Access Point functionality are called Mesh Access Points (MAPs). A MAP allows to support other wireless *Legacy Stations (STA)*, e.g. IEEE 802.11b/g/n devices, acting as bridge between the STA and the mesh network. Some other special MPs can act as Portal between the mesh network and other IEEE 802 networks, usually wired Ethernet networks. These nodes are called Mesh Portal Points and allow the extension of the mesh network coverage. Figure 3 illustrate the relationship between the different types of nodes in a mesh network. An interesting survey on WMNs can be found in [3].

**Fig. 3.** Elements of a 802.11s Network.

IEEE 802.11s Mesh Networks use a multi-hop wireless relaying infrastructure, where all nodes cooperatively maintain network connectivity. Data can be routed from the source node to the destination node using multi-hop communication. In a IEEE 802.11s Mesh Network, routing is performed at the data link layer and is given the name of *path selection*. Prior to draft version 1.06, every MP supported two routing protocols: the *Hybrid Wireless Mesh Protocol (HWMP)* [4] as the default routing protocol and the *Radio-Aware Optimized Link State Routing Protocol (RA-OLSR)* [5] as an optional one. Since draft version 1.07, the RA-OLSR was removed from the IEEE 802.11s specification. HWMP can work in both reactive and proactive modes. In reactive routing the route discovery is performed on-demand. In proactive routing, performed only on MPPs, a distance vector tree is used to avoid unnecessary routing path discovery and recovery messages. However, both HWMP and RA-OLSR have several shortcomings, namely in what concerns the scalability of the network. In HWMP, the proactive mode is centralized and constrained by the root node. Even when two MPs near for each other need to communicate, the proactive routing protocol routes the frames through the root node, which results in poor performance. At the same time, the reactive (on-demand) mode will initiate a path discovery process to search for a optimized path before sending the frames, resulting in an excessive number of broadcast messages. The problem with RA-OLSR is the overhead of control messages, even when the Fisheye protocol is used.

## 2.2 *Distributed Hash Table (DHT)*

A Distributed Hash Table (DHT) enables the application of the hash table concept to a distributed environment. It allows the efficient recovery/publication of data, through the association of a key to each data element. DHT networks have gained popularity as they are the underlying support for the organization of Peer-to-Peer (*P2P*) networks, such as Chord [6] and Pastry [7]. Basically the DHT uses a space of identifiers to guide the resource allocation process, where a resource can be a process to be executed or a information to be stored. The space of identifiers is divided among the elements that form the DHT and the resources are mapped into that space, typically using a hash function. Each network node is the responsible for all the resources mapped by its identifiers. Such identifiers are also referred as *keys*.

Obviously, the nodes participating in the DHT use a physical communication network, such as the Internet, to exchange messages. However, they also create a new network, superimposed upon this physical network, called the *overlay network*. This overlay network has its own topology and routing protocols that are specified by the DHT. Moreover, DHTs are typically multi-hop networks, where each node forwards the messages to the nodes that are nearest to their destination addresses.

## 2.3 RA-OLSR

The RA-OLSR protocol is a proactive, link-state wireless mesh path selection protocol based on the *Optimized Link State Routing (OLSR)* protocol [8]. It also include extensions like the *Fisheye State Routing (FSR)* protocol [9], and the use of radio-aware metrics for forwarding path computation and multipoint relay set selection.

OLSR protocol provides an optimized flooding mechanism based on MultiPoint Relay (MPR), used to diffuse topology information. MPR flooding optimizes flooding by minimizing the redundant retransmissions of Topology Control (TC) messages as the set of MPRs relays is a small set of neighbors through which a sender can reach all two hop neighbors. Another optimization can be done by using the FSR technique. In FSR protocol, information about closer nodes is exchanged more frequently than it is done about further nodes.

The RA-OLSR protocols also include an Association Discovery and Maintenance protocol to support non-mesh STAs both internal (associated with MAPs) and external (connected through MPPs).

The base mechanism of this protocol is the following: mesh points diffuse the whole set of mesh clients associated to themselves. It works in a proactive fashion, similar in spirit to the topology information exchange of OLSR: in both cases the information messages must be refreshed within a guaranteed interval. However, in addition, in case of topology/association change, this mechanism allows faster updates.

### 2.4 Clustering

In DCRP, is used a modified version of the *Efficient Clustering Scheme (ECS)* [10]. The clustering module is implemented as an OLSR plugin. To avoid sending signaling packets with medium contention, all cluster formation/maintenance messages are piggybacked into HELLO advertisements of the routing protocol. When a node receive a cluster message, it will analyze the membership cluster ID of the source node and attach this information in the OLSR neighbor database, where it can be later used for the forwarding process. Therefore, the clustering technique helps to reduce the flooding during route discovery phase, since the broadcast of routing messages can be constrained within the cluster.

## 3 DCRP

The DHT-based Cluster Routing Protocol (DCRP) integrates clustering with DHTs to enhance the scalability of routing in 802.11s networks. Clustering allows for the use of hierarchical routing and therefore to reduce the amount of routing traffic. The routing information that is not exchanged through the routing protocols is kept in DHTs, which supports a rather efficient access.



**Fig. 4.** DCRP architecture elements.

As shown in Figure 4, MPs physically close are grouped in clusters. Most MPs in a cluster communicate only with MPs in the same cluster, whereas a few MPs in a cluster communicate both with MPs in the same cluster and MPs in other clusters. We called the latter *border MPs (bMP)* and the former *internal MPs (iMP)*. All

MPs in a cluster maintain an intra-cluster routing table by executing an intra-cluster routing protocol. In addition, the bMPs of each cluster maintain an inter-cluster routing table by executing a mesh-wide inter-cluster routing protocol.

In order to reduce the traffic generated by the intra-cluster routing protocol, this protocol exchanges routing information only for MPs. Routing information relative to stations in the cluster, i.e. the proxy-MAP a station is associated with, is kept in a DHT. For each cluster, there is an *intra-cluster DHT (intra-DHT)*, whose nodes are the MPs of that cluster. Each MP is identified in the DHT by a random unsigned integer, $id$, which is the result of applying a hash function to its MAC address. Thus to learn the $id$ of a node, it is enough to know its MAC address. The DHT entry with routing information for a given station is stored in the MP whose $id$ is closest to the $id$ of that station. We called this MP the *intra-cluster key MP (intra-kMP)* of that station.

Likewise, the inter-cluster routing protocol exchanges routing information regarding only bMPs. Routing information relative to other nodes, i.e. stations and iMPs, is kept in a mesh-wide DHT, the *inter-cluster DHT (inter-DHT)*. As the inter-routing protocol is aware only of bMPs, each entry in this DHT must contain the MAC address of a bMP of the cluster to which the corresponding node belongs. We called this bMP the *proxy-bMP* for the node. The DHT entry for a given node is stored in the bMP whose $id$ is closest to this node id. This bMP is known as the *inter-cluster key MP (inter-kMP)* of that node.

As there is a one-to-one map between a routing protocol instance and a DHT, each node of a DHT knows the remaining nodes of that DHT, and therefore the DHT overlay network is a fully connected graph, i.e. the set of neighbors of a given node comprises all other DHT nodes. In order to allow routing of frames in the DHT, each DHT node maintains a DHT *neighbor table (NT)*, which just maps the $id$s of the DHT nodes to their MAC addresses. As the $id$ of a node can be determined from its MAC address, the maintenance of NT is for free: it is provided by the corresponding instance of the routing protocol. Insertion of the entries in the DHT is also rather efficient and straightforward: when a station associates with a MAP, the latter inserts an entry for that station both in the intra-cluster DHT, for the local cluster, i.e. in its intra-kMP, and in the inter-cluster DHT, i.e. in its inter-kMP. Note that most likely the two entries are different and are inserted in different nodes.

In addition to the already mentioned tables, an MP also keeps a *proxy cache (P-cache)*, which maps the MAC address of a station to the MAC address of a proxy-MP of that station: the proxy-MAP for MPs in the same cluster and the proxy-bMP for MPs in other clusters. The purpose of this cache is to avoid DHT-routing, which is very likely less efficient than physical routing.

Therefore, DCRP always tries to forward frames using physical routing information. Only if no information is available, it resorts to DHT-based routing, i.e. it forwards the frames to a kMP.

To populate the proxy cache, when a kMP receives a frame routed through the DHT and it has an entry for the frame's final destination in its DHT, it sends back a *redirect* message to the node at the other end of the DHT-hop, which in turn may forward it back to the other end of the previous DHT-hop, if any. This redirect message contains the address of the destination of the frame forwarded by the kMP and the MAC address of the proxy MP (either the proxy-MAP or the proxy-bMP) of the destination.

## 4   Related Work

The DCRP relies on the use of proxies, DHTs and clustering. Although these concepts are widely used in many different scenarios including mesh networks, in this work we proposed to use them in a integrated way, exploring their synergies, in order to obtain a more scalable routing protocol for IEEE 802.11s WMNs.

The proxy concept is used both in HWMP [4] and in RA-OLSR [5] protocols to reduce the amount of information required to route frames to non-mesh stations. In both protocols the discovery of a station's proxy relies on broadcast messages. HWMP uses broadcast messages for route discovery, whereas RA-OLSR uses it for route announcement. Our protocol does not require broadcast messages for proxy discovery, instead it uses DHT to maintain this information. Furthermore, the DCRP protocol extends the proxy concept to clusters, by means of bMP proxies.

Other works that use DHT on the wireless routing process leave the task of determining the route to the DHT. In [11–13] no additional routing protocol is needed, whereas in [14–16] the DHT is used to indicate the next node to which the frame must be send. As this next node may not be reachable within 1-hop, in the later case, an additional path selection protocol is required, e.g. HWMP, to discovery the route to this node. In contrast with DCRP, which was designed to be easily integrated with the IEEE 802.11s standard, those protocols require the modification of the IEEE 802.11s frame, and force the use of a specific routing protocol. As illustrated above, the DCRP uses the IEEE 802.11s frame format. Furthermore, although in DCRP description RA-OLSR was used for both intra and inter-cluster routing, it is clear that DCRP can be used with other proactive routing protocols.

Clustering is widely regarded as an effective approach to increase the scalability of WMNs. Several works, e.g. [17, 18] and [19] have proposed modified OLSR versions with clustering. Among other, a key difference between these proposals concerns inter-cluster frame forwarding. In [18], the authors propose the use of special nodes with multiple radios, called *cluster heads*, that are assumed to be able to communicate directly with neighbor cluster heads using one of their radios. In contrast, in [17] and [19], as well as in DCRP, nodes need only one radio and the transmission among neighbor clusters is done trough border nodes. However, whereas

in those proposals the exchanged messages among clusters include also information about internal nodes, in DCRP the inter-cluster protocol exchanges information about border MPs only. This is possible, due to the concept of proxy bMP and the use of DHTs. As a result, the DCRP generates much less inter-cluster routing traffic than other proposals.

## 5    Preliminary Evaluation and Results

In order to provide a rough assessment of the DCRP performance, an initial set of simulations was done based on the comparison between the number of *Topology Control (TC)* messages generated by DCRP vs. OLSR. The purpose of such simulations is to provide preliminary estimates of the overhead caused by broadcast messages in DCRP, when compared to OLSR. All simulations were carried out using the *Network Simulator 3 (ns-3)*.

The topology of the used network was a square grid of NxN nodes (MPs). The distance between neighbor nodes in the horizontal and vertical directions was 100m and constant for the entire grid. Furthermore, the radio ranges of all nodes were set so that a node can communicate in one hop only with its neighbors in the horizontal and vertical directions. As mentioned above, there was no stations (STAs) involved. In our experiments we considered different values for N (5 to 30, in steps of 5) and ran the model for 5 minutes, 5 times, for each value of N.

For the analysis of DCRP, we partitioned the grid in clusters of 25 nodes in a square grid of 5x5 nodes, with 4 bMPs. The size of the cluster was chosen taking into account that the 802.11s draft standard was designed for meshes up to 32 MPs. Under these assumptions, the routing traffic in number of messages generated by the DCRP protocol can be estimated by:

$$N_{DCRP} = C \times M_c + I_c \times H \tag{1}$$

where $C$ is the number of clusters, $M_c$ is the number of TC messages sent by the OLSR protocol in a 5x5 cluster, $I_c$ is the number of TC messages generated by a grid with as many nodes as bMPs in the mesh, and $H$ is the average number of hops between bMPs within a cluster, i.e. the *propagation delay* metric. Thus the first term in 1 is the total routing traffic generated by all instances of the intra-cluster routing protocol, whereas the second term is an estimation of the total routing traffic generated by the inter-cluster routing protocol. The factor $H$ is used to consider that communication between bMPs in the same cluster may require more than one hop.

The Figure 5 compares the results obtained with the ns-3 simulation model for the OLSR protocol and those obtained for DCRP (1). In computing the value for (1) we used for $M_c$ and $I_c$ values obtained also with the ns-3 simulation model. As for $H$ we used the value of 5, as in our simulation we have used a 5x5 grid.

**Fig. 5.** Comparison of the number of messages generated by the routing protocol.

As we would expect, DCRP generates a lower number of TC messages than OLSR. The results are promising, but the used models are still rather approximate. Furthermore, as we stated, these experiments consider only the clustering aspect of the DCRP. The benefits of the DHT-based service have yet to be evaluated.

## 6 Conclusions and work in progress

In this work we have proposed a path selection protocol that increases the scalability of 802.11s WMNs. We partitioned the network in clusters and use proxies to allow the communication within and among the clusters: one type of proxy is used to allow the communication from non-mesh STA nodes trough the mesh network and another type used by internal MAPs to communicate among clusters. We have also defined a DHT approach to find the proxy for a node without use of broadcast messages. Theses combined features allow each cluster to run its own routing protocol instance where only internal nodes are affected. The intra-cluster communication is done by another routing protocol instance that only takes into account the border nodes.

The key impact of this work is that it can meaningly reduce the number of routing protocol messages. The main reason for this behavior is that internal route

changes within a cluster will not be broadcasted to other clusters. Another improvement is the reduction of the message size as the messages exchanged among clusters will only contain information about the border nodes.

The obtained results from a preliminary evaluation already indicate that the DCRP protocol may be an adequate solution to deal with scalability issues in WMNs.

Presently, we are setting up a full ns-3 model in order to make the performance assessment of the DCRP. Such model encompasses a RA-OLSR version on ns-3 and a complete simulation model of DCRP. The performance metrics that will be analyzed include: network throughput and access delay. The first preliminary results are quite encouraging, but at the moment of writing of this paper, we still do not have a full set of results available to be presented.

In order to complete the performance and applicability assessment of DCRP, future works includes a field trial that will be conducted over real world applications scenarios. Because the lack of commercial IEEE 802.11s mesh devices will be necessary to implement our solution over traditional Wi-Fi devices. We expect to be able to load a modified firmware with a DCRP implementation.

## References

1. IEEE P802.11s/D1.10: Draft Amendment to Standard for Information Technology - Telecommunications and Information Exchange Between Systems - LAN/MAN Specific Requirements - Part 11: Wireless Medium Access Control (MAC) and physical layer (PHY) specifications: Amendment: ESS Mesh Networking. (March 2008)
2. IEEE P802.11s/D3.0: Draft Amendment to Standard for Information Technology - Telecommunications and Information Exchange Between Systems - LAN/MAN Specific Requirements - Part 11: Wireless Medium Access Control (MAC) and physical layer (PHY) specifications: Amendment: ESS Mesh Networking. (March 2009)
3. Akyildiz, I.F., Wang, X., Wang, W.: Wireless mesh networks: A survey. Computer Networks **47**(4) (2005) 445–487
4. Bahr, M.: Update on the hybrid wireless mesh protocol of ieee 802.11s. 2007 IEEE Internatonal Conference on Mobile Adhoc and Sensor Systems, MASS (2007) 1–6
5. Hiertz, G., Max, S., Zhao, R., Denteneer, D., Berlemann, L.: Principles of IEEE 802.11 s. In: Computer Communications and Networks, 2007. ICCCN 2007. Proceedings of 16th International Conference on. (2007) 1002–1007
6. Stoica, I., Morris, R., Liben-Nowell, D., Karger, D.R., Kaashoek, M.F., Dabek, F., Balakrishnan, H.: Chord: A scalable peer-to-peer lookup protocol for Internet applications. IEEE/ACM Transactions on Networking **11**(1) (2003) 17–32
7. Rowstron, A., Druschel, P.: Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer systems. In: IFIP/ACM International Conference on Distributed Systems Platforms (Middleware). Volume 11., Heidelberg (2001) 329–350
8. Clausen, T., Jacquet, P.: IETF RFC-3626: Optimized Link State Routing Protocol OLSR. The Internet Society http://www. ietf. org/rfc/rfc3626. txt (2003)

9. Gerla, M., Hong, X., Pei, G.: Fisheye state routing protocol (FSR) for ad hoc networks. IETF Draft (2002)

10. Yu, J.Y., Chong, P.H.J.: An efficient clustering scheme for large and dense mobile ad hoc networks (MANETs). Computer Communications **30**(1) (2006) 5–16

11. Zahn, T., Schiller, J.: DHT-based unicast for mobile ad hoc networks. Proceedings - Fourth Annual IEEE International Conference on Pervasive Computing and Communications Workshops, PerCom Workshops 2006 (2006) 179–183

12. Fuhrmann, T., Di, P., Kutzner, K., Cramer, C.: Pushing Chord Into the Underlay Scalable Routing for Hybrid MANETs. Univ., Fak. für Informatik, Bibl. (2006)

13. Caesar, M., Castro, M., Nightingale, E.B., O Shea, G., Rowstron, A.: Virtual ring routing: Network routing inspired by dhts. Computer Communication Review **36**(4) (2006) 351–362

14. Zahn, T., Schiller, J.: MADPastry: A DHT Substrate for Practicably Sized MANETs. In: Proc. of ASWN. (2005)

15. Takeshita, K., Sasabe, M., Nakano, H.: Mobile P2P Networks for Highly Dynamic Environments. Pervasive Computing and Communications, 2008. PerCom 2008. Sixth Annual IEEE International Conference on (March 2008) 453–457

16. Baccelli, E., Schiller, J.: Towards scalable manets. ITS Telecommunications, 2008. ITST 2008. 8th International Conference on (Oct. 2008) 133–138

17. Ros, F., Ruiz, P.: Cluster-based OLSR extensions to reduce control overhead in mobile ad hoc networks. In: Proceedings of the 2007 international conference on Wireless communications and mobile computing, ACM New York, NY, USA (2007) 202–207

18. Ge, Y., Lamont, L., Villasenor, L.: Hierarchical OLSR-a scalable proactive routing protocol for heterogeneous ad hoc networks. In: IEEE International Conference on Wireless And Mobile Computing, Networking And Communications, 2005.(WiMob 2005). Volume 3. (2005)

19. Baccelli, E.: OLSR Scaling with Hierarchical Routing and Dynamic Tree Clustering. In: IASTED International Conference on Networks and Communication Systems (NCS),(Chiang Mai, Thailand). (2006)

# Secure communications: The IPsec role

Jorge Pinto Leite

Faculty of Engineering, University of Porto,
Rua Dr. Roberto Frias, s/n, 4200.- 465 Porto, Portugal
pro09015@fe.up.pt

**Abstract.** Secure communications is a "must have" in today's world. The threats are huge and no one can assume that the data sent or received is private and secure and, on the other side of the communication channel, exactly as it originally was. Authenticity of received data should also be assured. Having this idea in mind as well as the enormous amount of low cost – yet insecure – communication links, a protocol to provide data security is of vital importance. The security it provides should be sufficient to comfort any user and keep him confident when communicating with a third part by encrypting and authenticating the data. Among several protocols for this subject, IPsec is one of the most used and well known for secure communications, but configuring it is not an easy task. And is it worthwhile? In this paper we pretend to show the role of IPsec as an important player of secure communications domain and its behavior when used. We will show that IPsec deserves the confidence that everyone is putting on it, in spite of some issues it has.

**Keywords:** IPsec, IP, secure communication, AH, ESP, tunnel, transport.

## 1 Introduction

It is impossible to imagine the world today without computers and communications. Computer usage is increasing rapidly and it is expected to grow up to approximately two thousand millions units worldwide in 2013 [1], according to a forecast of Computer Industry Almanac, shown in Figure 1. This forecast did not include neither embedded computers (which are used to control all types of electronic and electromechanical products) nor smartphones or PDAs.

And what kind of tasks are these computers doing? It is well known – we feel it! – that more and more services and organizations are moving to the cyber world, and using Internet to buy anything or to plan and buy tickets for travelling is not unusual nowadays. This reality varies for each country, and the data for Portugal covering the period from 2005 to 2009 can be found in [2]. By pure common sense, we can suspect that privacy might be in risk. It is also usual to interact with our bank, the tax authority or many other private services as well as with our employer. Under these specific situations privacy is not the only thing we must be careful about. Data confidentiality and integrity must also be present in our minds. The proliferation of

computer access, mobile or not, to everyone reinforces the assumption that in the future more man-in-the-middle[1] attacks will likely happen [3], with inherent threats to all communications [4, 5].

One of the ways to deal with this reality, among others, is to provide a security protocol, and several have emerged. But at the same time, an IP (Internet Protocol) layered infrastructure was implemented, so a security protocol to operate at the network layer was valuable, because it would operate over standard protocols, it wouldn't imply any infrastructure change, and could be used by higher-layer applications.



**Figure 1 -** Computers in use by regions worldwide

IPsec is that protocol. It started being developed at same time as IPv6, but as IPv6 has taken years to develop and roll-out, IPsec design kept in mind its use with IPv4, as well as IPv6. In this paper we analyze the benefits and disadvantages of IPsec nowadays. It must be kept in mind that as IPsec operates in the network layer of the OSI model, it is assumed that the format of the IP datagram is known.

This paper is organized as follows: Section 2 presents and describes the IPsec protocol. Section 3 presents a lab test to verify the security that IPsec offers. Finally we present some conclusions as well as some of the topics we intend to do in the near future.


## 2   The IPsec protocol

In this section we will describe the IPsec protocol starting with his background, its principals, the actual situation and some threats and problems it has.

---

[1] The man-in-the-middle attack and techniques is described in detail in [4, 5]

## 2.1 History

The definition considered to be the first one of IPsec was published on November 1998 [6]. In fact there was a kind of previous definition [7] published on August 1985, but it was almost a "declaration of principles" and not the protocol definition by itself. IPsec stands for Internet Protocol Security, "…*an architecture to provide various security services for traffic at the IP layer, in both the IPv4 and IPv6 environments*" [5].

Linked to IPsec as part of its operation, several other protocols were published in separate RFC[2], as stated by IPsec definition. So, November 1998 also saw the publication of several others RFC's without which IPsec could just not be considered as complete [8, 9, 10, 11].

During the next years changes and improvements have been made to IPsec until the actual version, version 3, that we will cover later in this paper.

## 2.2 IPsec principles

First of all, let us introduce one of the most elementary principles of IPsec: the peer. Any conversation implies a sender and a receiver. When this conversation is secured with IPsec, the peer is the end point of the IPsec channel, so each IPsec implementation implies two peers, one for each end point of the channel.

Another important principle of IPsec is the Security Association (SA). This is a channel that is established between the peers within which flows the data that should be protected. In fact, if security must be provided in both directions, two SAs are established because IPsec SA is unidirectional [12]. The SA is limited by a grace period and amount of traffic. Whenever one of these limits is reached, the SA is discarded and a new SA is established, aiming to prevent attempts to decipher the flowing data. The Security Parameter Index (SPI), a 32 bit number that is chosen to identify a particular SA for any connected device, the IP destination address and the security protocol identifier that specifies the security protocol is use, are the three variables that characterize a specific SA.

IPsec protocol defines a Security Association Database (SAD) to store all existing SA, in order to allow IPsec to process the datagram to be protected with the correct mechanisms. Also stored are the Security Policies in a Security Policies Database (SPD). This policy describes how different datagrams should be handled, for example, which ones should be subject to IPsec or not. These two concepts seems similar, but the main difference is that security policies are general while security associations are more specific. IPsec first looks in SPD and applies the security association stored in SAD that SPD points to.

IPsec aims to provide security to the IP datagram. But it is known that IP datagram changes while flowing from sender to receiver – just think of the TTL (Time to Live) field. So, IPsec acts on all fields of the datagram except the ones that change during traffic.

---

[2] RFC stands for Request For Comments, containing technical and organizational documents about the Internet, including the technical specifications and policy documents produced by the Internet Engineering Task Force (IETF)

The desired security is provided by authenticating the data and encrypting it. Also, as a major aspect of security, IPsec can deal with replay, that is, the resending of a datagram (nevertheless, this is an optional configuration of IPsec).

To achieve this, IPsec acts using two function modes and two security protocols. The function modes are the transport and the tunnel mode, and the security protocols are the Authentication Header (AH) and the Encapsulating Security Payload (ESP).

## IPsec architecture

IPsec defines three different architectures for implementation, the Integrated, the "Bump In The Stack" and the "Bump In The Wire" architectures.

The Integrated Architecture is the ideal and most elegant one and implies that IPsec functionalities are integrated directly into IP itself. This architecture implies that layer 3 of the OSI model applies IPsec together with the encapsulation of the IP datagram.

Another one is sometimes referred as "Bump In The Stack" (BITS). The IP datagram is captured by IPsec that process it and passes the resulting IPsec datagram to the layer 2. This architecture is generally used for IPv4 hosts.

The third architecture is sometimes referred as "Bump In The Wire" (BITW). This one requires additional equipment that contains a security processor to act as gateway between the hosts and the link. That equipment would be responsible for IPsec implementation, turning itself into a security gateway, but has the obvious disadvantage of cost and complexity.

It must be noted that whichever of the architectures are used the results are the same.

## Transport mode

Transport mode is a way of functioning in which the IP header is kept as it is, meaning that the IP addresses of the sender and of the receiver are the real ones. This means that if someone is sniffing the link he will see both addresses, which may be considered insecure. However, it must be noted that in some situations it is not possible to masquerade IP addresses, so this mode of function must exist.

In Figure 2 we can see an example of transport mode operation. The layer 3 is named "IP / IPsec" because IPsec operates in this layer. Being IPsec turned on or off, no change is made to the normal construction of IP header that this layer does.

This is clearly an integrated architecture as layer 3 must apply AH/ESP (described below) as the original IP packaging is performed.
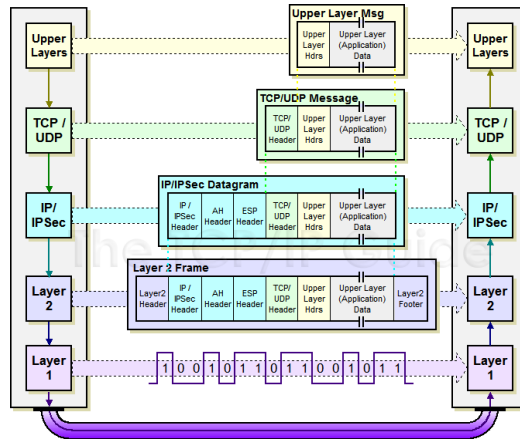
**Figure 2 -** Transport mode[3]

**Tunnel mode**

Tunnel mode looks like the most secure mode of IPsec operation, as in this mode the IP addresses of both sender and receiver are changed to different ones, so the mapping of internal addresses of sender and receiver networks becomes impossible to find out.

In this mode, it is almost as if a new layer is introduced between layer 3 and layer 2 of the OSI model, as we can see in Figure 3. In the usual layer 3, the normal IP datagram is created. But before layer 2 processes the datagram, IPsec treats the layer 3 output and a new datagram is created. This implies that all datagram – including the IP addresses of sender and receiver and the higher layer data – is encapsulated in the payload and a new IP header is created, containing the new IP header the addresses of the peers.

This represents an implementation of architectures BITS or BITW, as IPsec is applied after layer 3 has processed the higher layer message.

---

[3] Reference: http://www.tcpipguide.com/free/t_IPSecModesTransportandTunnel.htm, accessed on the 5th Nov 2009

**Figure 3 -** Tunnel mode[4]

## Authentication header (AH)

Authentication Header is a security protocol in which all or part of the IP datagram is authenticated, except for the fields that are changed during traffic. The authentication is made by applying a hashing algorithm and a shared key, known only by source and destination. A security association is established between the peers that specify these particulars so that the source and the destination know how to perform the computation but everybody else does not.

On the source device, AH performs the computation and puts the result – the Integrity Check Value, ICV – into a special header with other fields for transmission. The destination device does the same calculation which enables it to check if any of the fields in the original datagram was modified [13].

The operation is similar to the calculation of the CRC to handle error detection. Also as CRC, the AH and the ICV does not change anything in the original datagram, only appends a new header to it. The exact position of this header depends of the mode of operation and the version of IP (IPv4 or IPv6) [13].

Whichever the mode and the IP version, the AH header must be the first one after the datagram's final IP header and its length must be a multiple of 32 bits [13]. This header must contain the decimal value of 51 (fifty one) in its protocol (IPv4) or next header (IPv6) field.

---

[4] Reference: http://www.tcpipguide.com/free/t_IPSecModesTransportandTunnel.htm, accessed on the 5th Nov 2009

**Encapsulating Security Payload (ESP)**

While AH provides authenticity to datagram, Encapsulating Security Payload provides confidentiality and integrity[5]. This is achieved by encrypting the information sent. In the way that ESP operates, authentication can be also optionally achieved. ESP is computed by introducing three additional components to the datagram: the ESP Header, the ESP Trailer and the ESP Authentication Data.

The ESP Header contains two fields, the SPI (Security Parameter Index) and the Sequence Number. This header appears before the encrypted data. Whichever the mode and the IP version, the header that precedes ESP Header must contain the decimal value of 50 (fifty) in its protocol (IPv4) or next header (IPv6) field.

The ESP Trailer is placed after the encrypted data and contains padding that is used to align it to a multiple of 32 bits, through a Padding and Pad Length field. It also contains the Next Header field for ESP. This component is used not only to align the encrypted data but also allows an important feature as follows. This component is encrypted before the datagram is sent. As the Next Header is included here, the necessary data to tie all packets together is not visible to anyone that might be intercepting the conversation. Also, it becomes difficult to guess where the real data ends and the padding field begins. If this behaviour wasn't included the Next Header field would be visible to an attacker.

The ESP Authentication Data contains an ICV computed in a similar manner to AH, and is used when authentication is configured for ESP. The difference of ICV calculation when compared with the one built by AH is that this ICV is computed after the data has been encrypted, not before.

The symmetrical encryption algorithm used is specified by the SA. An interesting aspect of ESP is that it generates dummy data to avoid absence of data in the SA. This implies that if someone captures the packets that are flowing through a channel that is protected with ESP protocol, the split of real from dummy data is almost impossible, so the effort needed to make a successful brute force attack to decipher the data will be huge (more information on brute force attack can be found in [14]).

**AH and ESP in transport and tunnel modes**

It has already been said that AH and ESP can operate in both modes of operation of IPsec, being no surprise that the resulting datagram varies between the modes.
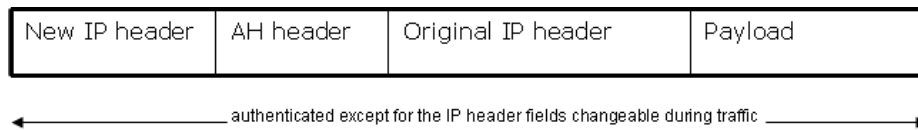


| Original IP header | AH header | Payload |

authenticated except for the IP header fields changeable during traffic

**Figure 4 -** Transport mode of AH

---

[5] On security purposes the terms *confidentiality* versus *privacy* versus *secrecy* can overlap. We share in this paper the assumptions on these terms described in [2].
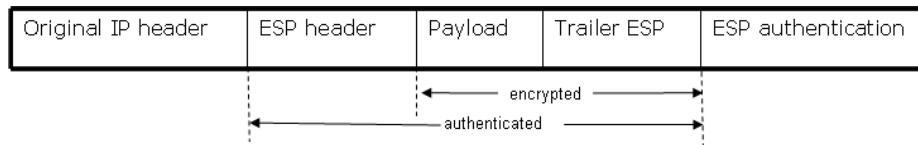
Figure 4 shows the resulting datagram of AH application in transport mode. As it can be noted, the AH header is placed between the original IP header and the payload, and all fields are authenticated except for the IP header fields that change during traffic.

In tunnel mode, the resulting datagram of an AH application is shown in Figure 5, where we can note that the AH header is placed immediately after the IP header, but this is a new one, not the original, that comes immediately before the payload.
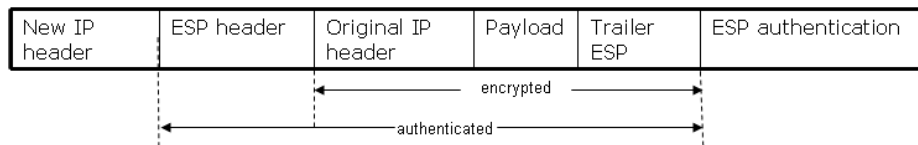
| New IP header | AH header | Original IP header | Payload |
|---|---|---|---|

←————————— authenticated except for the IP header fields changeable during traffic ——————————→

**Figure 5 -** Tunnel mode of AH

ESP principal concern is with integrity and confidentiality.. When applied in transport mode and with the optional authentication, the resulting datagram is as the one represented in Figure 6. Note the encryption of payload (that contains the higher layer data as well as the ESP Trailer) and the authentication of all data between the original IP header and the ESP authentication data.

| Original IP header | ESP header | Payload | Trailer ESP | ESP authentication |
|---|---|---|---|---|

←——— encrypted ———→
←——————— authenticated ———————→

**Figure 6 -** Transport mode of ESP

When ESP is used in tunnel mode, the final datagram that is passed to data link layer is the one represented in Figure 7. As with transport mode, all data between the (new) IP header and the ESP authentication data is authenticated, being the original IP header, the payload and the ESP trailer encrypted.

| New IP header | ESP header | Original IP header | Payload | Trailer ESP | ESP authentication |
|---|---|---|---|---|---|

←——— encrypted ———→
←——————— authenticated ———————→

**Figure 7 -** Tunnel mode of ESP

**Other aspects of IPsec**

IPsec can be used together with automatic key exchange, the IKE Protocol [15]. In spite of some changes in operation mode, our aim was to describe and test IPsec by itself, so IKE operation and functionality is not included here.

It must be noted that AH and ESP can operate independently or together, providing in this last mode almost all the security that can be achieved nowadays, authentication, confidentiality, integrity and anti-replay (although this functionality is optional).

### 2.3 IPsec version 3

The current version of IPsec is version 3, published in December 2005. In spite of being similar to the previous definitions, IPsec version 3 has a slight difference. It states that IPsec implementations may support AH [12]. It explains that "…*IPsec implementations MUST support ESP and MAY support AH. (Support for AH has been downgraded to MAY because experience has shown that there are very few contexts in which ESP cannot provide the requisite security services…*" [12]. Nevertheless, this definition still describes the two security protocols, AH and ESP, and their modes of operation. This version has no fundamental changes from the previous one, and focus especially on details of IPsec implementations (e.g., fragmentation of packets before applying IPsec) and implementation simplification. Other changes are due to a more detailed explanation in more detail of some concerns on security, the revoke of requirements of previous version, and minor updates due to high layer changes.

### 2.4 IPsec problems and limitations

As we can see from the previous sections of this paper, there are some disadvantages of IPsec. First off all, the overhead it implies can be a problem, especially in links with low bandwidth. Other problems lie on limitations of IPsec applications and are discussed in [16]. While agreeing with these authors, we consider that some of their conclusions aren't tied to IPsec but to higher layer applications. For example, the use of dynamic ports would be effectively an improvement to IPsec, but we consider that this functionality is not directly bounded to layer 3 of the OSI model. They also mention the lack of authorization in IPsec protocol. Again, we consider that this is not a function of the layer 3 stack.

Nevertheless, their conclusions should be considered but, in our point of view, not as an improvement to next version of IPsec but as complements to IPsec provided by other security protocols.

## 3   Experimental results

According to the previous description of functioning mode, the tunnel mode seems to be the most secure one. In spite of the MAY consideration of version 3,  IPsec was

tested using AH and ESP simultaneously to simulate an actual possible implementation. The decision on such a test scenario was supported by [17] that alerts that "…*Using encryption without a  strong integrity mechanism on top of it (either in ESP or separately via AH) may render the confidentiality service insecure against some forms of active attacks …*".

To achieve this a small network was built on a lab that has two Cisco® 2801 routers, two Cisco® Catalyst 2960 switches and two personal computers (PC). VLSM (Variable Length Subnet Mask) technology was used to interconnect all equipments. To capture the flowing packets, a third PC with its Ethernet card in promiscuous mode was inserted between the routers, with Wireshark® installed in it. The test network is represented in Figure 8, where the most representative IPv4 addresses are shown.

Wireshark® is the actual name of a packet sniffer application that was previously called Ethereal. Its purpose is to capture packets that are flowing in a channel and analyze them. Some additional features were introduced when the name changed, however the methods and its use are almost the same. More information about using Wireshark® can be found in [14].

The test was made by sending from PC with IP address 192.168.0.10/30 an ICMP message type 08, Echo request [18]. This was repeated without IPsec and with IPsec configured with the two security protocols in transport and tunnel mode.



**Figure 8 -** The testing network

All data was captured in the Wireshark® station and then analyzed. Without IPsec, the datagram is as expected, and shown in Figure 9 (further reading and analysis of the IP datagram can be found in [19]).

```
0000   00 1a a1 40 29 f5 00 1a   a1 f4 41 3f 08 00 45 00   ...@)... ..A?..E.
0010   00 3c 06 60 00 00 7f 01   b4 04 c0 a8 00 0a c0 a8   .<.`.... ........
0020   00 02 08 00 cc 5b 03 00   7e 00 61 62 63 64 65 66   .....[.. ~.abcdef
0030   67 68 69 6a 6b 6c 6d 6e   6f 70 71 72 73 74 75 76   ghijklmn opqrstuv
0040   77 61 62 63 64 65 66 67   68 69                     wabcdefg hi
```

**Figure 9 -** The unsecure packet

It can be noted here the IP header with the IP address of the sender (starting in position $1A_{16}$) and of the receiver (starting in position $1E_{16}$). In position $17_{16}$ – the protocol field – we can note the value 01, meaning that the protocol used is ICMP. Also note the ICMP message itself. It starts on position $22_{16}$ with the ICMP message

type ($08_{16}$, Echo request) with 1 byte length, followed by a 1 byte field for code ($00_{16}$), a 2 bytes field for checksum (CC $5B_{16}$), a 2 bytes field for the identifier (03 $00_{16}$), a 2 bytes field for sequence number (7E $00_{16}$) and finally the message data starting in position $2A_{16}$.

Then the routers were configured to secure the channel with IPsec using AH and ESP in transport mode and the sending was repeated, with result shown in Figure 10.

```
0000   00 1a a1 40 29 f5 00 1a   a1 f4 41 3f 08 00 45 00    ...@)... ..A?..E.
0010   00 6c 06 8c 00 00 7f 33   b3 76 c0 a8 00 0a c0 a8    .l.....3 .v......
0020   00 02 32 04 00 00 00 00   00 00 00 00 08 e2 9c        ..2..... ........
0030   91 da 00 0c 73 bd f8 46   a5 77 00 00 00 00 00 00    ....s..F .w......
0040   00 08 39 be 78 83 bb ab   0f a4 5b 5f ff 71 d1 e9    ..9.x... ..[_.q..
0050   bd e7 91 fd 84 47 48 25   13 8d 0b 49 09 e6 43 78    .....GH% ...I..Cx
0060   6d 5e c1 b6 b2 3a 1d fa   24 6a c9 d3 bd 88 63 8b    m^...:.. $j....c.
0070   a2 64 df c7 76 58 bc de   a5 11                       .d..vX.. ..
```

**Figure 10 -** IPsec in transport mode captured packet

The IP addresses of the sender and of the receiver remains exactly as without IPsec. But as AH and ESP are in use, some expected changes can seen.

First of all, let's take a look in the protocol field of IPv4 header. The IP header remains equal when comparing Figure 9 with Figure 10, except in position $17_{16}$, that contains, in Figure 10, the value $33_{16}$, i.e., the hexadecimal representation of the decimal value 51, meaning that the next protocol in the packet is AH. Its header starts in position $22_{16}$ (this position of Figure 9 contains the ICMP type of message) and starts immediately with the hexadecimal representation of the decimal value 50, $32_{16}$, meaning that next protocol will be ESP. The ESP header starts in position $3A_{16}$ with, in this case, a SPI of zeros and a sequence number of eight.

Let's look to the data itself. As the message sent was again the ICMP message "Echo request" it is certainly on the packet. The data from higher layers starts in position $42_{16}$. Comparing this field with the corresponding one from Figure 9, it is noted that the message is ciphered to what seems spurious data.

```
0000   00 1a a1 40 29 f5 00 1a   a1 f4 41 3f 08 00 45 00    ...@)... ..A?..E.
0010   00 7c 00 34 00 00 ff 33   39 bf c0 a8 00 06 c0 a8    .|.4...3 9.......
0020   00 05 32 04 00 00 d7 00   d6 8f 00 00 00 07 4f d9    ..2..... ......O.
0030   a7 75 07 1f b9 c1 ed d6   9a a5 79 e3 8b 3b 00 00    .u...... ..y..;..
0040   00 07 02 d7 22 d3 f4 ea   89 12 23 fe b5 4f fb c2    ...."... ..#..O..
0050   c9 4a 88 72 58 4b 94 0f   83 92 79 04 1e 16 35 a4    .J.rXK.. ..y...5.
0060   92 27 b0 a3 6a 44 47 cc   f9 94 48 8b a7 a1 c1 8c    .'..jDG. ..H.....
0070   81 92 43 e3 34 95 6d c0   cb 5d c2 c7 47 bd 8b 0a    ..C.4.m. .]..G...
0080   0f 80 4e cc 80 d9 71 cd   f4 81                       ..N...q. ..
```

**Figure 11 -** IPsec in tunnel mode captured packet

Finally IPsec was configuration was changed for tunnel mode (Figure 11). While the AH header and the ESP header remains unchanged, except for the SPI and the sequence number, it is interesting to note that the source and destination IP addresses have been switched to the external addresses of the routers, that are, in this design, the peers of the SA. As with the simulation in transport mode of operation, the ICMP message starts in position $42_{16}$, but again neither the 8 byte fields of message type, code, checksum, identifier and sequence number nor the message payload is understandable. This is due to the ESP encryption of to the higher layers data.

One of the issues mentioned above can be seen in these experimental results: the overhead. The increase of data at IP layer due to IPsec application is 80% (1).

$$\text{Increase of packet size} = (6C_{16} - 3C_{16}) / 3C_{16} = 80\% . \tag{1}$$

## 4  Conclusions

We could confirm and notice the two modes of operation of IPsec and the security provided by its security protocols. The message used for the test wasn't big enough to check other functionalities of IPsec, e.g., the injection of dummy data in the SA that we expect to do in future tests.

The results obtained on the test environment are too few for providing general conclusions. However it was evident the IPsec encryption and the authentication it provides by applying AH. The significant overhead due to IPsec was also pointed and should be kept in mind of everyone that is thinking of implement it anywhere.

We believe that IPsec implementation is worthwhile but a study on its implication especially on a small bandwidth channel should be previously done.

As a final conclusion, IPsec should be considered a major player in communications security domain.

## References

1. Computer Industry Almanac, http://www.c-i-a.com/compuseexec.htm (accessed on the 2nd January 2010)
2. Destaques, Instituto Nacional de Estatística, Portugal, http://www.ine.pt/xportal/xmain? xpid=INE&xpgid=ine_destaques&DESTAQUESdest_boui=56910276&DESTAQUESmod o=2 (accessed on the 2nd January 2010)
3. nderson, R.: Security Engineering (second edition), Wiley, 2008
4. Pfeelger, C.P., Pfeelger, S.L.: Security in Computing, Fourth Edition, Prentice Hall, 2006
5. Bishop, M.: Introduction to Computer Security, Prentice Hall PTR, 2004
6. Internet Society, RFC2401, http://www.rfc-editor.org/rfc/rfc2401.txt
7. Internet Society, RFC1825, http://www.rfc-editor.org/rfc/rfc1825.txt
8. Internet Society, RFC2402, http://www.rfc-editor.org/rfc/rfc2402.txt
9. Internet Society, RFC2403, http://www.rfc-editor.org/rfc/rfc2403.txt
10. Internet Society, RFC2404, http://www.rfc-editor.org/rfc/rfc2404.txt
11. Internet Society, RFC2406, http://www.rfc-editor.org/rfc/rfc2406.txt
12. Internet Society, RFC4301, http://www.rfc-editor.org/rfc/rfc4301.txt
13. Internet Society, RFC4302, http://www.rfc-editor.org/rfc/rfc4302.txt
14. Gregg, M., Watkins, S., Mays, G., Bandes, R., Franklin, B.: Hack the Stack – Using Snort and Ethereal to Master the 8 Layers of an Insecure Network, Syngress, 2006
15. Internet Society, RFC4306, http://www.rfc-editor.org/rfc/rfc4306.txt
16. Arkko, J., Nikander, P.: Limitations of IPsec Policy Mechanisms. In: Security Protocols, LNCS, vol. 3364, pp. 252--254, Springer, Heidelberg (2005)
71. Internet Society, RFC4303, http://www.rfc-editor.org/rfc/rfc4303.txt
18. Internet Society, RFC0792, http://www.rfc-editor.org/rfc/rfc792.txt
19. RFC0791, http://www.rfc-editor.org/rfc/rfc791.txt

# RPM: Anonymous Communication in P2P Networks

Robson Costa

Engineering Faculty – University of Porto (FEUP)
Rua Dr. Roberto Frias, 4200-465 – Porto – Portugal
robson@fe.up.pt

**Abstract.** Anonymity is a growing concern in recent Internet-based systems. Traditional mix and multicast based anonymity networks have a number of reliability, confidentiality, and performance issues. The large scale of P2P networks can be leveraged to minimize such issues, but these networks have to deal with node churn (a major factor in P2P systems) and the lower trustworthiness of individual nodes. In this paper we introduce RPM, a protocol for anonymous communication in P2P systems. In addition to anonymity, RPM aims at being resistant to churn and lowering the overhead usually found in anonymity systems. Our results show that RPM is effective in satisfying all these requirements, especially with respect to churn resistance.

**Keywords:** Anonymity, Routing, P2P networks

## 1 Introduction

Many applications in the Internet need to keep the confidentiality of its communications. Many times, these necessities go beyond secret of content of messages, being also needed the secret regarding who is the entities that communicate (*i.e.* final users whom they desire to keep to secret on the people with which they share to instantaneous messages and emails). In situations as these, end-to-end cryptographic mechanisms — they guarantee confidentiality of content — are not capable, by itself, of supply the desired protection; for this, **anonymity systems** are used.

Diverse architectures for anonymity systems exist. Simplest it uses only one node as intermediate between the origin and destination [1], with cryptographic chanels between nodes and intermediary. The degree of secret offered for this architecture is limited: if the intermediate node be compromised or apprehended, the secret of all the communications is violated. Moreover, if this node suffer an imperfection, the anonymous communication if it becomes impracticable. For to skirt these limitations, other architectures had been proposals based in networks (logical) of anonymity. In these networks, the messages follow a sequence of nodes before arriving at its destination. The idea is that each node that processes the message has information limited on its true origin or destiny, as form to prevent that one node

discovery the comunicante pairs in the network or to tie to a certain traffic specific nodes. Moreover, the dispersion of the traffic in several nodes make the task of external observers it difficult for try to identify standards of communication in the network.

Two basic variants of anonymity networks exist, the networks of mixers and the networks based on multicast. In the networks of mixers [2], the nodes of origin mounts a sequence of intermediate nodes until the destination and uses encryption layers in such a way (onion encryption) that each intermediate node only knows its predecessor and its successor in the route, and not the real origin or destiny of the traffic. In the networks based on multicast, a message is sent for a group that contains its real address; a external observer is incapable to infer which member of the group is effectively the destined. Although they is a significant advance in relation to an only intermediary, the majority of the anonymity networks also suffers from limitations in terms of reliability, confidentiality and performance, how the functions that guarantee the anonymity are generally centered in a small set of nodes. How much lesser is the network, more serious if they become these limitations.

An attempt to solve the problems with classic anonymity networks, the peer-to-peer (P2P) networks are used. The idea is to share the resources of the users who desire communications anonymous, using to advantage the scale ample and available computational capacity in P2P networks to mainly remove the restrictions of scalability of the classic networks. P2P networks are typically characterized for its great number of participants, voluntary machines, having a great amounts of nodes that they can be used as intermediate between origin and destination. Much even so systems of anonymity P2P are efficient in this direction, a challenge that they face is the constant join and leave of nodes of the network, a known phenomenon as churn.

In nets of mixers, churn makes with that the routes need be reconstructed frequently. In the networks based on multicast, aggravates the problem of management of the used cryptographic keys for communication with groups, that need to be changed to each change in the composition of the group. In both the cases, are necessary operations of high computational cost generating a delay, what have negative impact in the performance of the nets. Another challenge involving nets of anonymity P2P is that, as the route in these nets is made in the application layer, it becomes more easy to observe the traffic, what facilitates the task to spy on the other people's communications.

This paper considers RPM (Random Path + Multicast), a protocol for anonymous communication in P2P networks. The RPM has for objective to guarantee the anonymity of nodes from external observers, minimizing overhead and offering resistance to churn. The gotten results evidence that RPM is capable to supply anonymous communication with high level of same reliability with high rates of churn.

The paper is organized of the following form: the section 2 presents the objectives and premises used, the section 3 describes the protocol in details, section 4 present a evaluation of protocol is made, in a qualitative way and in a quantitative, and in the section 5 the related work are argued, finally, section 6 presents conclusions and future perspectives.

## 2 Objectives and Premises

### 2.1 Objectives

The main objective of the considered project is to make possible the anonymous communication between nodes in a P2P network. More specifically, it is desired to guarantee the following attributes in the terminology of [3]: 1) the anonymity of relationship between sender and receiver from others nodes in the network; 2) the anonymity of the sender from others nodes in the network and the receiver; and 3) the anonymity of the receiver from others nodes in the network.

Beyond these objectives in anonymity terms, other specific objectives are the resistance to churn, that it allows that nodes if communicate exactly that nodes exist diverse others entering and leaving P2P network, and low overhead of the solution, characterized for the restricted use of cryptographic mechanisms, especially those basing on public keys.

### 2.2 Premises

The project argued in the section 3 estimate the existence of a P2P network that allows the discovery of IP addresses of nodes that his constitute, and that it have some mechanism so that a node locates who is the receiver with who wants to communicate itself. Moreover, it is considered that a node can discover, anonymously, the public key of this receiver (*i.e.*, recouping a set $\mathscr{C}$ of certifyd of a repository such that the certificate of the receiver $C_R \in \mathscr{C}$).

## 3 Protocol Description

### 3.1 General View

The RPM (Random Path + Multicast) uses a random route to send messages from sender until the receiver. With each intermediate node choosing the next jump randomly amongst its neighbors in P2P network, in the origin, each message receives an accountant from jumps (hop count), that it determines how many intermediate nodes are used before being transmitted for the destination; this accountant is chosen randomly for each message, and decremented to each jump. To guarantee the

anonymity of the receiver, the sender choose a group of nodes (that it includes the receiver) for which the message will be transmitted when the hop count arrives to zero (figure 1(a)).



(a) Sending the message of the sender $E$ to receiver $R$ ($G$ is the reception group)

(b) Response of the receiver $R$ to the sender $E$ ($G'$ is the response group)

**Fig. 1.** General view of the protocol

The sending of responses of the receiver for the sender follows same process (figure 1(b)), being that the response group also is chosen and inserted previously by the sender in original message.

The data messages are ciphered using an established symmetrical key dynamically at the beginning of the communication between origin e destination. This key is modified periodically to minimize the risks. Control messages are used in the establishment and change of cryptographic symmetrical keys.

### 3.2   Sending of the messages

The algorithm 1 show the function RANDOM-SEND, used to transmit messages for one receiver; this function it is not invoked directly by an application, being used as block of construction for the functions that transmit data or control messages. A message have the format $\langle type, hc, G, id, ciphertext \rangle$, where: *type* is the type of message, that can be of data or control; *hc* is the hop counter chosen randomly for each message that have a variation from $hc_{min}$ to $hc_{max}$; $G$ is the reception group;[1] *id* is a unique identification used by receiver for locate the context of the message, and that also is defined in the beginning phase; *ciphertext* is the content to be transmitted, already ciphered with the appropriate key (that it depends on the type of the message).

---

[1] The reception group $G$, the response group $G'$ and the identification id is defined in the beginning phase, and argued in details in section 3.4.

---

**Algorithm 1** Random sending of messages

---

1: **procedure** RANDOM-SEND(*type*, *G*, *id*, *ciphertext*)
2:     $hc \leftarrow$ RANDOM($hc_{min}$, $hc_{max}$)         // it chooses by lot a number $hc \in [hc_{min}, hc_{max}]$
3:     $M \leftarrow \langle type, hc, G, id, ciphertext \rangle$
4:     send $M$ to random neighbor
5: **end procedure**

---

The algorithm 2 describes as the transmitted messages using RANDOM-SEND are processed. When receives a message $M$, the node verifies the value of $hc$. If $hc = 1$, the node assumes the role of **exit node**, transmitting the message for nodes who are part of the reception group $G$ (line 5). Case $hc > 1$, the node decrease $hc$ and send the message for one random neighbor who not it its predecessor (line 7). If $hc = 0$, the node belongs to the reception group $G$, and can be the receiver of the message. If the node be a possible receiver and the message be a data message, it tries to locate the symmetrical key corresponding to $id$ of the message (line 9 to 11). If obtain, the node deciphers *ciphertext* using this key and then processes its content (lines 12 to 14); of the opposite, the node concludes that it is not the receiver and discards the message.

---

**Algorithm 2** Processing of transmitted messages randomly

---

1: **upon receiving** $M = \langle type, hc, G, id, ciphertext \rangle$ **from** $x$ **do**
2:     **if** $M.hc > 0$ **then**
3:         $M.hc \leftarrow M.hc - 1$
4:         **if** $M.hc = 0$ **then**
5:             send $M$ to $M.G$
6:         **else**
7:             send $M$ to random neighbor $\neq x$
8:     **else**
9:         **if** $M.type =$ DATA **then**
10:             $K_S \leftarrow$ GETKEYFROMID($M.id$)
11:             **if** $K_S \neq \bot$ **then**
12:                 $payload \leftarrow$ DECRYPT($K_S$, $M.ciphertext$)
13:                 SETRESPONSEGROUP($M.id$, $payload.G'$)
14:                 deliver($payload.data$)
15:         **else**
16:             // process control message
17: **end do**

---

Each node keeps one cache of the symmetrical keys defined by the senders who that communicate with it (as receiving). *id* is used in the messages of data to lo-

cate the key that must be used for decipher; in case that the node does not have a corresponding key to *id* supplied, it is not the receiver of the message. With this, prevents that each node of the reception group executes a hard operation (as a decipher using its private key) to only verify if it is really the receiver, what improves the performance of the net.

## 3.3 Response messages

The sending of response messages follows the same principles of the original messages, but some restrictions must be considered. The main it is that the receiver does not know the sender of the original message, only the response group that must use to communicate itself with it. The used procedure to send response messages is shown in algorithm 3.

   This procedure receives as parameters the original message $M^o$ and the content from response to transmit (*response*). The response message uses the reception and response groups of the original message with inverted papers (lines 2 and 3). The symmetrical key is recouped from the same identification *id* of the original message (line 4), being used to cipher the content and the response group together (line 5). The resultant message then is transmitted using RANDOM-SEND (line 6).

---

**Algorithm 3** Transmission of response messages

---
1: **procedure** SEND-RESPONSE($M^o$, *response*)
2:     $G \leftarrow M^o.G'$
3:     $G' \leftarrow M^o.G$
4:     $K_S \leftarrow$ GETKEYFROMID($M^o.id$)
5:     *ciphertext* $\leftarrow$ ENCRYPT($K_S$, $\{response, G'\}$)
6:     RANDOM-SEND(DATA, $G$, $M^o.id$, *ciphertext*)
7: **end procedure**

---

## 3.4 Beginning of data flow

Before it is possible to transmit given of the sender for the receiver, it is necessary make the beginning of the data flow. In this phase, sender defines, beyond the used symmetrical cryptographic the key for data cypher, a reception group and a response group. The first group is a set of nodes pertaining to P2P net that includes the receiver, the second group is a set of nodes pertaining to P2P net that includes sender. Both the groups are formed choosing randomly a set of nodes of the net and including the receiver or the sender, as the case.

To discover nodes who can be part of the groups, can be used a mechanism of filiation (membership) of proper P2P network if this will be available. However, such mechanism is not indispensable, therefore each node can go forming its proper idea of the filiation from its neighbors and of nodes who appear in reception and response groups in the messages that it receives. Moreover, the neighbors also can change its lists of members periodically.

The algorithm 4 shows the behavior of the sender in the beginning phase. First, it randomly generates the reception and response groups, the symmetrical key $K_S$ and an unique identification $id_S$ (lines 2 to 5). To follow, the sender sends a control message for the receiver with $K_S$, $id_S$ and the response group $G'$ that will be used for the return of messages; these data they are ciphered with the public key $K_R^+$ of the receiver (line 6). This message is transmitted using RANDOM-SEND (line 7). The sender wait a confirmation of the receiver before starting to transmit given, thus preventing the wastefulness of computational resources and transmission band case the receiver cannot be reached. Multiple sendings or retransmissions can be used to become this communication trustworthy.

---

**Algorithm 4** Beginning of data flow: sender code

---

 1: **procedure** INIT-FLOW($R$)
 2:     $G \leftarrow$ GENERATERECEPTIONGROUP($R$)
 3:     $G' \leftarrow$ GENERATERESPONSEGROUP($R$)
 4:     $K_S \leftarrow$ GENERATEKEY()
 5:     $id_S \leftarrow$ GENERATEUNIQUEID()
 6:     $ciphertext \leftarrow$ ENCRYPT($K_R^+, \{id_S, K_S, G'\}$)
 7:     RANDOM-SEND(NEW-KEY, $G$, $\perp$, $ciphertext$)
 8:     set $id_S$ as pending
 9: **end procedure**

10: **upon receiving** $A = \langle$KEY-ACK, $hc$, $G'$, $id_S$, $\{id_R\}_{K_S}\rangle$ **from** $x$ **do**
11:     **if** $A.id_S$ is pending **then**
12:         SETRECEIVERID($A.id_S$, $A.id_R$)
13: **end do**

---

The procedure of the receiver is detailed in algorithm 5. When receiving a control message $\langle$NEW-KEY, $hc$, $G$, $\{K_S, id_S, G'\}_{K_R^+}\rangle$, the receiver uses its private key $K_R^-$ (line 2) to extract the symmetrical key $K_S$ and the response group $G'$ (in case that fails on decipher, the node is a member of $G$ that not real it receiver $R$). In the sequence, the receiver generates its identification $id_R$ (line 4), associates the key $K_S$ with this identification (line 5) and sends a message of confirmation with type KEY-ACK for the group $G'$ (lines 6 and 7). When the sender receives the confirmation, it

associates the identification of the receiver $id_R$ to data flow that was hanging for the identification $id_S$ (algorithm 4, lines 10 to 13).

---

**Algorithm 5** Beginning of data flow: receiver code

---

1: **upon receiving** $M = \langle$NEW-KEY$, hc, G, ciphertext \rangle$ **from** $x$ **do**
2:     $payload \leftarrow$ DECRYPT$(K_R^-, M.ciphertext)$        // $M.ciphertext = \{K_S, id_S, G'\}_{K_R^+}$
3:     **if** decryption is sucessful **then**
4:         $id_R \leftarrow$ GENERATEUNIQUEID$()$
5:         ASSOCIATEKEY$(id_R, payload.K_S)$
6:         $ciphertext \leftarrow$ ENCRYPT$(payload.K_S, id_R)$
7:         RANDOM-SEND$($KEY-ACK$, payload.G', payload.id_S, ciphertext)$
8: **end do**

---

## 4  Evaluation

The evaluation of this work was divided in two parts, a qualitative analysis of offered security (section 4.1) and an simulation analysis results (section 4.2).

### 4.1  Security

It is possible, through a qualitative analysis, evaluate how much the considered model satisfies its objectives, that are anonymity of relationship, sender and receiver from third, and anonymity of sender from the receiver (section 2.1). As the anonymity of relationship exists whenever it has anonymity of sender and of receiver [3], is argued only these two last, that automatically guarantee the first.

**Sender Anonymity**  The sender anonymity face the receiver is supplied by the use of response groups, and is intrinsically linked on the size of these groups. How much bigger is it, more difficult will be to identify who is the emitting; the problem to use very great response groups is that this harms the performance of the network, which had to the multicast carried through the exit node. Another aspect of this type of anonymity is that the beginning of data flow does not distinguish between an sender who goes to initiate one transmission of an sender who only wants to change the parameters of the flow; the only form of the receiver to differentiate the two cases is if it have few established flows and a new response group sent in the message NEW-KEY will have a intersection with the response group of some of the established flows (still thus, this inference has a imprecision degree). The anonymity of the sender before observers who belong to P2P network is given by

variation of hop counts *hc*; how a node never knows to the certain o initial value of *hc*, it cannot determine with clarity who is the sender, since none does not exist another information that of identifies in the message.

**Receiver anonymity** As well as it happens in the anonymity of the sender before the receiver, the anonymity of the receiver before others nodes in P2P network is given by the reception group, and by the size of this group. An adversary who has controlled a node that receives a message can act as malicious exit node, zeroing *hc* and sending the message separately for each node of the reception group, and observing if this message excites some response. Exist two difficulties in implement an attack of this type. The first is the randomly route: this attack will demands that the malicious node receives a message, and this is not guaranteed. To another difficulty is the necessity of adversary be capable to observe all the traffic of the target node, and exactly thus it can arrive to a wrong conclusions (if the volume of traffic in the network is high). Another attack that can be attemped, especially when it has little traffic in the net, is the intersection of reception groups. This attack is facilitated by the use of fixed *id* in each message, that can be used by observers to identify messages that belong the same flow. The solution adopted for this problem is keep the reception group exactly constant for all the messages with *id*, in way that if becomes impossible to make this intersection of the groups (in contrary case, the node that appeared in all the groups with *id* would be probably the receiver).

## 4.2   Results

**Simulation Description** Quantitatively to evaluate the degree of anonymity and the resistance to churn of the considered mechanism, had been carried through some experiments of simulation. For this was used the PeerSim [4], a simulator of P2P networks based on Java. The topology of P2P networks follows the ScaleFreeBA model with connection degree 3, that is, a topology scale-free [5] where each node have three neighbors at least. As the simulator is deterministic, has been created 100 distinct seeds in order to get the results of 100 different networks. For the calculation of average values had been discarded ten of the better samples and ten of the worse samples to prevent the influence of outliers. Three different sizes of networks ($2,000$ nodes, $6,000$ nodes and $10,000$ nodes) and three different sizes of reception groups had been simulated (4, 6 and 8 nodes). Moreover, churn was varied of 0% the 90%, at intervals of 10%. This rate of churn represents the percentage of nodes that are substituted in the topology during simulation.

**Simulation Results** First, one measured how much of the referring traffic to a transmission it was observed in each intermediate node. The objective of this is to evaluate resistance of the net the traffic analysis, a time that this type of analysis

depends intrinsically on the amount of traffic that can be observed. The graphs in figure 2(a) shows the average percentage of traffic observed in each node, gotten for the reason the traffic received in node and the total transmitted by the sender (considering only nodes without be the sender or the receiver). The results demonstrate that isolated nodes have access to a small fraction of the traffic, less than 0,6% in all the analyzed situations and below of 0,2% for networks with 2,000 nodes. This make of the analysis of traffic one attack extremely difficult for isolated nodes, and proves that the random route is efficient in spread the traffic of the network enter different nodes.

The trends observed in figure 2(a) are waited. The fraction of observed traffic is inversely proportional to the size of the network, therefore with the increase of the network, the number of available nodes grows for be used as intermediate, and each node starts to receive a ratio lesser from traffic. On the other hand, the observed traffic is directly proportional to the size of the reception group, therefore how much bigger the group more copies of the messages is sent (for the exit node) and observed.



(a) Average traffic    (b) Maximum traffic

**Fig. 2.** Observed traffic by node

The graph of figure 2(a) shows the observed average of traffic in each node. To evaluate how much this average is representative, and to verify which the situation of worse case, also analyzed the maximum percentage of traffic observed for a any node in all the simulated situations. The gotten results are shown in figure 2(b) (the scale in the axle $y$ is different of that one of figure 2(a)). The worst case of all occurs when we have 8 nodes in reception group and a network with 2.000 nodes; however, exactly in this in case that, the observed percentage of traffic was of 7,4%, this index that still becomes the analysis of traffic a difficult task.

In according to place, the average reliability of the net was evaluated front the different levels of churn, with intention to evaluate the resistance of the net to this phenomenon. The results are shown in figures 3(a), 3(b) and 3(c). The reliability is given by the reason between the number of delivery messages in the receiver and the number of messages sent by the sender. As waited, the reliability is of 100% for nets without churn and goes falling to the measure where this increases. The fall, however, is gradual, and same with 90% of churn, the reliability remains high, with approximately 75% of messages been delivery to the receiver in all the simulated scenes. These results shows that the random route used in this paper is sufficiently resistant to churn.

Analyzing the curves in each one of graphics 3(a), 3(b) and 3(c), and comparing the graphs between itself, perceives that it does not have a strong correlation between the reliability and size of the net or the reception groups. It seems to exist a small trend of that bigger nets are more reliable, but an sanalysis of the collected data discloses that this relation is not statistical significant.



(a) $|G| = 4$

(b) $|G| = 6$

(c) $|G| = 8$

**Fig. 3.** Reliability of the transmissions

# 5 Related Work

In the mixnets approaches [2], an sender determines the sequence of nodes to be followed by a message, in such a way that each node knows only its predecessor and its successor in the route. The anonymity is guaranteed by onion encryption: the emitting one message successive ciphers layers, and each intermediate node removes a layer to discover the next jump. Exist diverse implementations of this basic model [6–8], each one with small variations in relation to the original. Compared with the proposed approach, these experiences present as main disadvantages the susceptibility to churn, that it can make with that many routes have that to be reconstructed (an operation with high cost), and (in some cases) the performance, since high cost cryptographic operations are used extensively. The advantage of the nets of mixers is that they offer an stronger anonymity of receiver, since intermediate nodes not have notion of which nodes can be receiver; on the other hand, when the receiver be external of anonymity network [8], exists a exit node that knows the identity of this receiver. Some proposals [9, 10] works with the problem churn in networks of mixers making with that each cypher layer is addressed to a group of nodes, and not more to a node only. This approach introduces the problem to manage the filiation (membership) of the groups and keys to cipher messages for the groups (public or symmetrical).

The works more nearly of the considered approach are those that do not depend on the previous construction of routes through the networks [11, 12]. In the MuON [11], the sender generates a header each message to be transmitted. This header is spread out using an epidemic protocol (gossip); the node that can decipher an item of header (using its private key) is the receiver of the message, and requests the complete message to its owner, also identified in header. Nodes who are not the receiver also can requesting the complete message, what he guarantees the anonymity of the receiver. When receiving a complete message, the node passes if to announce as owner of this message in header that it spreads out in the net, what it guarantees the anonymity of the sender. The final result is that header are spread out in all the net, and the message (with payload) is only sent for some nodes. The disadvantage of the MuON in relation to the RPM is that this diffusion of headers makes with that all nodes need execute a decipher using its private key for each message received, independent package to be the receiver or not; in the RPM, beyond the messages to be spread out only for the reception group, the use of *id* prevents that others nodes process messages uselessly.

The Rumor Riding (RR) [12] is a protocol for anonymous operation in nets P2P not structuralized. In the transmission, each message is ciphered with a symmetrical key; the ciphered message and the key are sent separated for different neighbors, and uses random routes in the net. When these packages if find in one node any, this deciphers the message using the supplied key and transmits the message deciphered (will be a consultation, this is spread out using flooding; in the case of being a

message directed to a specific receiver, the node opens a connection TCP with one proxy that it answers for the receiver). So that the probability of that the pairs of packages (key and ciphered message) if find either significant, each node of net P2P keeps one cache contends the messages received recently; each package of key that arrives is compared with the content of this cache to see if this key if applies to some of the stored messages. Moreover, the experimental results disclose that, in practical, the RR requires the use of replication of traffic e of a TTL (time-to-live, that represents the maximum number of jumps of the packages) very high; a good probability of meeting is gotten with $TTL > 30$ and of 2 the 6 rejoinders of each package. The RPM, on the other hand, layoff the use of replication of traffic, uses routes considerably shorter, it does not have overhead associated to the maintenance of caches of packages and to the processing of the messages received against these caches, and reduces (through *id*) the necessity of high cost cryptographic operations in nodes that not is the destiny; all these factors represent a significant profit of performance for the RPM. The main disadvantage of the RPM in relation to RR is the lesser anonymity of the receiver.

## 6 Conclusion

This article presented RPM, a approach for anonymous communication in P2P networks. The RPM uses random route to confer resistance to churn characteristic of nets P2P and a set of mechanisms to reduce the high cost generally associated to the used protocols to guarantee the anonymity in the communications. The results demonstrate that the RPM reaches the objectives fully that if it considers, particularly in regards to churn.

As extension of this work, it is considered to investigate a method to improve the anonymity of the receiver, turning the groups opaque from intermediate nodes, and to extend the studies through experimentation.

## References

1. Anonymizer: The Anonymizer. `http://www.anonymizer.com/` (2007)
2. Chaum, D.L.: Untraceable electronic mail, return addresses and digital pseudonyms. Volume 24. (1981) 84–88
3. Pfitzmann, A., Hansen, M.: Anonymity, unlinkability, undetectability, unobservability, pseudonymity, and identity management – a consolidated proposal for terminology. Version 0.30 (26 de novembro de 2007) `http://dud.inf.tu-dresden.de/Anon_Terminology.shtml`.
4. PeerSim: Peersim: A peer-to-peer simulator. `http://peersim.sf.net/` (2007)
5. Albert, R., Barabási, A.L.: Statistical mechanics of complex networks. Volume 74., American Physical Society (Jan 2002) 47–97

6. Freedman, M.J., Morris, R.: Tarzan: a peer-to-peer anonymizing network layer. In: Proceedings of the 9th ACM Conference on Computer and Communications Security (CCS-9), Washington, DC, USA (November 2002) 193–206
7. Rennhard, M., Plattner, B.: Introducing MorphMix: Peer-to-peer based anonymous Internet usage with collusion detection. In: Proceedings of the 2002 ACM Workshop on Privacy in the Electronic Society, ACM Press New York, NY, USA (2002) 91–102
8. Dingledine, R., Mathewson, N., Syverson, P.F.: Tor: The second-generation onion router. In: Proceedings of the 13th USENIX Security Symposium, San Diego, CA, USA (August 2004) 303–320
9. Zhu, Y., Hu, Y.: TAP: A novel tunneling approach for anonymity in structured P2P systems. In: Proceedings of the International Conference on Parallel Processing (ICPP). (2004) 21–28
10. Zhuang, L., Zhou, F., Zhao, B.Y., Rowstron, A.: Cashmere: Resilient anonymous routing. In: Proceedings of the 2nd Symposium on Networked Systems Design and Implementation (NSDI). (2005)
11. Bansod, N., Malgi, A., Choi, B.K., Mayo, J.: MuON: Epidemic based mutual anonymity. In: Proceedings of the 13th International Conference on Network Protocols (ICNP). (2005) 99–109
12. Han, J., Liu, Y.: Rumor riding: Anonymizing unstructured peer-to-peer systems. In: Proceedings of the 14th IEEE International Conference on Network Protocols (ICNP). (2006) 22–31

# Speech Quality Assessment when Using VoIP Over ADSL

Altino Sampaio

Faculdade de Engenharia da Universidade do Porto, Rua Dr. Roberto Frias, s/n
4200-465 Porto, Portugal
pro09002@fe.up.pt

**Abstract.** Speech quality in VoIP calls is a major home users concern when sharing ADSL subscriber lines of limited bandwidth with other type of traffic. Using VoIP service providers that home users can easily access with open source telephony systems, we measure speech quality, by transmitting a sample speech signal to PSTN phones, through ADSL channel, while varying the audio codec used and injecting different type of traffic. Measurement's data are gathered by conducting MOS listening tests with listeners assessing speech quality. Taking advantage of the specific characteristics of different audio codecs, this paper investigates and proposes how speech quality can be improved, even when different type of traffic competes for channel bandwidth. The results show that some audio codecs respond better to the presence of different types of traffic, in an ADSL channel of limited bandwidth.

**Keywords:** Speech quality, Voice over IP (VoIP), ACR, MOS, Asterisk.

## 1 Introduction

Telephone systems are seen by humans as a way of interaction, through which it is possible to transmit voice data. The international telephone system, based on copper wires carrying voice data, is referred as the Public Switched Telephone Network (PSTN), that reliably facilitate telephone conversations at any time of day, year-round [1]. Today's PSTN offers many different services and one of the most important is Plain Old Telephone Service (POTS), which refers to the standard voice telephone service that most residential and small business subscribers use. When we are using a PSTN line, we typically pay charges according to the time usages.

Plain old telephone service will not be used for much longer [2]. With ADSL subscriber lines, users tend to relying on the power and flexibility of data networking and the Internet [3]. Overtime, voice services using Internet Protocol (IP) [4] will simply be another application in everyone's homes.

Voice over Internet Protocol (VoIP) is a revolutionary way of communicating, which stands for a transmission technology for delivery of voice communications over IP networks, such as data networking and the Internet. VoIP converts the voice signal from telephone into a digital signal, sends it through the Internet and then converts it back at the other end. Some of the most important advantages deriving

from the use of such technology are cost saving, talk with many people at the same time and, while we are talking, we can exchange data with people we are talking with. For this reason, VoIP is seen as the method of choice for carrying telephone calls. However, as IP networks are inherently nondeterministic and can suffer from packet loss and packet delay variation, it is almost impossible to guarantee the voice quality delivered to customers who are accustomed to reliable and stable circuit-switched performance [5]. Also, ADSL subscriber lines have limited bandwidth and different types of traffic are transmitted over it. As a consequence, the number of simultaneous calls and the speech quality are also limited, due to the channel bandwidth limitation. Audio codecs can be used to reduce bandwidth requirements to transmit VoIP data, thus augmenting the number of concurrent calls.

In this paper we will focus on the measurement of speech quality for a certain number of simultaneous VoIP calls to PSTN phone network, regarding the audio codec used to encode voice data, the presence of different type of traffic traversing the channel and bandwidth limitations of the ADSL subscriber line.

Rest of the paper is organized as follows, in section 2, we will introduce some related work about speech quality measurements. In section 3 of this paper, we describe typical scenario of the experimental setup. In section 4 we present results and we discuss them. Section 5 is conclusion and proposal for future work.

## 2 Related Works

In this section, we will introduce some related work about speech quality measurement and VoIP telephony systems.

### 2.1 Speech Quality Measurement

VoIP quality is difficult to expect because IP networks are inherently nondeterministic and can suffer from packet loss, packet delay, packet delay variation, packet error, etc. Although there is no physical definition for speech quality, this is vital to service acceptance by the user. To measure speech quality, there are two main approaches, which are objective testing and subjective testing [5].

In objective testing, the system physical properties are measured and results can be used to predict perceived performance. For telephony systems objective testing, there are two kinds of assessment methods, carried out either by injecting a test signal (intrusive or active measurement) or by monitoring live traffic (nonintrusive or passive measurement). Intrusive testing methods inject a reference acoustic speech signal into a system so they can be captured at a further point and distortion assessed (Fig. 1). Such method, standardized by ITU-T as P.862, is called Perceptual Evaluation of Speech Quality (PESQ) [6].

**Fig. 1:** Illustration of an intrusive test.

Nonintrusive or passive methods (Fig. 2) monitor live traffic directly to estimate the quality perceived by the user. Despite the nonintrusive monitoring measures are less accurate, they have a much lower operational cost, compared with intrusive method.



**Fig. 2:** Illustration of a nonintrusive test.

Nonintrusive method can further be classified as parametric model or signal-based model. In parametric model, results are constructed based on various important properties of telecommunications networks, such as packet loss, jitter and end-to-end delay of a call. This transmission quality model is defined in ITU-T G.107 and is commonly referred to as the E-model [6]. In the case of VoIP, such model is considered suitable for real-time evaluation of the call quality. Measurement results are repeatable, fast and efficient [7].

In subjective testing, quality of speech signals is ultimately adjudged by a group of human listeners, who assign an opinion score on an integral scale ranging between 1 (unacceptable) to 5 (excellent) [7]. Based on this assessment approach, there are two basic types of subjective test, namely conversational quality and listening quality tests [5]. As the name suggests, conversational quality tests involve two users conversing through a controlled communication channel. Such tests typically permit to find out the influence of conversational scenarios in echo and delay parameters, often associated with coding and channel error distortion. On the other hand, listening-only quality tests are based on unidirectional speech transmitted to users, permitting to evaluate the impact of coding and channel error distortions on quality.

There are several listening tests, namely Absolute Category Rating (ACR), Degradation Category Rating (DCR) and Comparison Category Rating (CCR) tests. In telecommunications area, Absolute Category Rating (ACR) tests are the most commonly used to assess integral quality. The usage of different sequences of two to five independent, short and meaningful sentences is recommended for ACR tests, with sentence durations of 2–3 seconds and overall sequence durations of about 12 seconds [6]. In this paper we will focus on ACR tests. In this kind of test, the user rates the quality of each speech sample on an absolute quality scale. Different scales

are recommended, such as listening quality, listening-effort or loudness preference [6]. The 5-point ACR quality scale is the most frequently used one. The average of these scores (votes of all users) obtained on this scale is called Mean Opinion Score (MOS), being calculated for a particular condition and provides a widely accepted measure for subjective speech quality.

There are various scales recommended, each focusing on different dimensions of quality [8]. The dimension of quality scale used in this paper is Quality of the Speech (Table 1).

**Table 1.** Relation among MOS level and user satisfaction.

| MOS Level | User Satisfaction Meaning |
|:---:|:---:|
| 5 | Very satisfied |
| 4 | Satisfied |
| 3 | Some users dissatisfied |
| 2 | Many users dissatisfied |
| 1 | Not recommended |

Traditionally, speech quality is estimated using subjective tests. In this paper we evaluate speech quality using subjective testing when transmitting some voice data (ACR listening test).

## 2.2 VoIP Telephony Systems

The idea of VoIP is to use the Internet as a telephone network with some additional capabilities. However, this technology can also be used to connect to the traditional telephone network. In such case, call travels across the IP network and terminates at the local phone network. Home users can use VoIP to call, nationally or even internationally, people who are served by phones, either fixed or mobile, for a fraction of the prices they would pay using PSTN.

VoIP Gateway is either software or hardware device that works as an entrance to another network, enabling communication between an IP network and legacy PSTN [9]. The main functions of VoIP gateways include voice compression/decompression, control signalling, call routing and packetization.

Asterisk is a complete Private Branch eXchange (PBX) in software, written in C programming language and it runs on Linux operating systems. Asterisk does VoIP in many protocols, PBX switching, codec translation and various other applications like voicemail, conference bridging, IVR and various others. The ability to load codec modules [10] allows Asterisk to set preferences for each Session Initiation Protocol (SIP) [11] channel as to which codecs should be used, or allowed. This can be necessary for Packet Voice over a low-bandwidth link while still providing high audio quality over less constricted connections.

A normal telephonic call requires 64 kbps, whereas the same can be handled an as low as 44.2 kbps using codecs such as GSM. Hence the bandwidth requirements will be considerably reduced and more calls will fit in the same given bandwidth. Table 2 shows the total Ethernet bandwidth consumed when using different audio codecs.

**Table 2.** Some codecs used by Asterisk [1], [12], [13].

| Codec | Voice Bits per Packet | Encoded Sound Bandwidth | Total Bandwidth | MOS |
|-------|----------------------|------------------------|----------------|-----|
| G.711 | 1280 bits | 64 Kbps | 95.2 Kbps | 4.3 |
| G.723.1 | 159 bits, 192 bits | 5.3, 6.3 kbps | 37.5 kbps | 3.8 |
| GSM | 160 bits | 13 Kbps | 44.2 kbps | 3.7 |

The MOS scores presented in the above table are merely indicative [13]. Asterisk can act as VoIP gateway or connecting to one. In this paper, we use Asterisk software to connect to a VoIP provider, which supports the audio codecs shown in table above.

## 3   Test Scenario Set-up

In this section, we describe the VoIP calls session. We first present the network set-up and later the average of votes of all listeners (MOS), for each network condition.

### 3.1   Network Set-up

In order to perform an assessment of the perceptual quality of VoIP traffic on ADSL access link, the system architecture depicted in Fig. 3 was implemented.



**Fig. 3:** Experimental set-up to measure speech quality.

Asterisk 1.6 was installed on Linux Ubuntu Server 8.10 operating system, running in a machine equipped with Intel Core 2 Duo 2.33 GHz processor and four gigabyte of memory. Asterisk was configured with an extension from which voice data was transmitted to thirteen different phone destinations. Asterisk extension was enabled in using GSM audio codec. A laptop computer was also introduced in order to generate the other type of traffic, hence consuming bandwidth on ADSL link. Laptop computer

had installed Windows 7, equipped with Intel Core 2 Duo 2.2 GHz processor and two gigabyte of memory. The other type of traffic consisted in downloads of video streamed in simultaneous with upload of large files. The amount of other type of traffic, together with VoIP traffic, consumed all the ADSL channel bandwidth.

The access to VoIP service provider, via ADSL line of 2048Kbps downlink and 512Kbps uplink, was established through modem router Netgear DG834. This modem router has no QoS support. In order to carry calls between Asterisk and PSTN phone network, we connected Asterisk to Voipraider [14] VoIP service provider, using the Session Initiation Protocol (SIP). As usual in VoIP applications, we assumed transmission over RTP [15] over UDP/IPv4.

### 3.2  Measurements Workflow

In this section we describe and explain the applied methodology in running the experiments. Our approach is to apply listening quality subjective test, using the most commonly method, ACR, and scale speech quality. Submitting a single speech sample, users adjudged the quality of the speech on a five-point scale: Excellent, Good, Fair, Poor and Bad, where Excellent equates to a score of 5 through to Bad which equates to a score of 1.

In our tests, six sets of experiments were performed. In each set of experiments thirteen VoIP calls were established with PSTN phones. Only a maximum of three concurrent calls were established at each moment. Sets were grouped in pairs, where, from pair to pair of sets, we varied the audio codec applied to transmit voice data from Asterisk to VoIP provider. In each pair of sets of experiments, the first set used ADSL line bandwidth entirely for VoIP data, unlike the second set where other type of traffic (OTT) was injected in the line. Table 3 resumes the sets of experiments done.

**Table 3.** Sets of experiments, according to different network conditions.

| Experiment | Codec | ADSL Traffic |
|------------|-------------|--------------|
| 1 | GSM | VoIP |
| 2 | | VoIP + OTT |
| 3 | G.711 (a-law) | VoIP |
| 4 | | VoIP + OTT |
| 5 | G.723.1 | VoIP |
| 6 | | VoIP + OTT |

The single speech sample transmitted in each VoIP call was about 12 seconds. Stored in a file, the speech sample was encoded using a standard GSM codec implementation. This packet stream was then sent to several PSTN phones with controlled network impairments. Trying different audio codecs and injecting different type of traffic in ADSL line, the output speech from every PSTN phone was listened by users which adjudged and scored the speech quality. Listening users were spatially separated from each other. The uplink channel constituted the link bottleneck, since it had lower bandwidth than downlink channel.

The idea was to measure the speech quality taking into account the ADSL line bandwidth limitation and audio codec used, in a typical residential access, where an environment of varied types of traffic apply. The speech quality should depend on audio codec used and network conditions.

### 3.3  Measurements Summary

The speech quality measurement took place on December 20th, 2009. We obtained 78 VoIP sessions, connecting to PSTN phone network, from ADSL subscriber line provided by a Portuguese ISP. The listeners considered for the test belongs to various age groups ranging from 18 – 35. The listeners were asked to rate the speech signal sample individually and scores were registered. The same speech signal sample was played for all the listeners. Each VoIP session was about 12 seconds of duration.

## 4  Speech Measurements and Analysis

In this section we present the results obtained in each set of experiments and MOS values for quality of speech, regarding each experiment.

### 4.1  Measurements Results

Based on listening ACR tests and accomplished all the sets of experiments, we gathered speech quality scores considered by listening users. The results collected are graphed in Fig. 4.



**Fig. 4:** Score values obtained from listening user's votes, regarding each experiment.

With the aim of estimate the perceived quality of speech, MOS was calculated for each particular network condition, by averaging the gathered votes of all users, previously illustrated in figure above. The MOS values for all experiments are presented in Fig. 5.



**Fig. 5:** MOS results for each experiment.

The coefficient of variation for each set of experiments, presented in Fig. 6, was calculated as the ratio of the standard deviation to the mean.



**Fig. 6:** Coefficient of variation obtained regarding each set of experiments.

The calculation of the coefficient of variation aimed to describe the dispersion of the gathered votes of all users around the MOS value, regarding the audio codec used and the type of traffic traversing the ADSL channel.

These codecs were chosen to ours experiments since they are largely implemented in telephone systems, as Asterisk. G.723.1 describes a compression technique that can be used to compress speech at a low bit rate. Despite it requires a license for the patents that cover the algorithm, it will be typically used and will be applied at a larger scale in future packet networks. The GSM is the main audio codec used by Asterisk system. The G.711 encoded voice is already in the correct format for digital voice delivery in the public phone network (PSTN) or through Private Branch eXchanges (PBX).

## 4.2 Results Analysis

The results show that GSM audio codec delivers the best speech quality, since it presents the higher MOS value obtained. When using such codec, the presence of other type of traffic in the same channel where VoIP traffic is transmitted causes almost imperceptible impact. However, considering the coefficient of variation for each pair of sets of experiments, we denote an increasing in the dispersion of speech quality scores around the MOS value, whenever other type of traffic is injected in ADSL channel. This behaviour could represent an impact in the stability of the transmission if other type of traffic competes for ADSL line bandwidth.

Unlike some authors [13] referred, in these experiments we verify that GSM audio codec presents the highest MOS score, followed by G.711 and then G.723.1. However, we must consider that the single speech sample transmitted to listening users was previously encoded with GSM audio codec. The extension through which voice data has been transmitted was configured in using GSM format. When connecting Asterisk to Voipraider VoIP provider using the G.711 and G.723.1 audio codecs, transcoding procedures are required. This explains why during the use of G.711 and G.723.1 audio codecs, listeners noticed deterioration in the quality of speech, even in absence of other type of traffic. This degradation is due to transcoding the speech sample signal to another format, suitable for transmission [16]. Taking as reference the obtained MOS value for GSM audio codec, a substantial decrease of 5.0% in MOS value occurs when speech sample signal is transcoded into the G.711 audio codec, regarding the absence of other type of traffic. Applying the same criteria, the decrease in MOS value is bigger, 15.0%, when transcoding into G.723.1. It is clear that the bitrate reduction operated by G.723.1 audio codec leads to signal quality degradation [16]. Considering these observations, we can expect deterioration in speech signal when calling outside parties whenever audio codecs interfacing the inner extension differs from the audio codec used to connect phone system to VoIP service provider.

Experimental results show that speech quality deterioration increases when using G.711 and G.723.1 audio codecs if other type of traffic competes for ADSL line bandwidth. Comparing the MOS values when using G.711 audio codec, we observe a stronger reduction of 14.0% if other type of traffic is injected in the ADSL transmission channel. This speech quality decrease is more severe than the 4.5% of

MOS reduction when using G.723.1 audio codec. Such observations show that signal degradation accents if it has been transcoded previously. We can conclude that, in stressed links, as ADSL subscriber lines, transcoding should be avoided whenever possible, since signal degradation increases. Moreover, results consistently demonstrate that G.711 audio codec is not suitable for transmission when ADSL link is being frequently used for other type of traffic than VoIP. Stressed transmission channels tend to lose packets, leading to increased latency and jitter. The G.711 audio codec requires more bandwidth resources, generates more voice bits per packets and longer packets suffer more loss [17]. Quality of speech degrades substantially if such codec is applied in such environment.

One cannot neglect the significant variation in MOS results for each network condition. Nevertheless, all experiments resulted in MOS above 3.0, where quality of speech signals is found to be acceptable. More important, our experiments demonstrated that we achieved optimal POTS quality since the utilisation of GSM audio codec in our VoIP telephony system resulted in a MOS of about 4.0 [18]. This conclusion applies even when different type of traffic traverses ADSL subscriber line of limited bandwidth. Thus, ADSL subscriber lines of limited bandwidth can reliably transport voice data, having speech quality identical to POTS telephony service.

## 5   Conclusions and Future Work

In this paper, we measured speech quality when using ADSL subscriber line of limited bandwidth.  To achieve this, we prepared sets of experiments, where distinct audio codecs were used and different type of traffic was injected into the channel used to transmit VoIP data.

Speech quality was assessed using ACR listening tests, a subjective testing technique, where a single voice signal was transmitted and adjudged by a group of thirteen human listeners, who assigned an opinion score on an integral scale ranging between 1 (unacceptable) to 5 (excellent). The average of these scores on this scale, called MOS, measured the quality of speech for each network condition.

Different type of traffic traversing the ADSL line used to transmitting VoIP data, associated with bandwidth limitation causing signal distortion, can cause serious negative impact on quality of speech. In order to reduce the speech signal deterioration, results consistently demonstrate that transcoding should be avoided whenever possible.

The obtained results lead us to conclude that the correct selection of audio codec can result in higher withstand to speech signal impairments. The experiments demonstrated that when different type of traffic is injected in VoIP transmission channel, the MOS value is severely reduced if G.711 audio codec is used. Nevertheless, MOS values obtained in all experiments are above 3.0, which stands for acceptable quality of speech signals.

The GSM audio codec, so common in Asterisk systems, presented stronger robustness to ADSL bandwidth limitations, achieving optimal POTS with a MOS of about 4.0. This paper showed that ADSL subscriber lines of limited bandwidth can

reliably transport voice data, having speech quality identical to POTS telephony service.

In the future work, we would like to use Perceptual Evaluation of Speech Quality, an intrusive testing method, to accurately measure the speech quality in ADSL subscriber lines.

## References

1. Wallingford, T.: Switching to VoIP. O'Reilly Media (2005)
2. Davidson, J., Peters, J., Bhatia, M., Kalidindi, S., Mukherjee, S.: Voice over IP Fundamentals (2nd Edition) (Fundamentals). Cisco Press (2006)
3. Noll, A.M.: Technical opinion: does data traffic exceed voice traffic? Commun. ACM 42 (1999) 121-124
4. Comer, D.E.: Internetworking with TCP/IP: Principles, Protocols, and Architecture, Vol. 1. Prentice Hall (2000)
5. Broom, S.R.: VoIP Quality Assessment: Taking Account of the Edge-Device. Audio, Speech, and Language Processing, IEEE Transactions on 14 (2006) 1977-1983
6. Raake, A.: Speech Quality of VoIP: Assessment and Prediction. John Wiley & Sons (2006)
7. Raja, A., Azad, R.M.A., Flanagan, C., Ryan, C.: An Evolutionary Approach to Speech Quality Estimation. Frontiers in the Convergence of Bioscience and Information Technologies, 2007. FBIT 2007 (2007) 757-760
8. Möller, S.: Assessment and prediction of speech quality in telecommunications. Springer (2000)
9. Alam, M.Z., Bose, S., Rahman, M.M., Abdullah Al-Mumin, M.: Small Office PBX Using Voice Over Internet Protocol (VOIP). Advanced Communication Technology, The 9th International Conference on, Vol. 3 (2007) 1618-1622
10. Qadeer, M.A., Imran, A.: Asterisk Voice Exchange: An Alternative to Conventional EPBX. Proceedings of the 2008 International Conference on Computer and Electrical Engineering. IEEE Computer Society (2008) 652-656
11. Chung-Hsin, L., Chun-Lin, L.: The Study of the SIP for the VoIP. INC, IMS and IDC, 2009. NCM '09. Fifth International Joint Conference on (2009) 1473-1478
12. Meggelen, J.V., Smith, J., Madsen, L.: Asterisk: The Future of Telephony. O'Reilly Media, Inc. (2005)
13. Rudkin, S., Grace, A., Whybray, M.W.: Real-time applications on the Internet. BT Technology Journal 15 (1997) 209-225
14. Voipraider, http://www.voipraider.com
15. H. Schulzrinne, S.C., R. Frederick, V. Jacobson: IETF AVT working group document RFC1889.
16. Duysburgh, B., Lambrecht, T., Turck, F.D., Dhoedt, B., Demeester, P.: An active networking based service for media transcoding in multicast sessions. Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on 34 (2004) 19-31
17. Black, U.N.: Voice Over IP. Prentice Hall (2001)
18. Duysburgh, B., Lambrecht, T., Dhoedt, B., Demeester, P.: On the quality and performance of active networking based media transcoding in multicast sessions. Telecommunications, 2003. ConTEL 2003. Proceedings of the 7th International Conference on, Vol. 2 (2003) 455-462 vol.452

# Data Mining

# Proposal of a Pattern-Mining Tool to Identify and Suggest Reconfigurable Units

Adriano Kaminski Sanches

Department of Informatics Engineering (DEI), Faculty of Engineering (FEUP),
University of Porto, R. Dr. Roberto Frias, 4200-465, Porto, Portugal
aksanches@gmail.com

**Abstract.** The advent of reconfigurable computing fabrics (e.g., FPGAs) makes now possible to implement entire and complex hardware/software systems in a single reconfigurable device that allow two distinct features: programmability and specialization. Parts of an application can be executed by specialized reconfigurable functional units (RFUs) tightly coupled to general purpose processors (GPPs). Specialized RFUs may be used to reduce execution time, power dissipation, energy consumption, etc. Typically, RFUs are designed based on empirical expertise without strong evidence on their efficiency over a broad range of real applications. Bearing this in mind, we are working on methods to analyze existing code repositories and recover valuable information from which part of the code would have a high potential for running in specialized RFUs and what shape those RFUs would need to have for more effective support. This paper shows our preliminary efforts on researching and developing software techniques able to address these issues.

**Keywords:** Reconfigurable Computing Fabrics; Pattern Mining; Information Retrieval

## 1    Introduction

The semiconductor companies face a challenge, each year harder, to keep the Moore's Law[1] true. The past technology miniaturization was making possible to increase the clock frequencies and to achieve the needed improvements in the hardware performance. Software developers were used to take advantage of this trend, developing each time more complex software that could execute in the newer platforms in a reasonable time, mainly because of the progress in the clock speed. Nowadays, the size of the electronic components is reaching critical limits and physical restrictions are to be faced in the upcoming years.

Currently, the main research focus to improve the performance of a computational system is the advent of multi-core systems. Nowadays it is common to see processors with 4 cores, such as the Intel Core 2 Quad [2] or 6 cores, such as the Six-Core AMD Opteron [3]. There are some companies with products available with hundreds of cores, like the PicoChip PC205 which integrates 273 processors in the same chip [4].
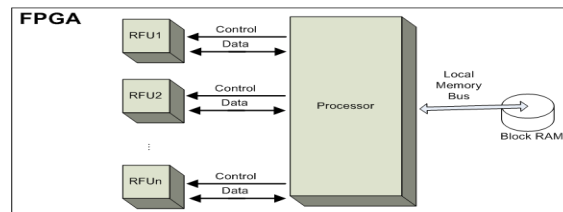
These systems are able to considerably improve the performance of multi-task operating systems that can distribute the tasks among the cores.

The main problem with multi-processor systems comes when one needs to allocate heavy tasks among several cores. To reach this purpose, the software developer must know in developing time how will the architecture benefit from the multi-core system and decide where each part of the code will be executed. However, it is not possible to parallelize all kinds of code because of data and control dependencies, and even when it is possible, sometimes it is not recommended because of the overhead assigned to the data distribution between the cores and the result aggregation. This means that, even adding more cores, in some specific applications, the overall performance is not improved.

One alternative to try to overpass this multi-core parallelization limitation is specialization. Specialization, also called customization, is a different approach that may complement the multi-core. It consists on trying to improve the performance, mapping some critical parts of the software application directly in best-fitted hardware implementations. Specialization has several benefits such as more power efficiency and less costly than general-purpose processors (GPPs) performing the same application [5]. The problem associated with specialization is the lack of programmability since the hardware implementation tends to be inflexible [6].

The size of the software mapped in hardware has an important factor in the flexibility of the overall system. There are some research efforts in trying to map a complete application in hardware [5, 7] in the context of Application-Specific Integrated Circuits (ASICs). This approach allows a high level of performance improvements, but the system generated is not able to execute a different application with the same hardware, even changing small features in the application. For instance, a video player that implements in hardware some specific video standard cannot take advantage of this logic and be compatible with a new standard that was not available when it was developed, being necessary to upgrade the hardware with the new specification to be able to play a new video standard.

However, the advent of reconfigurable hardware (e.g., FPGAs – Field Programmable Gate Arrays[8]) makes specialization an important solution as this kind of hardware brings the required flexibility and programmability. Reconfigurable hardware allows complex systems coupling GPPs to accelerator units (RFUs) that can be changed according to the application needs. This approach makes possible at the same time to improve the performance and to keep the flexibility of a GPP. In this kind of computing system, the user can detect the critical bottlenecks of the application and generate specific hardware to execute these special parts. The other parts can be executed in the usual way, i.e., using the GPP. Academic and commercial efforts on research and development of dynamically reconfigurable system-on-a-chip (SoC) architectures have shown diverse technologies and applications in various areas related to embedded systems [9]. The Fig. 1 illustrates one architecture based in this approach where is showed a GPP linked with one RAM memory and several RFUs coupled in the processor, each one performing a different task.

**Fig. 1.** RFUs coupled in one GPP

Energy savings and the satisfaction of performance requirements are usually very important features in mobile embedded systems (e.g., mobile phones, PDAs)[10]. These requirements are sometimes difficult to reach using traditional solutions, and reconfigurable hardware may become a prominent component in those systems. Because of the enormous design space to be explored, time-to-market pressure, and the ever increasingly complex applications, tools and methods that automatically exploit the synergies of these computational systems are necessary. Sometimes the non-existence of powerful tools makes impossible to evaluate and implement solutions which may efficiently benefit from the existing reconfigurable hardware resources.

Thus, methods for devising efficient RFUs are of paramount importance. These RFUs may have parameterized features (e.g., numbers of ALUs, granularity) needed to be explored. Similar software patterns can be mapped to a single RFU with some programmability to expose those different patterns by a time-multiplexing approach (e.g., improving hardware reuse). The minimization of hardware resources is usually an important factor and pattern-mining can be used to identify shareable resources [11].

We are researching methods aiming to achieve a tool to both identify the most appropriate code segments in software programs to run in RFUs with a certain structure, and to suggest RFU structures suitable for a large number of benchmarks. Our main goal is to optimize the overall performance of those benchmarks, also considering the sharing of hardware structures whenever possible. The methodology under development permits the user to select which features are important to be part of the patterns in the process of mining. Doing this, different levels of abstraction can be used (e.g., one might be interested in innermost loops with certain memory accesses shapes and/or in producer/consumer kernels). With our research we intend to propose software techniques able to give answers to the following questions:

- Is a given RFU efficient in terms of performance considering a set of applications or a benchmark repository?
- What should be the hardware structure for a performance efficient RFU considering a set of applications or for a benchmark repository?

The remainder of the paper is organized as follows. Section 2 gives the methodology used to identify the code segments; Section 3 discusses related work; Section 4 shows the current status and preliminary results and is followed by conclusion in Section 5.

## 2    Proposed Methodology in Pattern-Mining to Evaluate and Suggest RFUs

The main purpose of this work is both to evaluate a given RFU and to propose new RFUs capable of improving the performance for a large number of benchmarks using the minimal possible resources. To accomplish this purpose, we need to identify similar code sections in the set of applications expected to be executed in the target computing system. Such code sections will maximize hardware reuse by sharing the same RFU.

We believe that finding approximated similarities in high level source code could help us to diminish the huge exploration space of all the source code in the input applications/benchmarks. In this way, we could benefit from clone detection research [12] to locate these similarities and then, perform a post processing step to elect the possible candidates to generate RFUs using third party tools.

We intend to build a tool able to starting from the benchmarks analysis, perform RFU detection and RFU generation. To evaluate the techniques researched we will attach the RFUs to the GPP and make the changes in the original code to reflect the use of the RFUs (this requires communication and synchronization primitives). Bearing on mind this, we propose the steps showed in Fig. 2. Next sections describe each stage of this flow.



**Fig. 2.** Block diagram of the proposed tool

## 2.1    Preprocessing

This stage performs the first steps through the source code to be analyzed. Some uninteresting parts of each source code are deleted, e.g., the header of each file and the variable initialization. This is done to avoid the detection of false positives clones in the match detection block. After that it is important to determine the source units, which consists in split each code in several fragments. Each fragment is one particular piece of code that will determine a special functionality that can be the source code inside a loop, an if-else statement, or even one procedure or function. This stage is prepared to diminish the complexity of the clone detection algorithm that is strictly dependent on the size of these fragments.

## 2.2    Transformation

In this stage are applied some rules to transform the code in a suitable way to the match detection. All comments and unnecessary white spaces are removed. The identifiers are normalized to force the clone detection in examples where more than one source unit has the same functionality, being different just in the identifiers name. Other structural transformations are currently under study.

## 2.3    High level match detection

All the transformed source units are presented to the high-level match detection algorithm. In this stage, all similar codes of the benchmarks will be aggregated in one useful structure. This structure will contain the number of each approximated clone detection and the position in the original source unit of each one of these clones. The algorithm used to perform this operation is based in the Ukkonen suffix-trees [13] and there are some tools available that implements this algorithm that will be used in this project, such as the ones presented in [14-15].

## 2.4    Post-processing

This stage receives all the clones detected in the prior step and analyzes the possibility for each one to be an RFU candidate. This considers the possible performance acceleration to be achieved and the overall performance improvement based in the functionality of each pattern. All patterns detected that pass in this stage are promoted to RFU candidate to be analyzed in the next block.

## 2.5    Low level match detection

In this stage all RFU candidates will be deeply analyzed. Each piece of source code of the benchmark that will be implemented by an RFU is transformed from the High Level Intermediate Representation (HIR) to the Low Level Intermediate Representation (LIR) using the COmpiler INfra-Structure (COINS) [16]. The results

of these transformations are several graphs, one for each piece of source code that represents the data flow and operations between them.

These graphs are analyzed to recognize the maximal overlap among them and the merge effect on the critical path length. Then, they are classified according to the highest usability and critical path modification. The pieces of code that can take more advantage of the hardware implementation and are present in the benchmark repository more often will stay better classified than pieces of code that will not improve so much the performance when performed in hardware and are not common among other applications.

## 2.6    Suggest new RFUs

Once the classification is done, it is time to suggest the RFUs. The output of our tool will be a dataflow graph and control information obtained by the merge of all the dataflow graphs being considered for each RFU. The final implementation of the RFUs needs a conversion of the dataflow graph in a hardware description language (HDL) and the use of commercial tools to map the hardware structures in the target FPGA. Note that this is out-of-the-scope of this paper and the RFUs implementation for evaluation of the proposed approach can be done using tools able to generate application-specific architectures from software code.

## 2.7    Final Experimental Evaluation

In order to finish the entire process is necessary to validate and evaluate the new computing platforms. This is done executing all the benchmarks with the RFUs and comparing the execution time with a pure platform, without RFUs. Another analysis to be performed concerns the hardware resources needed for the new RFUs in the system. Once the RFUs are generated (using commercial tools), they must be attached to a softcore processor (this will be evaluated with the Xilinx Microblaze[17], RISC processor).

The code sections to be implemented by the RFU are replaced with communication and synchronization primitives. This will be done with the help of the tool that will be developed in the context of this PhD thesis, but is a step that will not be fully automated.

# 3    Related Work

Many research efforts are being applied to automate the efficient translation of source codes to hardware to reach gains in performance and energy savings. Relevant to our work are the approaches using pattern matching and merging of different data-paths, in order to produce hardware units for maximal reuse. In the next paragraphs, we briefly describe some of the most relevant efforts.

Some authors use as input dataflow graphs representing the input application(s). In [18] the authors use techniques to decrease the complexity of the NP-completeness of the pattern matching in several dataflow graphs. They use two main observations. First, we do not need to analyze all possible graphs of one special binary code. In the real world, the size of the patterns will be limited to the number of architectural registers in the ISA (Instruction-Set Architecture) and the size of the patterns produced must have complexity of ten or fewer nodes. Second, each graph is reduced into an unique hash-code and they can set a sample variable that will determine the percentage of sampling that will be used in the compare algorithm, reducing the execution time with a small decrease in accuracy. They used MediaBench[19] and Spec2000int benchmark repositories in the experiments. The results are evaluated showing the coverage of the 100 most popular patterns in all benchmark.

In [6] they show another graph isomorphism technique that is divided in 3 phases: First, potential subgraphs are identified using bounded enumeration. Subgraph isomorphism is then used to remove candidates that are not compatible with the computation acceleration. Finally, unate covering is used to select subgraphs that will be executed on the accelerator. The main contribution of this work is the new algorithms for identifying and mapping subgraphs optimally with intelligent pruning mechanisms. They evaluate the new algorithms' performance and compilation time across a variety of accelerator designs comparing with a traditional greedy approach. To validate they used 23 source codes from Mediabench, Mibench and AVG benchmark repositories. The main experiments were done using an acyclic accelerator with 4 inputs and 2 outputs with 4 intermediate layers(rows) and 15 function units (FUs). The FUs support the complete set of addition, subtraction, and bitwise operators on two inputs. The DFG to be accelerated must fit in this structure and must have only the FUs operations. Another approach is shown in [20]. The authors present a new heuristic method that controls the degree of resource sharing between a given set of custom instructions focused in less die area and energy efficiency losing a bit in latency. They assume that in a previous step, one compiler generated the instruction set extensions and they are represented as a collection of directed acyclic graphs (DAGs) annotated with execution frequency. They use 3 variables (alfa, beta and theta) to decide if each graph will be merged into another graph. The decision to merge two graphs is done based in the tradeoff between the increase in the latency and the decrease in the area (based in the 3 variables). The area and latency values for each path in the graph are calculated summing each operation in the path based in a pre-calculated table with the values for each operation. Changing the variable values, they can reach a wide tradeoff between area and speed. This permits to easily acquire the differences between possible implementations.

In [21] they present a method for achieving any desired balance between flexibility and efficiency by automatically combining any set of individual customization circuits into a larger compound circuit. They observed that the compound circuit cost does not increase in proportion to the number of target applications, due to the wide range of common data-flow and control-flow patterns in programs. They use the UTDSP benchmarks[22] and target accelerators coupled with an embedded PowerPC405 processor. They reached an average speed-up of 2.97× using 3× less

area than the sum of all individual accelerators for each target benchmark. They use loops until 3-deep levels, but, for 2 and 3-deep levels they just do for fixed-bound. They use several different configurations, to force more or less aggressively the circuit merge. When using a highly aggressive merge, less area is used with the disadvantage of longer critical paths.

Our work differs from the above in the use of a higher-abstraction approach able to deal with large sets of benchmarks and making possible both the evaluation of an existent RFU and the suggestion of new RFUs given a set of applications expected to be executed by the computing systems. In the case of the suggestions for new RFUs, our approach can be considered a front-end to the graph-based merge approaches as the ones previously referred.

# 4      Current Status and  Preliminary Results

Currently we are in the initial step of our research and we are putting efforts to find out mechanisms to mine specifics patterns to locate the RFUs candidates.

Figure 3 shows a block diagram of the tool under development. A first stage transforms the source code of each input program into a model represented by the features specified by the user (this step currently uses grammar rules input to JavaCC[23]). Currently, a string consisting of user-specified features (e.g., features can identify keywords such as FOR and WHILE) is generated for each program. Similarities between these strings are then detected.

The tool also permits assisted pattern-mining by using input code patterns, as represented by the boxes in the bottom of Fig. 3.



**Fig. 3.** Block diagram of the pattern-mining tool.

The preliminary version of the tool is able to analyze 234 source files (C language) of the MiBench[24]  and MediaBench[19] benchmark repositories. At the moment, the tool transforms the string of features for each input file in a suffix tree. Operations are performed in the suffix tree to recover the patterns that are repeated in this

collection of C programs, and probable candidates for RFUs. The tool is currently able to accept a set of features specified by the user and one or more patterns based on sequences of those features. For instance, in Fig. 4 a) are illustrated three different loop patterns with high parallelization possibility that could be RFU candidates and in Fig. 4 b) is showed the amount of this pattern mined in the analyzed benchmarks.



a)



b)

**Fig. 4.** a)Mined patterns in the benchmarks b) presence of the patterns in the Benchmark

In the Table 1 is summarized some information about the source codes analyzed and mined information. We can see that in the source codes analyzed the tool is able to find patterns, such as loops that are presented 1654 times in the 234 source codes. Another retrieved information are the patterns defined in the Fig. 4.a) that are presented 766 times. Now we are trying to detect other kind of shapes and study the effects of these shapes in the hardware generation and posterior performance improvements. Other efforts are being put in two clone detection tools[14-15] where we are trying to find patterns that happen with higher frequency in the benchmarks.

**Table 1.** Preliminar mined informations

|  | MiBench Automotive | MiBench Network | MiBench | MediaBench | MiBench + MediaBench |
|---|---|---|---|---|---|
| Total Files (*.c) | 17 | 4 | 1087 | 711 | 1798 |
| Files Analyzed | 17 | 4 | 140 | 94 | 234 |
| #Kernels | 17 | 3 | 392 | 192 | 584 |
| DNA Length | 17419 | 2201 | 84098 | 151986 | 236084 |
| WHILE | 16 | 5 | 181 | 132 | 313 |
| DO WHILE | 1 | 4 | 28 | 21 | 49 |
| FOR | 84 | 15 | 337 | 955 | 1292 |
| #Loops | 101 | 24 | 546 | 1108 | 1654 |
| FOR deep 1 | 55 | 13 | 266 | 586 | 852 |
| FOR deep 3 | 4 | 0 | 6 | 66 | 72 |
| FOR deep 5 | 4 | 0 | 4 | 28 | 32 |
| Pattern 1 | 28 | 2 | 67 | 125 | 192 |
| Pattern 2 | 23 | 3 | 110 | 319 | 429 |
| Pattern 3 | 15 | 2 | 44 | 101 | 145 |

## 5    Conclusion and Future Work

We believe that pattern mining can be extremely important for the mapping and design of reconfigurable functional units. Although in a preliminary stage, our tool is currently able to suggest pattern similarities among collections of many benchmarks.

Currently, we are working on more advanced methods to find code segments with high potential for being mapped to RFUs tightly coupled to a softcore microprocessor and both implemented in an FPGA. This stage also needs models to estimate the hardware resources and the latency achieved by possible RFUs for the patterns being suggested.

The next step is to devise layers of representation models that can be easily used to find similarities at different abstraction levels using, e.g., clustering algorithms. We believe that a top-bottom approach will permit to isolate highest potential candidates using a first rough and less computational intensive step. Then refinement steps can be done to the candidates already identified in previous stages.

Our future research work will focus on the following issues:
- Devise layers of representation that can be easily used to find similarities at different abstraction levels using, e.g., clustering algorithms;
- Analyze how certain aspects of a program can be extracted and represented in a model ready to be efficiently used by pattern-mining tools;
- Propose new techniques for pattern-mining aware of execution time and hardware resources and evaluate results obtained using existent data mining tools.

- Elaborate a domain specific language (DSL) to specify the patterns suitable for a specific target reconfigurable functional unit.

# References

1. Moore, G.E., Cramming more Components onto Integrated Circuits. Electronics, 1965. 38, Number 8: p. 4.
2. Intel.      [cited 2010 January]; http://www.intel.com/products/processor/core2quad/-index.htm.
3. AMD.    [cited 2010 January];  http://www.amd.com/uk/products/server/processors/-six-core-opteron/Pages/six-core-opteron.aspx.
4. picoChip.  [cited 2010 January]; http://www.picochip.com/page/76/.
5. Mihai Budiu, et al., Spatial computation, in 11th International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS 2004. 2004: Boston, MA, USA. p. 14-26.
6. Nathan Clark , A.H., Scott Mahlke ,  Sami Yehia, Scalable Subgraph Mapping for Acyclic Computation Accelerators, in International Conference on Compilers, Architectures, and Synthesis for Embedded Systems. 2006: Seoul, South Korea.
7. Jelena Trajkovic and D.D. Gajsk, Custom Processor Core Construction from C Code, in Sixth IEEE Symposium on Application Specific Processors. 2008: Anaheim, California
8. Hartenstein, R., A decade of reconfigurable computing: a visionary retrospective, in Design, Automation and Test in Europe, 2001. Conference and Exhibition 2001. 2001: Munich, Germany. p. 642-649.
9. Becker, J., Configurable Systems-on-Chip (CSoC), in 15th symposium on Integrated circuits and systems design (SBCCI), Invited Tutorial. 2002, IEEE Computer Society: Los Alamitos, CA, USA.
10. Martin Vorbach and J. Becker, Reconfigurable Processor Architectures for Mobile Phones, in International Parallel and Distributed Processing Symposium (IPDPS'03). 2003: Nice, France.
11. Philip Brisk, Adam Kaplan, and M. Sarrafzadeh, Area-efficient instruction set synthesis for reconfigurable system-on-chip designs, in 41st annual Design Automation Conference. 2004, ACM: San Diego, CA, USA.
12. Chanchal K. Roy, James R. Cordy, and R. Koschke, Comparison and evaluation of code clone detection techniques and tools: A qualitative approach. Science of Computer Programming, 2009. 74(7): p. 43.
13. UKKONEN, E., On-line construction of suffix trees, in Algorithmica. 1995.
14. Copeland, T., PMD Applied, ed. C. Books. 2005. 221.
15. Florian Deißenböck, Markus Pizka, and T. Seifert, Tool-supported realtime quality assessment, in 13th International Workshop on Software Technology and Engineering Practice 2005 (STEP-05). 2005, IEEE Computer Society. p. 127-136.
16. COINS.  [cited 2010 January]; http://www.coins-project.org/international/.
17. Xilinx. Microblaze Soft Processor Core.    [cited 2010 January];   http://www.xilinx.-com/tools/microblaze.htm.
18. Peter G. Sassone and D.S. Wills, On the extraction and analysis of prevalent dataflow patterns, in 7th Workshop on Workload Characterization,. 2004. p. 8.

19. Chunho Lee, M. Potkonjak, and W.H. Mangione-Smith, MediaBench: a tool for evaluating and synthesizing multimedia and communications systems, in International Symposium on Microarchitecture. 1997: Research Triangle Park, NC, USA. p. 330-335.
20. Marcela Zuluaga and N. Topham, Resource Sharing in Custom Instruction Set Extensions, in Symposium on Application Specific Processor - Sasp. 2008: Anaheim, CA, USA. p. pp.7-13.
21. Sami Yehia, et al., Reconciling Specialization and Flexibility Through Compound Circuits, in International Symposium on High-Performance Computer Architecture (HPCA). 2009.
22. Lee, C.G. The UTDSP Benchmark Suite. [cited 2010 January]; http://www.eecg.-toronto.edu/~corinna/DSP/infrastructure/UTDSP.html.
23. JavaCC. [cited 2010 January]; Available from: https://javacc.dev.java.net/.
24. M. R. Guthaus, et al., MiBench: A free, commercially representative embedded benchmark suite, in Workload Characterization, 2001. WWC-4. 2001. 2001, IEEE Computer Society: Austin, TX, USA. p. 3-14.

# Automatic Identification of Stage Transitions in Protein Unfolding Simulations

Md. Anishur Rahman[1],

[1] FEUP, Universidade do Porto
Rua Dr. Roberto Frias s/n, 4200-465 Porto, Portugal
pro090024@fe.up.pt

**Abstract.** The study of the unfolding process of proteins may highlight valuable information for the cure of a wide range of diseases. In this paper we have focus our attention on the TTR protein which, when degenerated, may cause one type of amyloid disease called the doença dos pézinhos. We have analysed data from the dynamic simulation of the unfolding process of the TTR protein. In the experiments we have compared the behaviour of the wild-type TTR (non malignant type) with the L55P-TTR (that causes the disease). We have applied a clustering algorithm to identify stage transitions in the unfolding process. Preliminary results obtained so far are not yet satisfactory according to a domain expert.

**Keywords:** Data Mining, Clustering, Protein Unfolding.

## 1 Introduction

A protein is a long chain molecule made up of amino acids joined by peptide bonds. There are more than 50,000 different proteins in our body. Proteins participate in almost all chemical reactions in our body and their function is highly determined by their 3-Dimensional shape. It is known however that small changes in the linear sequence of amino acids (also designated as residues) may change the 3-D shape of a protein and can cause incurable diseases. It is therefore of great importance to study and understand the process of protein folding. Protein folding is one of the most fundamental problems in modern Molecular Biology. The protein misfolding has been related with a huge number of human diseases such as Alzheimer`s, Parkinson`s and Huntington`s. The protein folding problems concerns the finding of rules to predict the 3D shape of a protein based on the linear sequence of amino acids (residues) of the protein (1D structure). The functionality of proteins is determined by their 3D structure and proteins are responsible by more than 99 percent chemical reaction in leaving organism. Solving the protein folding problems will be a major step in solving diseases like Alzheimer, be able to design from scratch the medicines, etc.
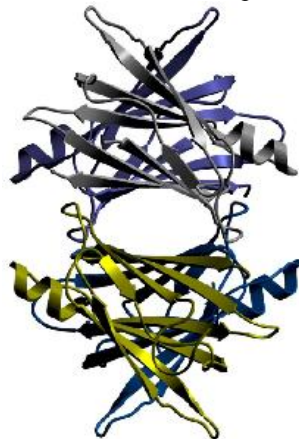
The protein unfolding problem is an alternative approach to address the protein folding problem. In protein unfolding a protein is put into boiling water or in an acidic container and the process of how its structures breaks may provide valuable information for the mechanism of protein folding.

In our study we will be analysing data collected during the simulation of a protein unfolding. During the simulation of proteins there is a large amount of data that is generated that is very hard (if not impossible) to be analysed by human experts. We need Data Mining algorithms to extract useful information from such simulation data. We will be analysing data from the unfolding simulation of a protein (TransThyRetin - TTR) that (when degenerated) gives rise to the "feet disease" (and several others)[1][2][3]. From the data collected during the simulation we want to automatically discover stage transitions of the unfolding processes. It is also a goal of our study to find out if there is significant differences in the unfolding of the \normal" protein (wild-type) WT-TTR and the malignant one L55P-TTR.

The rest of the paper is organised as follows. In Section 2 we present the domain of protein unfolding simulation. In Section 3 we define the events used to detect the stage transitions in the unfolding process. We relate our work with previous work done by others in Section 4. The experiments we have carried out, together with the results obtained, are presented in Section 5. Conclusions and further work are presented in Section 6.

## 2 Protein Unfolding Simulation

K TTR is a homotetrameric protein (see Figure 1) with a total molecular mass of 55 kDa, 127 amino acid residues per subunit, and a high percentage of Beta-sheet. One way of exploring the unfolding events that may be responsible for TTR aggregation is through the use of molecular dynamics (MD) protein unfolding simulation. In MD simulations, one tracks the atoms in only one protein molecule as a function of time. To explore the unfolding routes of monomeric species of TTR, several molecular dynamics simulations of the TTR subunits of wild-type (WT-TTR) and the variant L55P-TTR (with a Proline replacing a lleucine in position 55)were performed, at high temperatures. Among the numerous pathogenic variants, L55P-TTR is the most amyloidogenic, and V30M-TTR is one of the most prevalent.



**Fig. 1.** TTR protein. It is composed of 4 chains, each one with 127 residues.

It is possible to predict protein structure from basic physical-chemical interactions by using molecular dynamics (MD) simulations. In MD simulation studies of proteins, the time-dependent behaviour of the molecular system is obtained by integrating Newton`s equations of motion (classical mechanics) and the potential energy function (a.k.a. force fields). The result of the simulation is a time series of conformations; this is called a trajectory or the path followed by each atom in accordance with Newton`s laws of motion. MD simulations of proteins folding have been successfully performed for small proteins (less than 50 amino acid residues).

These MD simulations are computationally expensive and generate a huge amount of data, making the comparison of different trajectories a difficult task. Data Mining techniques will provide the variation of molecular properties related with each simulation.

In this paper we have used clustering analysis techniques to detect clusters of events that suggest stage transitions during the unfolding process. We have applied k-means clustering to events on three measured properties: the variation of the accessible surface area (SASA) of each amino-acid residues along each MD simulation; breaks of secondary structure and; events associated with the distance of the residues to the Mass Centre of the protein. We now describe the events in more detail.

## 3 Events in the Protein Unfolding Process

To investigate the possibility of defining stage transitions in the unfolding process we have defined a set of events that could be useful to establish the stages boundaries. The events were defined according to the properties computed from the simulation traces. In this study we used the three properties mentioned in Section 2: Solvent Accessible Surface Area (SASA); secondary structure and; distance from each residue to the Mass Centre of the protein.

For the secondary structure we looked at instants where existing structures ($\alpha$-helices and $\beta$-sheets) are severely damaged or new ones appear. Figure 2 shows the evolution of secondary structures in the simulation of the wild-type (WT-TTR). We will call this type of events ss-events (SSE).

For the MC property we collected two types of events. One type, that we will call cm-crossing-event (CCE) occurs when the trajectory of two residues cross each other. That is when one of them is approaching the MC of the protein and the other is moving away from the MC and they cross. A second kind of events, that we will call tendency-change-event (TCE) occurs when a residue changes its trajectory. either its was approaching the MC and then moved away or the revers change in the movements.

Regarding the SASA we have defined one type of events that we call sasa-event (SASAE). A sasa-event is signalled by the change in the trend of SASA values. Either a change from an sequence of increasing values to a sequence of decreasing value or the reverse situation.

**Fig. 2.** Evolution of the secondary structure of WT-TTR during one unfolding simulation. The xx axis represents simulation time. The yy axis represents the position of a residue in the protein (only positions between 10 and 51 are represented). Thick lines indicate that the residue belongs to a beta sheet, thin lines indicate that the residue belong to a alpha-helix. We can see that (on top of the picture) the beta sheet between position 41 and 50 looses a substantial amount of residues near simulation time 5000. We can also see that a alpha-helix appears near simulation time 9000 between positions 17 and 23 (near the bottom right side of the picture).

## 4 Related Work

In [2] Brito et al. discuss important issues of protein stability, folding and aggregation have become central in several pathological conditions and in particular in amyloid diseases. Here, Brito et. al. have reviewed the recent developments on the molecular mechanisms of amyloid formation by transthyretin (TTR), in particular, in what concerns to protein conformational stability, protein folding and aggregation. Today, more than 80 TTR mutations throughout the TTR sequence are known. In [4] Rodrigues et al. discuss that the association between amyloid filbril formation deposition with a series of diseases, including Alzheimer`s, Spongiform encephalopathies, and several systemic amyloidosis. In most of these amyloid diseases, it has been shown that the normal precursor protein, due to proteolysis, mutation or molecular environment stress, undergoes misfolding, leading to molecular species with a high tendency for ordered aggregation into amyloid. However, the structural nature of theses amyloidogenic intermediates is the subject of debate. The necessity of applying Data Mining to the analysis of Protein folding and unfolding simulations was suggested earlier by Brito et al. [5]. Using molecular dynamics in unfolding simulations of an amyloidogenic protein-transthyretin as an example, they

put forward a series of ideas on how simulations of this type may be used to infer rules and unfolding behaviour in amyloidogenic proteins, and to extrapolate rules for protein folding in different structural classes of proteins. They argue that it could help in the development of protein structure prediction methods.

Analysing the data from the Unfolding simulations has received a lot of attention recently. Several Data Mining perspectives have been adopted.

In [6] Azevedo et. al. demonstrated the use of Association Rules applied to the analysis of the variation profiles of the Solvent Accessible Surface Area of the 127 amino-acid residues of the protein Transthyretin, along multiple simulations.

In [7], Ferreira et al. looked for motifs taking the simulation logs as a time series data.

In Camacho et al.[8], ILP was used to extract rules that help in the explanation of the break down process of secondary structures. In that study Inductive Logic Programming was used, in a classification problem, to discover rules that could be useful in the explanation of the breaking the secondary structures. Although there is a similarity with the work reported in this paper in the sense that they addressed the breaking of secondary structures but they looked at it as a classification problem not a clustering one and they were not studying stage transitions.

In [9], Fernandes et al. presents a method of knowledge discovery in data obtained from Molecular Dynamics Protein Unfolding Simulations. In the line of our work Fernandes et al. also used clustering and defined events. They worked only with the Mass centre events and used cluster methods as a process to achieve data reduction. They were not concerned with stage transitions of the unfolding process.

Vilaça et al.[10], used graph mining techniques to analyse the simulation traces of the unfolding of TTR protein. They used such techniques to trace the persistence of small fragments of protein during the unfolding process.

Brito et al[11], propose the application of an augmented version of hierarchical clustering analysis to detect clusters of amino-acid residues with similar behaviour in protein unfolding simulations. These clusters hold similar global pattern behaviour of solvent accessible surface area (SASA) variation in unfolding simulations of the protein Transthyretin (TTR). Classical hierarchical clustering was applied to build a dendrogram based on the SASA variation of each amino-acid residue. The dendrogram was enriched with background information on the amino-acid residues, enabling the extraction of sub-clusters with well differentiated characteristics.

## 5 Experiments

### 5.1 Experimental Settings

The experiments we have designed have two objectives. First, to find clusters identifying unfolding stage transitions instants. Second, identify differences between the break down process in the WT (wild type) and L55P protein variants. This later goal is very important since it may contribute to an explanation for the malignant behaviour of L55P-TTR.

We have used data from simulations of the protein transthyretin (TTR)[1]. The simulations [5] include 5 runs using the wild type (WT) and 5 others using the amyloidogenic type (L55P). The simulations were done using the GROMACS [12][2] program and simulate the breaking of the TTR protein in boiling water as described in [5]. Each simulation lasted for 10000 instants of time and therefore we have 127*10000*NProps[3] values recorded for each simulation. In each simulation run we collected, for each residue and instant of time, information concerning the secondary structure it belongs to, its Solvent Accessible Surface Area (SASA) value and the distance from the carbon-$\alpha$ to the Mass Centre of the protein. The amount of data resulting from the all simulations is nearly 1.5 GB of log traces. We have used the SimpleKmeans algorithm for clustering. SimpleKmeans is the WEKA [13] implementation of the k-means clustering algorithm.

We have developed a script that collects all events defined in Section 3 from the simulation traces of WT-TTR and L55P-TTR. We have planned 12 experiments. Three experiments for each type of event and 3 others for the all set of events. The number three in each group considers running both variants of the protein and then the two variants in separate runs.

## 5.2 Experimental Results

At the moment we have only used events of type SSE in WEKA. There are 1255 events, 482 were found in the L55P variant and 773 in the WT variant. We have run Weka's k-mean implementation with K ranging from 2 to 7. We were investigating if there are 2, 3, 4, 5, 6, 7 or 8 stages during the unfolding process. We have repeated the experiments with all the SSE events, with the SSE of the WT variant alone and with SSE events of the L55P variant alone. The quality of each clustering was assessed by the Weka's Within cluster sum of squared errors measure. The best K was the k that produced the smaller value for the error measure and none of the clusters has less than 9% of the data. We have used the Euclidean Distance as one of Weka's SimpleKmeans parameter.

So far we have just manage to analyse the SSE events. The results are sumarised in Table 1 and are not yet totally satisfactory according to a domain expert.

---

[1] Reference 1TTA in the PDB (http://www.rcsb.org/pdb/home/home.do)
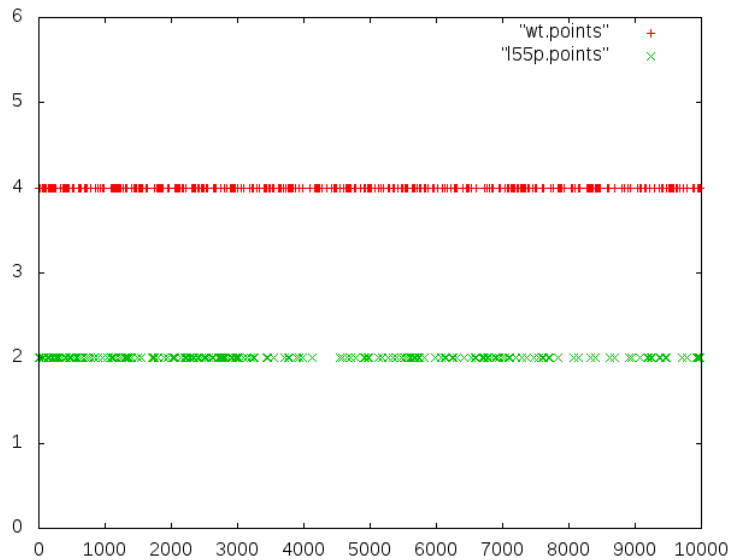[2] Groningen Machine for Chemical Simulations
[3] 127 is the number of residues and NProps is the number of properties measured

**Table 1.** Results Results of the analysis of the Secundary Structure break events. (a) is the results on the data set with all (WT and L55P TTR) the events. (b) is the results on the WT-TTR data set and (c) is the results on the L55P-TTR data set. WCSSE is Weka's error mesure \Within Cluster Sum of Squared Errors".

| K | WCSSE | Smaller group | K | WCSSE | Smaller group | K | WCSSE | Smaller group |
|---|-------|---------------|---|-------|---------------|---|-------|---------------|
| 2 | 354.3% | 38% | 2 | 197.2% | 46% | 2 | 96.6% | 46% |
| 3 | 237.6% | 11% | 3 | 170.4% | 23% | 3 | 54.0% | 9% |
| 4 | 196.8% | 11% | 4 | 149.4% | 20% | 4 | 31.7% | 9% |
| 5 | 165.5% | 11% | 5 | 141.3% | 15% | 5 | 28.5% | 9% |
| 6 | 140.6% | 11% | 6 | 136.1% | 12% | 6 | 23.7% | 9% |
| 7 | 128.3% | 9% | 7 | 133.5% | 8% | 7 | 22.1% | 3% |

We have not yet found clear stage transitions that would be found biologically significant by the domain expert. The results show, however, two promising results: analysing a data set with the two variants of TTR leads to worse results and; the two variants seem to actually behaving differently as the secondary structure is concerned. This later result masy explain the former one but is quite consistent with the biological explanation of the amyloid deseases. In amyloid deseases there is an excess of fibers that deposit in undesirable places and are a result of TTR misunfolding process as happens with the L55P-TTR. We, therefore, find it encouraging to have found significant differences between the breaking of secondary structures in the two variantes of the TTR protein. The distinction is highlighted in Figure 3 where you can see a significant differences in the number of events and distribution after the instant 3200 of the simulation, between WT and L55P variants.



**Fig. 3.** Events signaling the breaking of secondary structures of WT-TTR (in the upper level) and L55P-TTR (in the lower level). The xx axis represents the instant in the simulation ([0-10000] nano-seconds). The yy axis is artifitialy set just to differentiate the events of the two variants.

# 6 Conclusions and Future Work

We have applied Data Mining techniques to automatically discover stage transitions in data from the simulation of TTR protein unfolding. To that purpose we have defined 4 types of events occuring associated with the 3 properties measured in the simulations: secondary structure, SASA and; MC distance. We have developed a script that analyses the simulation logs and automatically extracts the events. We used the k-means implementation in Weka to search for meaningful stage transitions (groups of events). At this point our finding are not yet satisfactory according to a domain expert.

The ongoing and future work concerns the completion of the planed set of 12 experiments. We also want to increase the attributes given to each event to find out if better generalisations are possible. One last line of work will consider a Relational clustering algorithm in a comparative study with Weka's implementation of k-means.

# References

1. Quintas, A., Vaz, D.C., Cardoso, I., Saraiva, M.J., Brito, R.M.M.: Tetramer dissociation and monomer partial unfolding precedes protofibril formation in amyloidogenic transthyretin variants. J. Biol. Chem. 276 (2001) 27207--27213

2. Brito, R., Damas, A., Saraiva, M.: Amyloid formation by transthyretin: From protein stability to protein aggregation. Current Medicinal Chemistry Immun.Endoc. & Metab. Agents 3: (2003) 349--360

3. Rodrigues, J., Brito, R.: Amyloid formation by transthyretin: How much unfolding is required ? in "amyloid and amyloidosis". J. Biophys. 86(1) (2004) 323--325

4. Rodrigues, J.R., Brito, R.M.M.: How important is the role of compact denatured states on amyloid formation by transthyretin? Amyloid and Amyloidosis. CRC Press. (2004) 323--325

5. Brito, R., Dubitzky, W., Rodrigues, J.: Protein folding and unfolding simulations: A new challenge for data mining. OMICS: A Journal of Integrative Biology. 8(2) (2004) 153--166

6. Azevedo, P.J., Silva, C.G., Rodrigues, J.R., Loureiro-Ferreira, N., Brito, R.: Detection of hydrophobic clusters in molecular dynamics protein unfolding simulations using association rules. Proc. 6th International Symposium ISBMDA 2005, Lect. Notes in Comput. Sc. 3745 (2005) 329--337

7. Ferreira, P.G., Azevedo, P.J., Silva, C.G., Brito, R.M.M.: Mining approximate motifs in time series. Lect. Notes Art. Intl. 4265 (2006) 77--89

8. Camacho, R., Alves, A., Silva, C., Brito, R.: On mining protein unfolding simulation data with inductive logic programming. 2nd International Workshop on Practical Applications of Computational Biology and Bioinformatics (IWPACBB2008) (2008) 175--179

9. Fernandes, E., Jorge, A., Silva, C., Brito, R.: A knowledge discovery method for the characterization of protein unfolding processes. 2nd InternationalWorkshop on Practical Applications of Computational Biology and Bioinformatics (IWPACBB2008) (2008) 180--188

10. Vilaça, J.: Extracção e análise de fragmentos frequentes em simulação de dinâmica molecular de desnaturaçãoo proteica. Master's thesis, MSc thesis, Universidade do Minho (2009)

11. Ferreira, P., Silva, C., Brito, R., Azevedo, P.: A closer look on protein unfolding simulations through hierarchical clustering. Proceedings of the 2007 IEEE Symposium on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB07) (2007) 461--468

12. van der Spoel, D., van Buuren, A.R., Apol, E., Meulenhoff, P.J., Tieleman, D.P., Sijbers, A.L.T.M., Hess, B., Feenstra, K.A., Lindahl, E., van Drunen, R., Berendsen, H.J.C.: Gromacs User Manual version 3.1, Nijenborgh 4, 9747 AG Groningen, The Netherlands. (2001)
13. Witten, I.H., Frank, E.: Data Mining: Practical machine learning tools and techniques. 2nd edn. Morgan Kaufmann (2005)

# The impact of Pre-Processing in Clustering MEDLINE Documents

Carlos Adriano Oliveira Gonçalves

FEUP - Faculdade de Engenharia da Universidade do Porto
carlos.adriano@fe.up.pt

**Abstract.** The amount of information available in the Medline database makes humanly impossible to select the relevant documents for a researcher to read. Clustering of documents may be a valuable technology to help handling such amount of documents. To accomplish this process it is of capital importance to use appropriate pre-processing techniques on the data. In this study the main goal is to analyse the impact of pre-processing techniques in Clustering. We will use Clustering algorithms available in the WEKA tool to group the documents and evaluate the results according to MeSH terms. Our first experiments show that the application of pruning, stemming and wordnet reduces significantly the number of attributes without affecting the accuracy of results.

**Key words:** MEDLINE, Clustering, Machine Learning, Pre-Processing

## 1 Introduction

Molecular biology and biomedicine scientific publications are available (at least the abstracts) in Medical Literature Analysis and Retrieval System On-line (MEDLINE). MEDLINE is the U.S. National Library of Medicine's (NLM) premier bibliographic database that contains over 16 million references to journal articles in life sciences with a concentration on biomedicine. A distinctive feature of MEDLINE is that the records are indexed with NLM's Medical Subject Headings (MeSH terms). MEDLINE is the major component of PubMed [1], a database of citations of the National Library of Medicine in the United States. PubMed comprises more than 19 million citations for biomedical articles from MEDLINE and life science journals. The result of a MEDLINE/PubMed search is a list of citations[1] to journal articles. With this huge amount of information, it is humanly impossible to read and select relevant information from a database such as MEDLINE. As this is a very interesting and actual topic of investigation the main idea of this work is to download a

---

[1] including authors, title, journal name, the abstract of the paper, keywords and MeSH terms

MEDLINE sample (freely available at the NCBI site) and to apply text clustering to that huge amount of information. In this paper we explore the effects of using different pre-processing techniques in classifying MEDLINE documents. The main objectives of this work are:

– to make a pre-processing of the original *Data Set*
– to apply Clustering using the WEKA tool to the different Data Sets obtained in the previous step.

The rest of this paper is structured as follows. In Section 2 we overview the pre-processing techniques used. In Section 3 we present the text classification problem and existing techniques. We report on related work in Section 4. Section 5 describes the experiments and Section 6 the results obtained. Section 7 concludes the paper presenting the conclusions and future work.

## 2 Pre-Processing Techniques

In this section we will make a small resume of the most common pre-processing techniques used.

Before applying any data analysis technique it is necessary to pre-process the collection of documents (our data set). The pre-processing techniques are the ones also used in Information Retrieval, namely:

– **Tokenization** is the process of breaking a text into tokens. A token is a non-empty sequence of characters, excluding spaces and punctuations.

– **Stop Word Removal** removes words that are meaningless such as articles. conjunction and prepositions (e.g., a, the, at, etc.). These words are meaningless for the evaluation of the document content.

– **Stemming** is a widely used technique in text analysis. Stemming is the process of removing inflectional affixes of words reducing the words to their stem.

– **Pruning** discards terms either appearing rarely or "too frequently". Terms that rarely appear in a document or terms that appear too frequently do not contribute to identify the topic of the document. The most common techniques are term frequency and document frequency.

– **Treating synonyms**: the possibility to take care of synonyms may be seen as another pre-processing technique. If two words or terms mean the same thing, e.g, if they are synonyms we could replace them by one of them without taking the semantic meaning of the term.

– **Document representation:** it is necessary to follow a model representation of the document. In the standard Information Retrieval literature [2] there are three known approaches: the Boolean Model, the Probabilistic Model and the Vector Model. All these approaches represent documents as vectors of unique terms:

$$D_i = < T_1, T_2, ..., T_N > \tag{1}$$

A collection of documents is represented as a set of vectors that can be written as a matrix, that is called document-matrix. We provide an example in Table1. $W_{ij}$ (i $\in$[1,M]; j $\in$ [1,N]) represents the weight of term $T_j$ in document $D_i$. In the vector space model there are several variations to attribute the weights to the terms: TF (Term Frequency) and TFIDF (Term Frequency Inverse Term Frequency) are the most widely used and the last one is the one we used. In the Vector Space Model documents are represented as vectors in a high dimensional space. As already mentioned we have used the TFIDF method where according to [3]

$$TF(term, document) = the frequency of term in document \tag{2}$$

and

$$IDF = \log \frac{number of documents in collection}{number of documents with term} + 1. \tag{3}$$

**Table 1.** Matrix-document

|       | $T_1$    | $T_2$    | ...  | $T_N$    |
|-------|----------|----------|------|----------|
| $D_1$ | $W_{11}$ | $W_{12}$ | ...  | $W_{1N}$ |
| $D_2$ | $W_{21}$ | $W_{22}$ | ...  | $W_{2N}$ |
| ...   | ...      | ...      | ...  | ...      |
| $D_M$ | $W_{M1}$ | $W_{M2}$ | ...  | $W_{MN}$ |

These pre-processing techniques are of utmost importance once they prepare the dataset for Clustering. Stop word removal, stemming and pruning improves cluster quality once they remove the meaningless data that leads to a reduction in the number of dimensions in the term space. Document representation concerns the estimation of the importance of a specific term in the document. As the number of features (attributes) are very huge in a collection of documents, the pre-processing techniques help in reducing these huge number of features. Pre-processing is used as a traditional Machine Learning feature selection technique.

## 3   Text Classification

Text Classification attempts to automatically determine whether a document or part of a document has particular characteristics of interest, usually based on whether the document discusses a given topic or contains a certain type of information [4]. Text Categorisation involves two main research areas: Information Retrieval and Machine Learning. The first of step of Text Classification is to transform documents into a suitable representation for the Classifier. For this, and before applying the Classifier documents must be pre-processed using Information Retrieval techniques mentioned in the previous section. Some other pre-processing techniques that can be applied to the collection of documents, in order to reduce the huge amount of terms in the collection is called *Feature Selection*.

A Feature Selection or feature extraction phase is needed to reduce the dimensionality of the document.

There are three Classification techniques: supervised learning, unsupervised learning and semi-supervised learning.

Supervised Learning is based on a training set of examples, where the key idea is to learn from a set of labelled examples (the training set).

In unsupervised learning there is no a priori output. The goal of unsupervised learning is to learn a model that explains well the data. Usually, the result of unsupervised learning is a new explanation or representation of the observation data, which will then lead to improved future decisions.

Semi-supervised learning makes use of both labelled and unlabelled data (typically a small amount of labelled data and a large amount of unlabelled data. Semi-supervised learning [5] [6] is a machine learning paradigm in which the model is constructed with a small number of labelled instances and a large number of unlabelled instances. One key idea in semi- supervised learning is to label unlabelled data using certain techniques and thus increase the amount of labelled training data.

## 4   Related Work

The authors in [7] propose clustering MEDLINE documents using a method that integrates the semantic information embedded in MeSH thesaurus and the content information of texts for enhancing the performance of document clustering. In this work the authors also apply pre-processing techniques before applying clustering. to the data (MEDLINE documents). The pre-processing procedures used are the removing of stop words and the Porter's Stemmer algorithm; they remove the stemmed words that occur in less than three documents. Clustering is based on the three Medline fields: title, abstract and MeSH main heading. Spectral clustering methods are used for clustering. The authors, based on their experiences, claim

that the use of a combined similarity measure between two documents recurring to MeSH thesaurus improves document clustering.

The authors in [8] describe a system that automatically clusters PubMed query results into various groups where each group contains relevant articles, extracts the most common topic for each group, and ranks the articles in each group. Pre-processing techniques were applied: tokenization, stop word removal, stemming and the standard term frequency inverse document frequency was used to represent the document-matrix that was normalised using cosine normalisation. The authors use the title, abstract and the Mesh terms fields Before applying the clustering algorithms the author's returned from a PubMed query result. The authors used a their own novel algorithm (*spectroscopy*) that predicts the clustering characteristics of a text collection estimating the number of clusters. The authors used CLUTO [2] software using the bisecting k-means algorithm.

In [9] the authors present an algorithm for large-scale document clustering of biological text. The algorithm is based on statistical treatment of terms, stemming, the idea of a âgo-listâ, unsupervised machine learning and graph layout optimisation.

## 5 Empirical Study

Our base of work will be a MEDLINE sample. As we do not know the Data Set and as we do not have a set of positive and negative examples for training, the first idea is to apply clustering in order to obtain clusters (some algorithms such as K-Means will be tested with different values of K and an analyse of the best results will be made).

The focuses of this research is to study the impact of pre-processing techniques in a MEDLINE sample Data Set.

### 5.1 Data Set Characterisation

The data set that is subject of our study, is a MEDLINE sample that was downloaded from the NLM site at (ftp://ftp.nlm.nih.gov/nlmdata/sample/medline/). This sample has 53.2 MB, and contains 30000 citations. The sample is in the XML format. Each citation contains several information namely: the pmid (the pubmed id), the journal title, the pubmed date, the article title, the abstract of the paper if available, the list of authors, the list of keywords and the list of Mesh terms. A MeSH (Medical Subject Headings) is the (U.S.) National Library of Medicine's controlled vocabulary thesaurus. It consists of sets of terms naming descriptors in a hierarchical structure that permits searching at various levels of specificity. A MeSH term is a medical

---

[2]  CLUTO is a software package for clustering low and high dimensional data sets and for analysing the characteristics of several clusters.

subject heading, or descriptor, as defined in the MeSH thesaurus. We have used the ten MeSH terms with higher frequency.

## 5.2  Pre-Processing techniques applied

The MEDLINE sample we used in our study was in the XML format. So the first pre-processing step was to read the XML file and to filter the information we need namely the the pmid (the pubmed id), the journal title, the pubmed date, the article title, the abstract of the paper if available, the list of authors, the list of keywords and the list of Mesh terms. We have developed our application using the JAVA Programming Language which is familiar to us. We have also used a MySQL Database to store all the information useful for further pre-processing and Classification. We have applied the following pre-processing techniques to our original Data Set.

- Tokenization: which is the process of splitting a text document into a stream of words;
- Stop Words removal: removes words that are meaningless (such as articles, conjunctions, prepositions, etc.);
- Stemming: is the process of reducing a word to it's root (we used the Porter's Stemmer Algorithm [10]
- We have used WordNet [11] to search for synonyms of terms and replace each term for the respective synonym;
- We have implemented the standard term-frequency inverse document frequency (TFIDF) function to assign weights to each term in the document.

Besides this most common pre-processing techniques we have applied some other pre-processing features with the objective of reducing the number of attributes, namely:

- Pruning:
  - We have removed the words that appear once or twice in a document, because if their frequency is so low is because they are not discriminative of the document topic;
  - We have also removed the words that appear less than 10 times in the all collection, for the same reason, e.g., because they are not discriminative of the document content;

In the vector space model there are several variations to attribute the weights to the terms: TF (Term Frequency) and TFIDF (Inverse Term Frequency) are the most widely used and the last one is the one we used. In the Vector Space Model documents are represented as vectors in a high dimensional space. As already mentioned we have used the TFIDF method because in this weighting scheme terms that appear too rarely or too frequently are ranked lower than terms that balance

between the two extremes. And also because higher weight terms signify that the term contributes better to clustering results.

The Figure 1 summarises our approach.



**Fig. 1.** Summary of steps

### 5.3   Classification through the use of Clustering

The fundamental idea of Text Clustering is to divide a collection of documents into different group categories so that documents in the same category describe the same topic.

The main reason for using a clustering (unsupervised learning) is because we do not know the Data Set, so it is impossible to make a *a priori* classification. Besides unsupervised learning adjusts itself to the cases when we do not know the data set. Text clustering involves [12]:

  – document representation (may involve feature selection or extraction)
  – definition of a document similarity measure
  – a clustering algorithm

The similarity measures measure how similar are two documents or how similar a query is to a document. The most known similarity measures are the Euclidean Distance, the cosine similarity and the Manhattan Distance.

In this step, Data Mining Algorithms namely text clustering algorithms will be applied to classify the documents into categories.

Most of document clustering approaches are based on the vector space model and apply several clustering algorithms to the representation that can be divided into hierarchical and partitional [13]. Hierarchical clustering algorithms successively merges the most similar objects based on the pairwise distance between objects until a termination condition is satisfied. According to [14] Hierarchical clustering is limited because of its quadratic time complexity. Partitional clustering algorithms (especially K-Means) are the most widely used algorithms in document clustering. K-means and its variants have a time complexity which is linear in the number of documents, but are thought to produce inferior clusters. Clustering algorithms will be applied using the WEKA tool.

## 6  Experimental Results

**Table 2.** Machine Learning algorithms used in the study.

| Algorithm | Type |
|-----------|------|
| K-means | Clustering |
| EM | Expectation Maximization |

Table 2 lists the algorithms used in our experiments. We have generated several datasets combining the different pre-processing techniques (pruning, stemming and wordnet). These first experiences shows that the application of pruning, stemming and wordnet reduces significantly the number of attributes without affecting the accuracy of results. The results obtained are in Table 3.

## 7  Conclusions and Future Work

At this point and as this is a Phd research work in progress, only a few considerations can be made. This paper focuses on the study of the impact of pre-processing techniques in clustering MEDLINE documents. The amount of information available on the MEDLINE database can (and probably will) be an issue. Although we are using a MEDLINE sample we have a dataset with 30000 MEDLINE citations. We have generated several datasets combining the different pre-processing techniques.

**Table 3.** Clustering Results

| Prunning | Stemming | Wordnet | Attributes | Acc | |
|---|---|---|---|---|---|
| | | | | **K-Means** | **EM** |
| > 10∗ | yes | yes | 4986 | 69.41 | 74.14 |
| > 100∗ | yes | yes | 985 | 68.48 | 71.26 |
| > 10∗ | yes | no | 5193 | 70.57 | 61.70 |
| > 100∗ | yes | no | 950 | 68.30 | 72.20 |
| > 10∗ | no | no | 6716 | 70.77 | 75,01 |
| > 100∗ | no | no | 1033 | 68.59 | 71.29 |

We will work further on the pre-processing stage and the algorithms to improve the results.

As a future work and to achieve a better clustering we may: 1) incorporate more information 2) optimise the MeSH terms selection for each document and 3) test with other MEDLINE samples.

# References

1. Wheeler, D.L., Barrett, T., Benson, D.A., Bryant, S.H., Canese, K., Chetvernin, V., Church, D.M., Dicuccio, M., Edgar, R., Federhen, S., Geer, L.Y., Helmberg, W., Kapustin, Y., Kenton, D.L., Khovayko, O., Lipman, D.J., Madden, T.L., Maglott, D.R., Ostell, J., Pruitt, K.D., Schuler, G.D., Schriml, L.M., Sequeira, E., Sherry, S.T., Sirotkin, K., Souvorov, A., Starchenko, G., Suzek, T.O., Tatusov, R., Tatusova, T.A., Wagner, L., Yaschenko, E.: Database resources of the national center for biotechnology information. Nucleic Acids Res **34**(Database issue) (January 2006)
2. Baeza-yates, R., Ribeiro-Neto, B.: Modern information retrieval (1999)
3. Zhou, W., Smalheiser, N.R., Yu, C.: A tutorial on information retrieval: basic terms and concepts. Journal of Biomedical Discovery and Collaboration **1** (March 2006)
4. Cohen, A.M., Hersh, W.R.: A survey of current work in biomedical text mining. Briefings in Bioinformatics **6**(1) (2005) pp. 57–71
5. Chapelle, O., et al.: Semi-supervised learning. Cambridge MIT Press (2006)
6. Zhu, X.: Semi-supervised learning literature survey (2006)
7. Zhu, S., Zeng, J., Mamitsuka, H.: Enhancing medline document clustering by incorporating mesh semantic similarity. Bioinformatics **25**(15) (2009) pp. 1944–1951
8. Lin, Y., Li, W., Chen, K., Liu, Y.: Text classification using machine learning. Journal of the American Medical Informatics Association : JAMIA **14**(5) (2007) pp. 651–661
9. Iliopoulos, I., Enright, A.J., Ouzounis, C.A.: Textquest: Document clustering of medline abstracts for concept discovery in molecular biology. In: Pacific Symposium on Biocomputing. (2001) pp. 384–395
10. Porter, M.F.: An algorithm for suffix stripping. (1997) pp. 313–316

11. Miller, G.A.: Wordnet: A lexical database for english. Communications of the ACM **38** (1995) pp. 39–41

12. Jain, A.K., Dubes, R.C.: Algorithms for clustering data. Prentice-Hall, Inc., Upper Saddle River, NJ, USA (1988)

13. Kaufman, L., Rousseeuw, P.: Finding Groups in Data: an Introduction to Cluster Analysis. John Wiley and Sons (1999)

14. Steinbach, M., Karypis, G., Kumar, V.: A comparison of document clustering techniques. In: In KDD Workshop on Text Mining. (2000)

# Spatio-temporal collaborative filtering for personalized photowalk recommendation

Sofia Torrão,

University of Porto, Faculty of Engineering, Rua Dr. Roberto Frias,
4200-465 Porto, Portugal
sofia.torrao@fe.up.pt

**Abstract.** In this paper is discussed the use of spatio-temporal collaborative filtering in an application that suggests photowalks to photographers, based on their location, their photostream and the data available from fellow photographers with similar space and time tags. We start with a brief introduction to the social media service for photography sharing Flickr and to the information that is available to the recommender system. Then, an introduction to collaborative filtering (CF), the relevance of spatio-temporal data and to some CF methods and algorithms is presented. Finally, a description of the framework for an application using this recommender system is provided, taking into consideration aspects relevant for the application usability, the relevance and personalization of the suggestions. The conclusions point to some specific requirements for the analysis and synthesis of spatio-temporal CF.

**Keywords:** collaborative filtering, recommender system, spatio-temporal information, personalization, photography, Flickr, social networks.

## 1 Introduction

As more and more information is made freely and easily available online in the Internet the need for tools to help on filtering and choosing what information is relevant and/or important has increased and improved [1]. For many years recommender systems have been developed and refined [2] to be applied in areas like e-commerce – case of Amazon that has one of the most famous item-to-item recommendation system [3] - and more recently within social networks [4] or even e-Learning [5], to help on this task. These recommender systems' main role is to help users to choose something or someone based on what their friends, colleagues, peers or someone like-minded have chosen in similar situations or to make predictions based on what this particular group has chosen for other related items [6]. The key aspects in these systems are relevance, similarity and influence that combined can generate more accurate results and make more reliable predictions.

The purpose of this paper is to discuss the use of a recommender system based on spatio-temporal collaborative filtering and the photo-sharing social media service Flickr [7], to suggest photowalks and interesting spots to photographers based on factors like themselves, their current location, their mobility, their photostream, their tag cloud and interesting photographs, fellow photographers or photowalks in the

same area. The application works in the user profile (both local as online information from the user Flickr's profile), fellow photographers who have taken photos in the area the user's in (radius or area to be configured by user), photos taken in the area, spaced with a determined interval (area, terrain, walking distance and mobility are some factors that represent constraints), and tags of the images found and their similarity to the user tags can help refine the filtering of results. For example, a user that has a majority of black and white photos in his photostream is more likely to follow the advices of a fellow photographer with same preferences for black and white photos.

## 1.1 Motivation

In the web 2.0 flourishing of online applications we've recently seen an increase of applications that connect to different social networks using APIs such as Flickr API, Facebook Connect, YouTube, Twitter API or Google OpenSocial (see Fig. 1).



**Fig 1.** Top APIs for Mashups. Programmable Web.

Many of these applications are born for mobile devices and users now not only can take photographs or finds musics with their mobile phones but they can immediately post them on their blog, share them on Flickr or Facebook or just send them with a tweet to show where they are and some applications can even feed several of these networks. The simplest way to upload is the e-mail and several applications generate a private e-mail address for each user to upload items (see Flickr example on Fig. 2).



**Fig 2.** Uploading by email to Flickr.

Many are also connected to leisure activities where social engagement and contacts through the network are maintained with much more frequency due to availability and

mobility of the devices. On another hand we can see also an increase in the do-it-yourself over the Internet like buying or selling, making traveling arrangements like booking flights or hotel, seeing the news, asking for information, etc.

A photowalk is a path usually walked by the photographer with a camera and thought as a sequence of good and interesting photographic places to shoot. It can be a path made of a collection of monuments, traditional streets, touristic attractions, or anything the photographer wants and is interested in shooting. A photowalk estimates a limited duration, both total as in between photo locations, bringing time vs space restrictions to the filtering process.

The purpose of this application is to suggest photowalks to users taking in consideration the user's characteristics, the place where the user is, the photowalk duration and data provided by Flickr. The application searches for photowalks, as a set of sequential paths, determined by the photos available in the area, within the time constrains and ranks each one according to the user preferences. In any point of the photowalk, the user should be able to reset the photowalk and ask for new suggestions or just finish the photowalk and return.

This approach implies not just the common item-item or user-item CF but also needs information about spatio-temporal proximity relating both the photos with the time they were taken. Recent related research in the spatio-temporal CF topic [1][8] don't use the time factor the same way our application needs.

In the next section we will briefly present Flickr, the social media service used for photos exchange, we will proceed in section 3 with the brief analysis of recommender systems and collaborative filtering techniques and in section 4 the design of the application framework is proposed. Finally the conclusions about requirements for the analysis and synthesis of spatio-temporal CF are presented in section 5.

## 2 About Flickr

Flickr [7] is a social media service dedicated to photo-sharing where basically any member can upload photographs that show in a personal photostream, create sets and collections with those photos, join groups, add contacts, comment and made annotations to other members photos. Photos have a unique id, can be public or private and are presented to Flickr users in its page with related information (see figure 3). The interaction within the network is based on relations between members: add as contact, add as family or add as friend or block someone, between members and groups of members: there are thematic groups related to lots of different subjects, cameras, photographic techniques, places, etc., and between members and content (photos, videos, illustrations) by means of adding comments, adding notes, adding tags or adding to personal favorites or galleries.

**Fig 3.** Example of a Flickr's photo. For each photograph the user can add information like: title, description, tags, location and include the photo on sets and/or Flickr's groups. If the photo is public and accepts comments, then other users can comment (you can see one comment to this photo) on the photo or add it as a favorite (this photo has been added as a favorite by 8 people). The photograph information also includes the dates in which it was taken and when it was upload to flickr.

Comments, views and favorites support a rewarding system for members and a rank of interestingness for photographs. There is a rank of the most interesting 500 photos for each day, called Explore (Fig. 4), built with an algorithm that takes in consideration factors like number of times a photo has been marked as favorite by someone, who has marked the photo as favorite and his/her rank, the period and frequency of the collection of comments and favorites, the number of pools (groups) the photos is in, the ratio of views to favorites, and other factors Flickr doesn't disclose.



**Fig 4.** Flickr's Explore photos for January 2009.

For each day there is a top 500 more interesting photos that go into Explore. These photos can enter in a rank and go up or down or even get dropped, since the rank is constructed and actualized during the day.

## 2.1 Relevant information for the recommender system

Since we're trying to build a recommender system that suggests photowalks and places of interest based on our user/member, a place and a duration, it is most important to know or infer what are the personal preferences regarding those subjects. For instance it's important to know some user preferences like photographing people or landscapes, nature, buildings or street photography, modern or traditional architecture, the age, sex, profession and other interests besides photography, the place he's in and what will be the duration of the photowalk besides other characteristics that can help the suggestions to be more accurate. Some of these factors can be asked by the application or just inferred and then ask the user to confirm. That's the case with the location that the application can get through GPS coordinates and just ask the user to confirm with more accuracy. Duration can be a parameter with a default value that can be altered by the user and these leave the application with the task to retrieve information about the user (to complement with the local information) and the photos that make the photowalk suggestion.

Flickr members have a profile where they can fill in personal information and the photograph's information can be added by means of title, description, location, tags, camera, exif information, groups the user is member of, the contacts and their personal information. Most of this information is available through the Flickr API [9] and can be accessed by our application (like the time stamp when the photo was taken). Examples of methods used to retrieve information can be seen on Figure 5 and an example of a response to a flickr.photos.getWithGeoData method can be seen on Figure 6.

**Fig 5.** Flickr's API methods. The figure shows methods for photos, comments, geo location, places, preferences and tags. There are several methods for retrieving information given a specific user, a specific location or a specific photo and these methods can be combined to help us retrieve the information needed for the filters.

```
<photos page="2" pages="89" perpage="10" total="881">
        <photo id="2636" owner="47058503995@N01"
                secret="a123456" server="2" title="test_04"
                ispublic="1" isfriend="0" isfamily="0" />
        <photo id="2635" owner="47058503995@N01"
                secret="b123456" server="2" title="test_03"
                ispublic="0" isfriend="1" isfamily="1" />
        <photo id="2633" owner="47058503995@N01"
                secret="c123456" server="2" title="test_01"
                ispublic="1" isfriend="0" isfamily="0" />
        <photo id="2610" owner="12037949754@N01"
                secret="d123456" server="2" title="00_tall"
                ispublic="1" isfriend="0" isfamily="0" />
</photos>
```

**Fig 6.** flickr.photos.getWithGeoData method. Example response.

Since Flickr API is public lots of different applications [10][11] can be found using this API. That's the case of Flickriver [12] - an new way to explore and view Flickr photographs, Flickr Hive Mind [13] - is a data mining tool for the Flickr photography database, allowing search by: tags (keywords); Flickr photography groups; Flickr users, their contacts, and favorites; free text; the Flickr Explore algorithm for

interestingness, Retrievr [14] - an experimental service which lets you search and explore in a selection of Flickr images by drawing a rough sketch, Multicolr Search Lab [15] - where you can browse through 10 million of Flickr's most 'interesting' Creative Commons images, and find ones that share the same colours or Bighugelabs that provides several web applications like Scout or Flickr DNA [16].

## 3   Collaborative filtering

It is most usual when looking for information to search it ourselves or ask someone we think knows where to find it. This is true in real life talking about human relations, someone that can be trusted or is recommended by someone we trust, or about things we are looking for or we found and need some help on how to handle them [6][17]. Recommendations play a major role on the act of choosing, more even if they come from someone that is like-minded or has some degree of similarity with us.

Collaborative filtering (CF) is one of the most successful approaches to building recommender systems [6][18][19]. These systems are built to make recommendations not only based on what they know about the users but also predict and infer likes from the data they know. There are lots of applications for recommender systems like photo recommendations, friend recommendations, forum recommendations, books recommendations, movies recommendations, ads matching to user profiles, and many more. These recommendations are more efficient if they can be both personalized as well as contextualized and are based on relevance of persons/friends, profiles and activities [2]. Some persons have more influence and their opinion on another person or an item can trigger a decision when the recommendation appears.

Social networks connections can provide immediately links between users and also the kind of these links. Through examining data one can obtain evidence of similar likes and dislikes, frequency of contacts or interaction, similar choices, mutual friends, and more user-to-user relations. Based on the users activities we can also obtain data about user-to-items relations and their preferences, likes and dislikes. The user-user CF is particularly relevant within social networks since these technologies are based on and provide the tools for user recommendations [4]. We can find relations like those presented [in table 1]:

**Table 1.**  Example of types of relations on Flickr [7].

| a->b | Sofia | Paulo | Mafalda | Ricardo |
|------|-------|-------|---------|---------|
| Sofia | | Friend | Family | Friend |
| Paulo | Contact | | Friend | Contact |
| Mafalda | Family | Friend | | Contact |
| Ricardo | Contact | Friend | Contact | |

Because the amount of data is usually very large, considerations related to scalability and computational aspects arise and the CF algorithms have evolved to better address these questions.

Collaborative filtering systems can be user-based, like suggesting friends based on the user profile – it's a user to user aggregation type of CF – or item-based when the CF aggregates items – for example 2 items sold at the same time can be offered as a bundle for another user browsing one of those books [1]. In the user-based CF scenario first operation is to aggregate users for a chosen item based on users' similarity (profile based) whereas the item-based CF first aggregates similar items based on their ranking from different users

**Table 2.** Example of user-item relations.

| a->b | Photo 1 | Photo 2 | Photo 3 | Photo 4 |
|---|---|---|---|---|
| Sofia | | like | fav | |
| Paulo | | | like | like |
| Mafalda | fav | fav | | like |
| Ricardo | like | fav | like | |

A basic function of a recommender system is to suggest an item where there is no information (some empty cells in the table above) based on similarity, previous knowledge given by user profile and other similar items rating by the same user and others.

Collaborative filtering algorithms are usually divided into three categories: memory-based CF, model-based CF and hybrid recommenders [see figure 5].

| CF categories | Representative techniques | Main advantages | Main shortcomings |
|---|---|---|---|
| Memory-based CF | *Neighbor-based CF (item-based/user-based CF algorithms with Pearson/vector cosine correlation) <br> *Item-based/user-based top-N recommendations | *easy implementation <br> *new data can be added easily and incrementally <br> *need not consider the content of the items being recommended <br> *scale well with co-rated items | *are dependent on human ratings <br> *performance decrease when data are sparse <br> *cannot recommend for new users and items <br> *have limited scalability for large datasets |
| Model-based CF | *Bayesian belief nets CF <br> *clustering CF <br> *MDP-based CF <br> *latent semantic CF <br> *sparse factor analysis <br> *CF using dimensionality reduction techniques, for example, SVD, PCA | *better address the sparsity, scalability and other problems <br> *improve prediction performance <br> *give an intuitive rationale for recommendations | *expensive model-building <br> *have trade-off between prediction performance and scalability <br> *lose useful information for dimensionality reduction techniques |
| Hybrid recommenders | *content-based CF recommender, for example, Fab <br> *content-boosted CF <br> *hybrid CF combining memory-based and model-based CF algorithms, for example, Personality Diagnosis | *overcome limitations of CF and content-based or other recommenders <br> *improve prediction performance <br> *overcome CF problems such as sparsity and gray sheep | *have increased complexity and expense for implementation <br> *need external information that usually not available |

**Fig 5.** Overview of collaborative filtering techniques [1].

## 4 Application

The application should be simple and easy to use, can run in a laptop or a mobile device such as a 3G mobile phone or wi-fi smart-phone and to be available anytime and anywhere, and the user should enter is username and password of his Flickr account. After authentication and authorization, the application will construct a profile for the user and retrieve the necessary information both from Flickr as well as from the context, asking for the user intervention when necessary. For instance, some mobile devices have GPS location that can be used as automatic input to the recommender system but others don't. In this last scenario the user will be asked to enter the name of the place where he's at or eventually to select it from a list or a map. The application can also ask for a more complete profile that will be kept local and/or can override the settings on Flickr.

With this information the application will then ask for some more optional parameters like how much time the user wants to spend – photowalk duration, the "hardness" of photowalk (plain, steep, with stairs, climbing, etc.), the weather, and other possible restrictions the user has at that time. Some options can also be configured by this time, like selecting the type of photos or photowalks to search. This operation can be presented as a tag choosing where the user is presented with his common tags and where he selects those relevant to him. The application would then use that information to input in the CF algorithm to retrieve all the information/photowalks.

Finally the user will be presented with a few choices, ranked and preferably one major and two-second choices that reflect not only his preferences regarding photowalks, within the available time, but also the place he's in and the local conditions or restrictions.

## 5 Final considerations

With this paper we propose an application to study and understand the spatio-temporal collaborative filtering in detail and how it can be used to improve or extend recommender systems. We expect our research and expertise in spatio-temporal collaborative filtering to contribute to a new generation of recommendation technologies.

## References

1. Spindler, A., Norrie, M. C., Grossniklaus, M., Signer, B.: Spatio-Temporal Proximity as a Basis for Collaborative Filtering in Mobile Environments. Ubiquitous Mobile Information and Collaboration Systems, 912-926 (2006)
2. Adomavicius, G., Tuzhilin, A.: Toward the Next Generation of Recommender Systems: A Survey of the State-of-the-Art and Possible Extensions. IEEE Transactions on knowledge and data engineering, vol. 17, n.6 (2005)

3. Linden, G., Smith, B., York, J.: Amazon.com recommendations: item-to-item collaborative filtering. IEEE Internet Computing 7 (1) 76–80 (2003)
4. Zheng, R., Wilkinson, D., Provost, F.: Social network collaborative filtering. Working Paper CeDER-8-08, Center for Digital Economy Research, Stern School of Business, NYU. (2008)
5. Bobadilla, J., Serradilla, F., Hernando, A., MovieLens: Collaborative filtering adapted to recommender systems of e-learning. Knowledge Based Systems, doi:10.1016/j.knosys.2009.01.008 (2009)
6. Su, X., Khoshgoftaar, T. M.: A Survey of Collaborative Filtering Techniques. Advances in Artificial Intelligence, Volume 2009, Article ID 421425, doi:10.1155/2009/421425 (2009)
7. Flickr from Yahoo, http://www.flickr.com/about/
8. Lu, Z., Agarwal, D., Dhillon, I. S.: A spatio-temporal approach to collaborative filtering. In Proceedings of the Third ACM Conference on Recommender Systems RecSys '09. ACM. doi:10.1145/1639714.1639719 (2009)
9. Flickr API, http://www.flickr.com/services/api/
10. 60+ Tools To Enhance Your Flickr Experience, http://www.hongkiat.com/blog/60-tools-to-enhance-your-flickr-experience
11. Flickr API Mashups. http://www.programmableweb.com/api/flickr/mashups
12. Flickriver, http://www.flickriver.com/
13. Flickr Hive Mind, http://fiveprime.org/flickr_hvmnd.cgi
14. Retrievr, http://labs.systemone.at/retrievr/
15. Multicolr, http://labs.ideeinc.com/multicolr
16. Bighugelabs, http://bighugelabs.com
17. Liu, F., Lee, H. J.: Use of social network information to enhance collaborative filtering performance. Expert Systems with Applications, doi:10.1016/j.eswa.2009.12.061 (2009)
18. Chen, A. Y., McLeod, D.: Collaborative Filtering for Information Recommendation Systems. Encyclopedia of Data Warehousing and Mining (2005)
19. Kim, H.-N., et al.: Collaborative filtering based on collaborative tagging for enhancing the quality of recommendation. Electronic Commerce Research and Applications, doi:10.1016/j.elerap.2009.08.004 (2009)

# Image Analysis

# Binary Images Clustering with K-means

João Ferreira Nunes,

Escola Superior de Tecnologia e Gestão do Instituto Politécnico de Viana do Castelo
Avenida do Atlântico, s/n 4900-348 Viana do Castelo, Portugal
joao.nunes@estg.ipvc.pt

Faculdade de Engenharia da Universidade do Porto
Rua Dr. Roberto Frias, s/n 4200-465, Porto, Portugal
pro09001@fe.up.pt

**Abstract.** Nowadays, with the large amount of image databases available, there is the need to access them in the most various and easiest ways. For many applications, an expeditious manner to fulfill such need encompasses partitioning images by their content both for indexing or retrieval purposes. In this paper we propose a method to group binary images in respect to their content by means of an unsupervised learning technique, k-means clustering. The paper describes image pre-processing and feature extraction. A clustering algorithm is then applied over the extracted feature vectors. A set of clustering quality criteria is used to assist the selection of the best number of clusters. To have a better understanding of the effectiveness of the method, the obtained clusters were evaluated against an a priori available partition of the dataset. Achieved results are encouraging and demonstrate the ability and effectiveness of the proposed approach.

**Keywords:** image clustering, cluster analysis, data mining, k-means.

## 1    Introduction

With the fast growth of multimedia data, mainly due to the wide spread of digital devices, multimedia repositories became very common, and some of them extremely large. It was natural that with this amount of archived information it would come out the need of indexing and retrieving this unstructured data. There are available some tools for managing and searching within these collections, however we also need tools to extract the hidden useful knowledge embedded on them. For example, tools for discovering relationships between objects, classifying images based on their content and tools for extracting data patterns. This paper presents some experiments on clustering binary images. We intend to demonstrate the effectiveness of the k-means clustering algorithm for grouping a set of binary silhouette images into the *best* number of sets, based on some features extracted from those images. Finally we also intend to evaluate the results of this method using internal and external measures.

The experiments are supported on a dataset that gives some knowledge *a priori*, where every image is labeled as belonging to a given set. Because these images are in a binary format and are silhouettes of objects, it is expected that after the clustering process, images from different original sets should be grouped together, and

consequently the number of resulting sets it is also expected to be lower than the initial value (e.g., silhouette images of guitar objects can be grouped together with silhouette images of spoons).
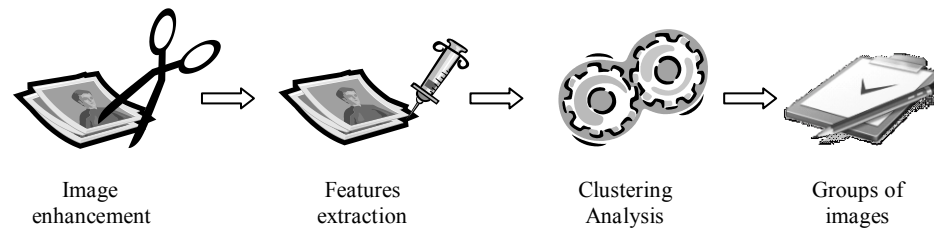
The paper is organized in four sections: the first one, which is the Introduction, is where we present the state of existing knowledge in the area of Clustering algorithms, Clustering validation and also Image Clustering. We also present the purpose of the paper, and a prediction of the expected results. The section 2 is where we can find the descriptive part, which includes the sections of Data Collection, Preprocessing and Cluster Analysis. The final results of ours experiments are reported in the third section, and finally, conclusions and future work appear in the last section.

## 2 Data Collection, Preprocessing and Cluster Analysis

In order to classify silhouette binary images, we have chosen the dataset "MPEG-7 Core Experiment CE-Shape-1 Test Set". The main reason to use this image collection was the fact that it is a public dataset, and because of that, it has been used in some other studies [1] [2] [3] [4]. Another aspect that made us to decide on this dataset was the previous knowledge on the images' labeling classification. From the beginning we knew for each image, what was its label, and that information could be useful to validate our clustering process.

With this dataset we have experimented grouping *similar* images into clusters through an unsupervised data mining technique of clustering [5], implemented by the k-means algorithm. However, before this step could be accomplished, previously we had to apply some image enhancements operations, so that we could work with *cleaner* images. In this case we have cut the images removing irrelevant information and we have also filtered morphologically them.

The workflow of our model includes the preprocessing phase, followed by the extraction of images features and finally, the clustering analysis with the k-means algorithm. The Figure 1 illustrates an overview of our clustering process.
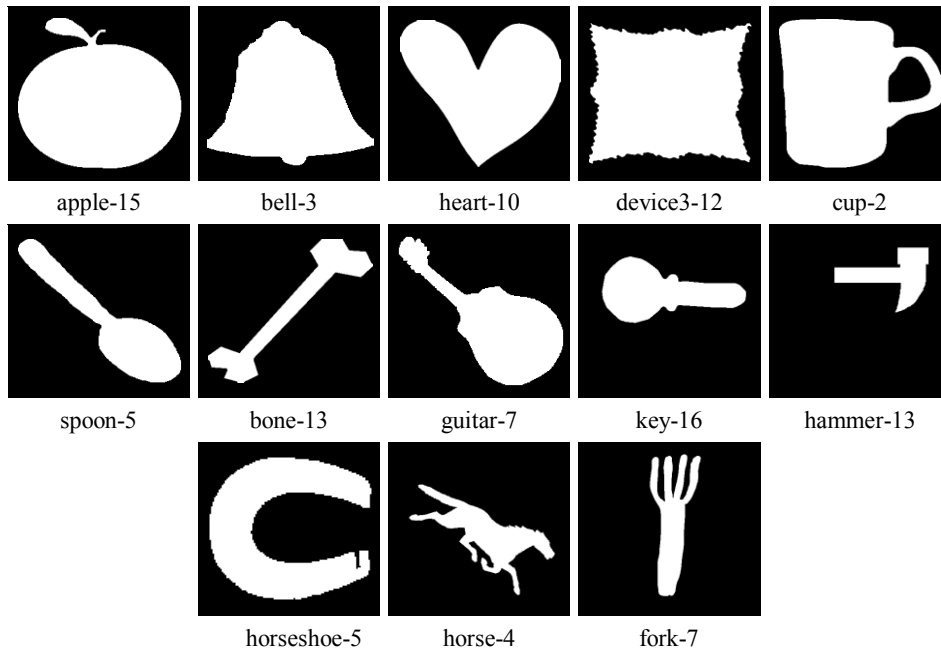


| Image enhancement | Features extraction | Clustering Analysis | Groups of images |

**Figure 1.** Image Clustering process.

### 2.1 MPEG-7 Images Dataset

The dataset that was used in our experiments is called "MPEG-7 Core Experiment CE-Shape-1 Test Set" and was created by the MPEG-7 committee. The Motion Picture Expert Group (MPEG) [6] is a working group of ISO/IEC[1] and has defined the MPGE-7 standard for description and search of audio and visual content. This dataset contains a large collection of binary images, which by its definition only allocates two possible values for each pixel. In this case, the two values allocated that correspond to the image colors are the black (zero) for the background and white (one) for the foreground.

Since these silhouette images represent 2D objects that are projections of 3D objects, their silhouettes may change due to: (1) change of a view point with respect to objects; (2) non-rigid object motion (e.g., people walking or fish swimming); and (3) noise (e.g., digitization and segmentation noise). Also, few additional characteristics of the dataset to be mentioned are that some images have holes in them, while others do not, and some images have experienced a number of transformations, such as scales, cuts and rotations and, at last, the image resolution is not constant among them. Figure 2 illustrates some sample images of the MPEG-7 dataset.

| apple-15 | bell-3 | heart-10 | device3-12 | cup-2 |
| spoon-5 | bone-13 | guitar-7 | key-16 | hammer-13 |
| horseshoe-5 | horse-4 | fork-7 | | |

**Figure 2.** Sample images included in "MPEG-7 Core Experiment CE-Shape-1 Test Set". The image files with the same name prefix are classified as belonging to the same set. In this case, each image represents a different set.

---

[1] International Organization for Standardization/International Electrotechnical Commission

The "MPEG-7 Core Experiment CE-Shape-1 Test Set" is accessible from various sources in the World Wide Web [7] [8], since it has been used on other researches, and as a result some of their authors also publish it. The dataset includes 1400 images grouped into 70 sets, and each set contains 20 samples. In our experiments, we wanted to work with a smaller dataset and therefore, using a visual and subjective criterion, we have chosen a MPEG-7 subset filled with the following 13 classes: bone, guitar, horse, horseshoe, heart, apple, bell, device3, cup, fork, hammer, key and spoon. This new set is composed of visually similar classes between them (guitar, spoon, key), as well as clearly distinct classes (horse, horseshoe).

### 2.2 Preprocessing and Features Extraction Phases

Preprocessing is always a necessity whenever the data to be mined is noisy, inconsistent or incomplete and preprocessing significantly improves the effectiveness of the data mining techniques [5]. This section of the paper introduces the preprocessing techniques that we have applied to the images before the feature extraction process. We intended to reduce the images' noise by removing their irrelevant information. This was accomplished by detecting and extracting the images' region of interest, cropping the images through their bounding box. Another preprocessing technique that we have also applied is the close morphological filter. This filter closes morphologically the binary image and it is defined as the dilation of the image followed by the erosion of the dilated image. The closing filter operation smooth boundaries, reduce small inward bumps, join narrow breaks and fill small holes caused by noise.

In order to compute these two preprocessing techniques, a procedure in MATLAB was developed and it is listed bellow:

```
for i = 3 : nfiles
    %opens the image file
    img = imread(filename);
    %applys the morphological filter to the binary image
    se = strel('disk',4);
    img = imclose(img, se);
    %REGIONPROPS expects a label matrix.
    LabeledBWImg = bwlabel(img);
    %extracts img features available from REGIONPROPS.
    stats = regionprops(LabeledBWImg);
    %crops the image through its bounding box.
    img = imcrop(img, stats.BoundingBox);
    %saves the resulting new image
    imwrite(img, filename, 'png', 'bitdepth', 1);
end
```

After accomplished the first step on the Image Clustering Process, we started to conclude which image features we intended to extract. Our main concern at this stage was to assure that all the features' values should be normalized between zero and one, so that at the time of clustering the images, all the features would be weight balanced.

Therefore we developed another MATLAB procedure that computes, for each image in the dataset, the required features based on the images attributes. Those attributes were acquired invoking the *regionprops* function from the MATLAB Image Processing toolbox and the *momentsupto3* function from the Lifting Scheme on Quincunx Grids (LISQ) toolbox [9]. The first function measures a set of properties of the image, such as its area, Euler number, bounding box, perimeter, centroid, etc., while the second one computes the images moment invariants. The resulting features were stored on an $n$ x $f$ matrix, where $n$ is the number of images and $f$ is the number of features and this matrix was used as input during the clustering phase. In our experiments we decided to extract seven features, which are:

F1: *Solidity*. This feature results from the ratio between the image Area and its ConvexHullArea, where Area is the number of pixels in the foreground region and the ConvexHullArea is the number of pixels of the area of the smallest convex polygon that can contain the same region.

F2: *Axis Ratio*. It's the ratio between the MinorAxisLength and the MajorAxisLength. The MinorAxisLength gives the length (in pixels) of the minor axis of the ellipse that has the same normalized second central moments as the region, while the MajorAxisLength gives the length (also in pixels) of the major axis of the ellipse that has the same normalized second central moments as the region.

F3: *Areas Ratio;* It's the ratio between the image Area and the image FilledArea. The attribute Area gives the number of pixels in its foreground region while the FilledArea gives the number of on pixels in FilledImage. This ratio feature gives a notion if the image has holes on it or not, where values close to zero indicate that the image has very few holes.

F4: *Perimeter-Area Ratio*; It's the ratio between the image Perimeter and the image Area.

F5: *Eccentricity*; Specifies the eccentricity of the ellipse that has the same second-moments as the region. The eccentricity is the ratio of the distance between the foci of the ellipse and its major axis length. An ellipse whose eccentricity is zero is actually a circle, while an ellipse whose eccentricity is one is a line segment.

F6: *Extent*; Specifies the ratio of pixels in the foreground region with the pixels in the total bounding box. It is computed as the Area divided by the Bounding Box area.

F7: *Invariant moment*; This is a useful measure to describe objects because image properties that are found via image moments are invariant under translation, changes in scale, and also rotation. From the Hu [10] set of invariant moments, we've chosen the *skew invariant* and to compute this feature we used the *momentsupto3* function from the LISQ toolbox. This function returns all the Hu seven moments, however, we chose the one that statistically seemed to establish more differences between classes, and consequently highlight their dissimilarities.


## 2.3    Cluster Analysis

Cluster analysis is traditionally considered as an unsupervised method that has been widely used in numerous applications, including market research, pattern recognition, data analysis, and image processing. The main potential of clustering is to detect the underlying structure in data, not only for classification and pattern recognition, but for

model reduction and optimization. It consists in the process of grouping a set of objects into classes of similar objects [5]. Thus a cluster is a collection of data objects that are similar to one another within the same cluster and are dissimilar to the objects in other clusters. Similarity and dissimilarity can be measured for two objects based on several features variables. In the context of our experiments, similarity is a quantity that reflects the strength of the relationship between two images, and dissimilarity measures the discrepancy between two images. Euclidean Distance (Eq.1) is the distance measure that we've used to measure dissimilarities between two images $i$ and $j$. It examines the root of square differences between all the attributes $f$ of the pair of objects (images):

$$d_{ij} = \sqrt{\sum_{f=1}^{n} \left( x_{if} - x_{jf} \right)^2} \; . \qquad \textbf{(1)}$$

In order to group *closest* images into the same sets through clustering, we've implemented the k-means algorithm, which will be explained in detail in the following section.

### K-means Algorithm

K-means is one of the simplest unsupervised learning algorithms that solves the clustering problem. It was developed by J. MacQueen [11] and then by J. A. Hartigan and M. A. Wong. This algorithm is used to group objects into $k$ number of classes based on a set of their attributes/features. It is sensitive to initialized partition. The main idea of how this algorithm works is the following: it starts by randomly picking $k$ objects defining them as the centroids of the clusters, and then, repeatedly does, for each object, place the object *inside* the cluster to whose centroid it is closest, calculating again the centroids for the cluster which has gain the object and also for the cluster which has lost the object. After that, the algorithm repeats this last step until there is no change in clusters' composition between two consecutive iterations. The most common and perhaps the *best* use of this algorithm requires the previous knowledge of the number of classes to split the initial set of objects, which is the $k$ number.

In our experiments we intended to find the *best k* number, in order to group the dataset images into an optimal number of classes based on their features' similarity. Defining the *best* number of classes has been an open problem in recent times with a considerable research activity and it can be validated using appropriate internal and/or external criteria and techniques [12].

### Evaluation of clustering

For the purpose of validating the clustering solutions and consequently assess the *best* number of clusters we have ran the k-means algorithm varying $k$ from 3 to 20 and then we used two different methods to evaluate the results. In one method we computed some internal criteria for every $k$ tested while in the other method we computed some external criteria. The internal criteria offer an idea of the solution cohesion (how closely related are objects in a cluster) and also of the solution separation (how well-separated a cluster is from other clusters) while the external

measures are related to how representative are the current clusters to the known classes. They compare a clustering result with a known set of class labels to evaluate the degree of consensus between the two. All the measures were calculated through the Cluster Validity Analysis Platform [13] in MATLAB.

It's important to notice that all these results obtained by these methods provide some guidelines and do not indicate an exclusive "correct" number of clusters. They are at the disposal of the *expert* in order to evaluate the resulting clustering.

Thus, the first method that was used to validate the clustering solutions was based on the fact that it had no knowledge a priori, and therefore it could only be computed internal criteria. Among the criteria available, we chose to compute the Silhouette index [14] where the largest silhouette value indicates the optimal $k$, the Calinski-Harabasz index where also the maximum value indicates the optimal $k$, the C index [15] where the minimal C-index indicates the optimal $k$ and finally the weighted inter-intra index that searches forward ($k=2,3,4,...$) and stop at the first down-tick of the index, which indicates the optimal $k$.

On the second method, it was used the available image labeling information in order to evaluate the clustering solutions, becoming possible to compute external criteria. In this case there were computed the Jaccard index [16] where the maximum value indicates the optimal $k$, the Fowlkes-Mallows index [17] where also the maximum value indicates the optimal $k$ and finally the Adjusted Rand index where also the maximum value indicates the optimal $k$.

The results of these methods are presented in the next section on this paper.

## 3    Results

In our experiments, we considered 13 classes of images from the MPEG-7 dataset and then we wanted to group them according to their seven previously extracted features, using the k-means clustering algorithm. So, the k-means algorithm was applied for $k$ varying from 3 to 20, and then two methods were used to validate the clustering results. The first one computed some internal criteria, which are graphically represented in Figure 3 and the second one computed some external criteria that are in Figure 5.

As we can observe on Figure 3, only the C index suggests the number five as the optimal number of clusters, while the others indexes indicate the number six, which was actually the $k$ number we considered. The following Table 1 expresses the images distribution among the six clusters based on the initial label information and Figure 4 illustrates three sample images that were grouped on the same cluster.
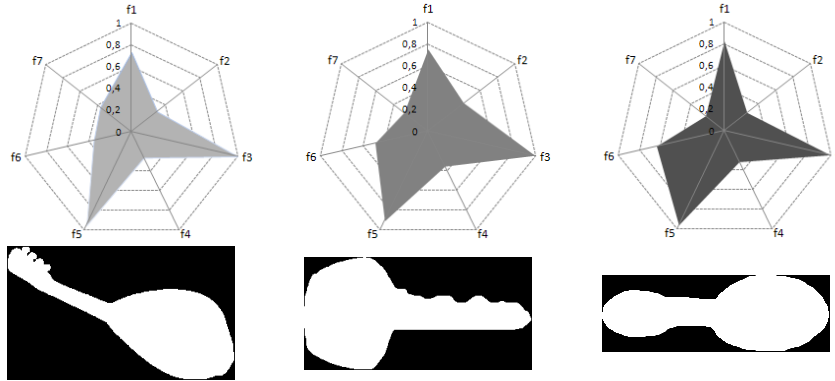
**Figure 3.** Graphical representation of the Silhouette index, the Calinski-Harabasz index, the C index and the weighted inter-intra index.

**Table 1.** Distribution of the images through the six clusters versus the a priori partition (labels). Computed clusters are represented in columns and initial labels in rows.
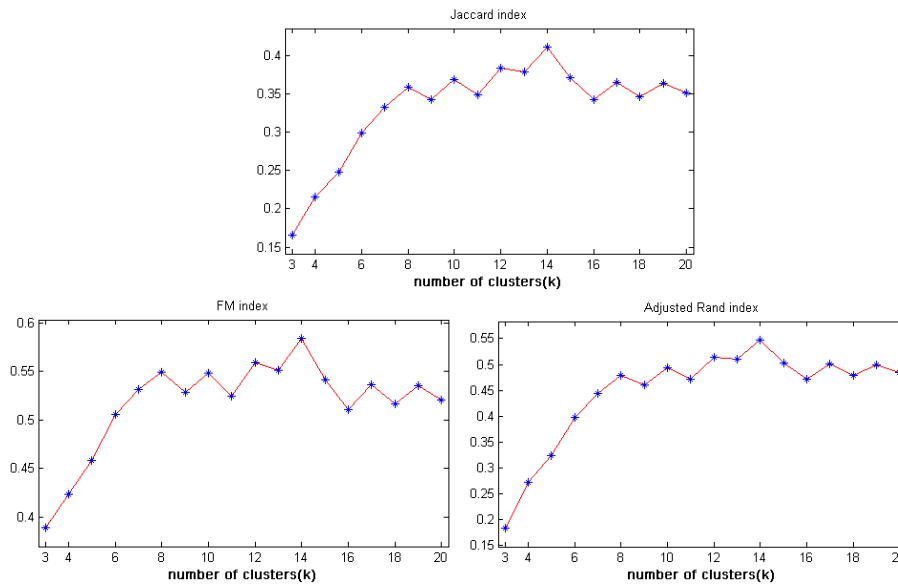
|  | C1 | C2 | C3 | C4 | C5 | C6 |  |
|---|---|---|---|---|---|---|---|
| apple |  | 20 |  |  |  |  | 20 |
| bell |  |  |  |  |  | 20 | 20 |
| bone |  |  | 18 | 2 |  |  | 20 |
| cup |  | 12 |  |  |  | 8 | 20 |
| device3 |  | 15 |  |  | 5 |  | 20 |
| fork |  |  | 15 | 5 |  |  | 20 |
| guitar |  |  | 1 | 19 |  |  | 20 |
| hammer | 8 |  | 12 |  |  |  | 20 |
| heart |  | 1 |  |  |  | 19 | 20 |
| horse | 20 |  |  |  |  |  | 20 |
| horseshoe |  |  |  |  | 20 |  | 20 |
| key |  |  |  | 19 |  | 1 | 20 |
| spoon |  |  | 8 | 12 |  |  | 20 |
|  | 28 | 48 | 54 | 57 | 25 | 48 |  |

**Figure 4.** This figure shows three images that were grouped in the same cluster (cluster 4). Each image has on its top a radar plot graph of their features, represented with f1, f2,...,f7. With this type of representation it is noticeable how the image features shapes are quite similar and therefore grouped together.

The next Figure 5 illustrates the graphical representation of the external criteria computed within the second method. In this case the suggested value for the *best k* was the fourteen.



**Figure 5.** Graphical representation of the Jaccard index, the Fowlkes-Mallows index and the Adjusted Rand index.

To better understanding how images were clustered in respect to their initial labels, Table 2 illustrates the distribution among the fourteen clusters versus initial label information. As it can be observed, there are some clusters that perfectly match the initial label, as for instance the horseshoe and horse. Some other images are well concentrated in only one cluster, such as the key, bone, heart and apple. Some other images are grouped with images having different labels. An example is the spoon set of images that results distributed among four different clusters (C5, C9, C10 and C11). An observation of the images grouped within these clusters suggests that the spoon images are clustered with visually similar images (e.g. guitar and key).

**Table 2.** Distribution of the images through the fourteen clusters versus the a priori partition (labels). Computed clusters are represented in columns and initial labels in rows.

| | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 | C10 | C11 | C12 | C13 | C14 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| apple | | | 15 | | | | | | | | | | 5 | | 20 |
| bell | | | | 20 | | | | | | | | | | | 20 |
| bone | | | | | | | | | 1 | | 2 | 17 | | | 20 |
| cup | | | 7 | 2 | | | | | | | | | 11 | | 20 |
| device3 | | | 2 | | | | 5 | | | | | | | 13 | 20 |
| fork | | | | | 5 | | | 7 | 2 | 1 | 4 | 1 | | | 20 |
| guitar | 1 | | | | | | | | 14 | 4 | 1 | | | | 20 |
| hammer | 8 | | | | 4 | | | 4 | | | 4 | | | | 20 |
| heart | | | | 19 | | | | | | | | | 1 | | 20 |
| horse | | 20 | | | | | | | | | | | | | 20 |
| horseshoe | | | | | | 20 | | | | | | | | | 20 |
| key | | | | 1 | | | | | 1 | 18 | | | | | 20 |
| spoon | | | | | 4 | | | | 4 | 4 | 8 | | | | 20 |
| | 9 | 20 | 24 | 42 | 13 | 20 | 5 | 11 | 22 | 27 | 19 | 18 | 17 | 13 | |

The used set of evaluation criteria, internal and external, does not suggest the same number of clusters. As the ultimate goal of this work is not to achieve a solution that recovers the truth (in form of the initial labeling), we consider from the conducted results analysis that the set of extracted features and selected algorithm is able to cluster images, in an unsupervised manner, in respect to their visual content.

We notice, that if some hint (a priori knowledge) about the initial number of classes would be known, as for instance in the form of a narrower range of expected clusters (e.g. from 10 to 16), the selected cluster would better match the initial partition, as the internal criteria would suggest a similar number of clusters. It is the case for the silhouette index, illustrated in Figure 3, which for that range suggests the 14 as the best number of clusters.

## 4    Conclusions and Future Work

In this paper, we presented a process to enable grouping of binary images by means of a learning technique over a set of extracted features. We assume that the process can be conducted without or with only few a priori information. The chosen

learning technique was the k-means clustering algorithm that needs the input of number of clusters (k). As, in our experience design, the number of clusters is not available or unrevealed on purpose, several clustering iterations were conducted for a wide range of the k value. Internal criteria were then used to choose the best number of clusters among all experiments. Results were then compared with the a priori partition in order to obtain an objective assessment of the ability of the unsupervised clustering match the represented object labeling.

Given the dataset, selected features and clustering algorithm, internal criteria suggests a best choice would occur for k = 6 clusters. The comparison with the (meanwhile revealed) initial partition, suggests a higher number of clusters, namely k = 14. This may indicate that the selected features are not enough or the more appropriate to obtain a clustering that closely matches the a priori partition. However, as illustrated, the images are closely grouped in respect to their visual silhouette. We notice that this is not an undesired result as the ultimate goal is not to match the given initial partition. The goal is to group images in a way that they can be indexed by their content (namely their visual properties). The fewer number of clusters suggest a reduction of the number of original concepts (initial labels) which is a feature/goal of this kind of clustering algorithm.

Achieved results are encouraging and suggest adequacy of the selected features and algorithm in order to group images by their visual content. Despite the former, we intend to explore new features and conduct an analysis to refine their selection. Instead of using the Euclidean distances to measure *similarity* we also consider the use of weighting attributes, according to their relevance.

We also intend to apply a supervised method (e.g.: k-nearest neighbor) to develop an information retrieval system. Using a query image, the system will extract the image features and classify it according to the sets previously defined with our clustering model. Then, it will return a set of images ranked according to similarity measures. Eventually this process could be evaluated by an *expert* and therefore it could provide more information/knowledge to the system. These are some of the future goals to improve our work presented.

# References

1. Bai, X., Yang, X., Latecki, L. J., Liu, W., Tu, Z.: Learning Context Sensitive Shape Similarity by Graph Transduction. IEEE Trans. Pattern Analysis and Machine Intelligence (PAMI), (2009)
2. Ling, H., Jacobs, D.: Shape classification using the inner-distance. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), 29, 2, p. 286-299 (2007)
3. Alajlan, N., Kamel, M., Freeman, G.: Geometry-based image retrieval in binary image databases. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), 30, 6, p. 1003--1013 (2008)
4. Super, B.: Retrieval from shape databases using chance probability functions and fixed correspondence. International Journal of Pattern Recognition and Artificial Intelligence, 20, 8, p. 1117-1137 (2006)
5. Jiawei, H.J., Kamber, M., Pei, J., Data Mining: Concepts and Techniques: Morgan Kaufmann Publishers (2006)

6.  The Moving Picture Experts Group (MPEG), http://www.chiariglione.org/mpeg/,2009.12.01

7.  MPEG-7 Core Experiment CE-Shape-1 Test Set, http://www.imageprocessingplace.com,2009.12.01

8.  MPEG-7 Core Experiment CE-Shape-1 Test Set, http://www.ehu.es/ccwintco/uploads/d/de/MPEG7_CE-Shape-1_Part_B.zip,2009.12.01

9.  Zeeuw, P.M.:A toolbox for the lifting scheme on quincunx grids (LISQ) (2002)

10. Hu, M.K.: Visual Pattern Recognition by Moment Invariants. IEEE Transactions on Information Theory, IT-8, 179-187 (1962)

11. MacQueen, J.B. Some Methods for classification and Analysis of Multivariate Observations. in 5-th Berkeley Symposium on Mathematical Statistics and Probability. Berkeley, University of California (1967)

12. Halkidi, M., Batistakis, Y., Vazirgiannis, M.: On Clustering Validation Techniques. Journal of Intelligent Information Systems, 17:2/3, 107--145 (2001)

13. Wang, K., Wang, B., Peng, L.: Validation for Cluster Analyses Data Science Journal, 8, 88--93 (2009)

14. Rousseeuw, P.J.: Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. Journal of Computational and Applied Mathematics, 20, 53--65 (1987)

15. Hubert, L., Schultz, J.: Quadratic assignment as a general data-analysis strategy. British Journal of Mathematical and Statistical Psychologie, 29, 190--241 (1976)

16. Jaccard, P.: The distribution of flora in the alpine zone. New Phytologist, 11, 37--50 (1912)

17. Fowlkes, E.B., Mallows, C.L.: A method for comparing two hierarchical clusterings. Journal of the American Statistical Association 78(383), 553--569 (1983)

# Exploiting eInfrastructures for Medical Image Storage and Analysis: A Grid Application for Mammography CAD

Raúl Ramos Pollán[1] , Miguel Ángel Guevara López[2]

[1] CETA-CIEMAT Centro Extremeño de Tecnologías Avanzadas,
Calle Sola 1, 10200 Trujillo, Spain,
raul.ramos@ciemat.es
[2] INEGI Instituto de Engenharia, Mecanica e Gestão Industrial, Universidade do Porto,
Campus da FEUP, Rua Roberto Frias 400, 4200-465 Porto, Portugal
mguevaral@inegi.up.pt

**Abstract.** This paper presents a Digital Repositories Infrastructure (DRI) software platform that simplifies the management of digitalized content and metadata stored on Grid, exploiting its features such as strong security contexts, data federation and large storage and computing capacities. The DRI proposed here includes repository browsing tools, a custom made viewer and a mammograms analysis graphical interface. In general, DRI offers an easy and intuitive interaction model with content stored on the Grid and it also constitutes the foundation to build specific domain applications. In this context, we developed an experimental Grid application in Mammography Computer-Aided Diagnosis (CAD) including a general framework that supports all its stages (image processing, segmentation, training, classification, etc) allowing semiautomatic classification of digital mammograms. This CAD achieved 90% true positive detection on the 322 images of the MIAS database.

**Keywords:** Biomedical Computing, Health Care Information Systems, Grid, CAD, Mammography

## 1  Introduction

Grid infrastructures are powerful tools providing large and federated storage and computing capacities for many user communities [1]. However, their usage is still cumbersome for many users due to a lack of easy-to-use tools to interact with them. The Digital Repositories Infrastructure (DRI) is a software platform created at CETA-CIEMAT aimed at reducing the cost to host digital repositories of arbitrary nature on Grid infrastructures, providing both users and repository providers a set of graphical and conceptual tools to easily define repositories and manage their content.

The digital repository here presented is composed of a set of units of digitalized content annotated with metadata [2] described through an entity-relationship model. With DRI, a repository provider describes his repository data model in an XML file and has immediately available a set of standard graphical user interfaces for  browsing

and managing repository content stored on a Grid infrastructure. On top of that he could also develop custom tools to provide specific functionality for his repository (for content viewing, data analysis, etc). This way, a repository of mammograms studies is composed of digital content (the mammograms images) and metadata (patient info, diagnoses, etc.). The aim of this work was to validate DRI usage to support two scenarios: (1) managing large collections of federated mammograms studies and (2) building CAD systems. The scope of previous works and pilot systems [3] was enlarged to a great extent to provide a full-featured production system to manage and analyse repositories of medical images through the exploitation of Grid infrastructures.

## 2  The DRI Platform

### 2.1 DRI Overview

Figure 1 shows the architecture of the DRI platform. A repository is defined by a *repository provider* in a *repository description file*, which is a XML representation of the data model of the repository. Fields defined in this file (such



**Fig. 1:** Simplified DRI architecture

as *patient name*, *age* or digitalized mammograms) are marked up as either *metadata* or as *large digital content* (see Section 3). The DRI engine parses the repository description file and creates the appropriate storage structures so that (1) metadata is stored on a regular local database and (2) large digital content is stored on a Grid infrastructure. DRI also provides web and stand-alone applications for browsing the repository and managing its content. Repository providers can also describe how users navigate the repository. For instance, a repository provider might want users to navigate a medical image repository by choosing between browsing studies classified by diagnosed pathology or browsing patients alphabetically. For this, the repository provider includes in the *repository description file* a description of *navigation trees* which are needed for the users. Note that a repository provider needs only to define a repository description file and, with this, he can start using the underlying Grid storage and database to host his repository and offer to repository users graphical tools to interact with the content.
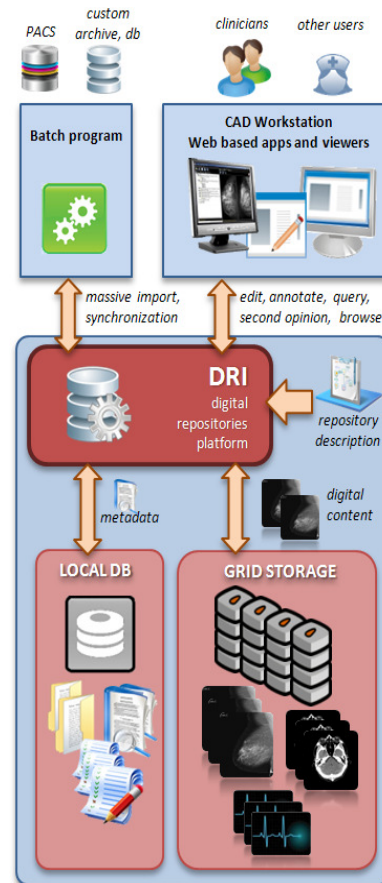
We are currently using DRI to host repositories of medical images (such as the one described in this paper) and repositories of digitalized ancient manuscripts. In a general way, we expect to use DRI to provide easy management of large scientific data collection stored in Grid infrastructures.

## 2.2 Using Grid infrastructures

The default DRI implemented uses gLite Grid middleware for authentication, storage of large digital content and computing power. A Grid infrastructure is typically made of a federation of sites (such as different data centres, hospitals, etc.), each one providing different amounts of Grid storage and computing power to the federation. Figure 2 shows a sample deployment scenario of several DRI instances over a shared Grid infrastructure. Note the following:

(1) Three sites participate in the federation: a *hospital* (deploying one DRI instance, a local database for metadata and a Grid site for storage), a *small clinic* (deploying one DRI instance and a local database) and one *data centre* (offering its Grid storage to the federation).

(2) When hospital users enter new repository items (patients and/or new mammograms), their DRI instance will store the metadata (such as patient information and/or diagnoses) on its local database, while storing the digitalized mammograms in the Grid storage space. Note that Grid policies established within the federation will determine where the digitalized mammogram will be physically stored, not necessarily on its local Grid storage.

(3) When clinic users enter new repository items, their DRI instance will store the metadata in their local database and mammograms in the Grid storage space according to the established Grid policies. Since they provide no Grid storage (maybe because they could not afford it) their mammograms will always be stored somewhere else.

(4) There might be many different Grid storage policies. For instance, *only files larger than 50Mb will be stored in the Data Centre Grid site*. Or, *the Hospital and the Data Centre Grid Sites will maintain replicas of files of the last three months*. Or, *data generated in the Clinic will only be stored in the Data Centre Grid site*, etc.

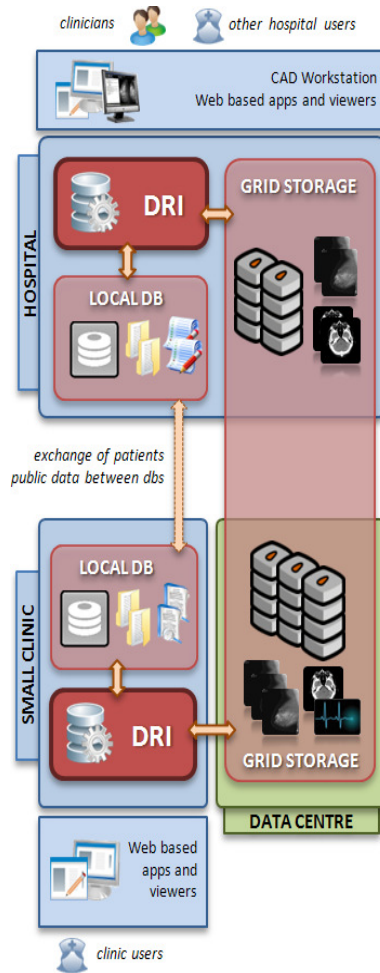(5) Both *local databases are isolated*. Users



**Fig. 2**: DRI deployment example

logging into the hospital DRI instance will only access patient data held in its local database (although mammograms might be stored somewhere else), and analogously for clinic users. The metadata spaces are not shared and therefore, no patient information is visible outside the institution scope.

(6) However, the clinic and the hospital *might agree to share some data*. For instance, since the clinic has few expert physicians for diagnoses, it might decide to synchronize regularly patient data from their local database to the hospital database, so that hospital expert physicians can occasionally diagnose patients. However, the clinic might decide not to expose patient sensitive data (name, etc.), but only anonymous records (patient studies). In any case, since the mammograms are stored in the federated Grid storage, hospital experts will be able to see and annotate them once they access to the patient data.

## 2.3 Open standards based technology

DRI exposes its functionalities through a well defined *web service* [4] so that expert users can develop their own applications on top of it. This allows developing a diversity of tools to exploit content offered through DRI (such as the web application mentioned above). For instance, we one can create a batch application to import or synchronize repository data with any external source or develop a specialized viewer application for a complex repository. This is the case for the mammograms repository described in section 3 below, where we adapted and improved an open source DICOM workstation [5] to work against DRI to support managing mammograms and building CAD systems using the repository content.

DRI is built entirely as a Java J2EE application [6] offering standard HTML/javascript web interfaces for users. It uses gLite middleware [7] to interface with Grid infrastructures for authentication and storing large digital content and regular JDBC drivers to access local databases (such as MySQL, Postgress, Oracle, etc.) for storing metadata.

DRI is formed by pluggable modules so that new behaviors can be added into the platform to allow its interaction with other resources. For instance, the *DRI Default Storage Module* defines how metadata is separated from digital content and stored differently in a *gLite Storage Element* and in a *JDBC database*. This has been our choice for its default behaviour. However any programmer can use the DRI APIs (*Application Programming Interfaces*) and create new storage modules to interact with other kinds of resources such as a different Grid infrastructure based on Globus [8] or UNICORE [9], for instance), database or custom data storage (PACs, etc.)

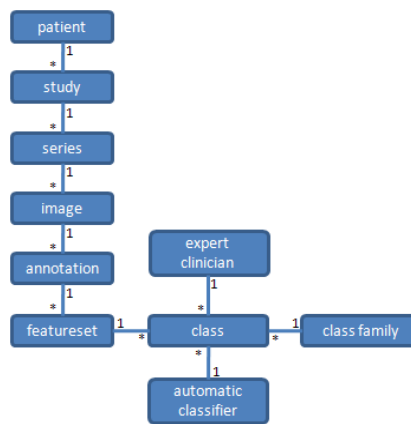The DRI platform is to be offered under a GPL-like open source licensing schema.

## 3 The Mammograms Repository

The mammograms repository data model here presented (Figure 3) is a subset of the DICOM medical file format [10] that was modified to store and manage specific patient information related to digital mammography images. In short, it deal with

many patients, each one might undergo one or more studies, each study is made of one or more series, each series contains one or more images and each image contains one or more annotations. An annotation corresponds to a mark made on the mammogram image, such as a circle, a text note or a segmentation of a certain region.

In turn to support manual or automated classification, each annotation may have one or more feature sets associated. A feature represents a certain characteristic of the annotation, such as brightness, elongation, etc. Each feature set is assigned by an expert clinician or an automatic classifier (for instance an artificial neural network – ANN) to a certain class belonging to a class family such as the BI-RADS [11] family of classes.



**Fig. 3:** The mammograms data model

The model supports the same feature set (and thus the same annotation) to be given several classifications by different clinicians and automatic classifiers under different class families. This supports storage of a variety of sets of experiments of classification runs performed both by human experts and automatic classifiers, so that later become available for statistical analysis. Figure 3 shows this data model.

This data model is represented in a repository description file in XML format. Each entity (patient, study, etc.) is described in a **TableName** XML tag.

The following is an excerpt of the definition of the *Image* entity:

```
<TableName name=MGImage
 forgeinIdAttr=SeriesID>
    <attr name="ImageID">
        <dbAttrName>ImageID</dbAttrName>
        <type>int</type>
        <dbAttrType>int</dbAttrType>
    </attr>
    <attr name="xSize">
        <dbAttrName>xSize</dbAttrName>
        <type>int</type>
        <dbAttrType>int</dbAttrType>
    </attr>
    <attr name="Mammogram">
        <dbAttrName>Mammogram</dbAttrName>
        <type>LNF</type>
```

```
        <dbAttrType>Varchar(255)</dbAttrType>
    </attr>
</TableName>
```

Note how the **Mammogram** field is marked as *LFN*. With this, DRI knows that this is *large digital content*, will manage its persistence in Grid storage and will provide appropriate graphical tools for users to interact with it.

Finally, a collection of computer vision algorithms were added in the mammograms viewer such as: region of interest (ROI) selection, image enhancing, livewire and snake segmentation techniques, features extraction facilities and functionality to generate and train ANNs that include feedforward back propagation and self organizing features maps among others models (see Section 4 below). With this and the data model described above, a platform that supports the entire process for building a wide range of mammography classifiers systems based on ANNs is provided.

## 4   Mammography Image Analysis Method

As it was mentioned before the DRI developed is a general framework to store and process digital images repositories, with focal point in medical imaging. Due to this, to evaluate and test DRI capabilities in mammography image analysis an experimental CAD was implemented using as start point the method developed by Lopez et al [12], which includes five main steps: the selection of region of interest (ROI), image equalization, segmentation, features extraction and classification of possible pathological lesions (PL).



**Fig. 4.** Mammogram analysis sequence

### 4.1   ROI Selection

ROI Selection allows that specialized users to select and process only the suspicious regions to contain PLs (see Figure 4a-b).

## 4.2    Image Equalization

Contrast - Limited Adaptive Histogram Equalization (CLAHE) is a subtype of class equalization algorithms known as Adaptive Histogram Equalization, which divide the images into contextual regions (called tiles) and applies the histogram equalization method to each one. This evens out the distribution of the used grey values highlighting the key features of the image. The neighbouring tiles are then combined using bilinear interpolation to eliminate artificially induced boundaries. The contrast, especially in homogeneous areas, can be limited to avoid amplifying any noise that might be present in the image (see Figure 4c). This step allows enhancing image details, which are critical in the image segmentation step.

## 4.3    Segmentation

Deformable models (snakes) have been used successfully in diverse image segmentation tasks [13; 14]. Two interactive segmentation techniques were implemented and offered: a variant of gradient vector flow snakes and livewire [15]. But finally the specialized medical personnel selected used the livewire as default segmentation technique.

Livewire (or intelligent scissors) is an interactive boundary tracing technique, considered as a competing technique to snakes. This technique allows (with minimal user interaction) to exercise online control over the segmentation process.

The livewire technique is applied first to produce a curve contour approximation (edge points) of the PL in a selected ROI. But because, resulting curve is not a continuous curve, it is then interpolated using a spline function to generate a better (more precise) PL contour (see Figure 4d).

## 4.4    Features Extraction

Lopez et al [12] demonstrated that a features vector formed only by four features (see Figure 4e) can be enough to obtain a correct classification of PL in mammograms, in the performed experimental study the same features vector was used, which was composed by: object area, brightness (mean of gray levels inside of segmented PL), object shape and elongation. Mathematical model used to compute the features vectors for each PL was as follows:

$$Object \equiv \text{pixel set of segmented PL} \qquad (1)$$

$$Edge \subset Object \text{ edge pixels} \qquad (2)$$

$$Area = \left| Object \right| \qquad (3)$$

$$Perimeter = length(E) \qquad (4)$$

$$Elongation = \frac{diam}{DIAM} \qquad (5)$$

$$Shape = \frac{Perimeter \cdot (diam + DIAM)}{8 \cdot Area} \qquad (6)$$

$$Brightness = mean(Object) \qquad (7)$$

where *diam* and *DIAM* represent minimum and maximum diameters respectively.

## 4.5    Classification

The classification step is intended to offer a "second opinion" (diagnosis) about possible PL present on the mammography images. Two classifiers were developed and implemented on our DRI platform: a first one (supervised) based on the feedforward backpropagation (FFBP) ANN model and the second one (unsupervised) based on the self organizing maps (SOM) ANN model. Figure 4f shows the set of classes recognized by the classifiers. The FFBP model is one of the more studied ANN by the scientific community and the most common used in many medical applications. Morphologically, the FFBP is formed by a set of organized neurons in layers: input, hidden and output layers. Network architecture is determined by the number of neurons in the hidden layers. The learning process of a FFBP network is characterized to be supervised; the network parameters, known as weights are estimates from a group (pairs) of training patterns composes for input and output patterns $\{(x^t, y^t)\}t = 1..n$. The backpropagation algorithm [16] is a generalization of the proposed rule delta by Widrow-Hoff [18]. The term "backpropagation" refers to the form in that the cost function (gradient) is calculated for the FFBP network. Therefore, the network adjust takes place as a result of the estimation of weights parameters. The learning involves an adjustment of the weights comparing the desired output with the network answer so that the error is minimized. A FFBP-based classifier was implemented formed by four layers: an input layer with 4 neurons, two hidden layers with 14 and 8 neurons respectively and an output layer with 13 neurons. Each neuron from output layer represent one (benign or malignant) pathological lesion class (calcifications, well-defined/circumscribed masses, spiculated masses, ill-defined masses, architectural distortions and asymmetries) including a normal image. SOM is an ANN model developed in 1980 by Teuvo Kohonen [18], consisting of several map modules (neurons) that have been used for tasks similar to those to which other more traditional neural networks have been applied: pattern recognition, robotics, process control, and even processing of semantic information. The spatial segregation of different responses and their organization into topologically related subsets results in a high degree of efficiency in typical neural network operations. It has a strong physiological inspiration, as it is based on the topological map that exists in the brain cortex. The cortex is organized so that topologically closer neurons tend to produce answers to the same kind of stimulus; this is one of the reasons why it is largely employed in visual pattern recognition [19]. Details about SOM algorithm (model, training, etc.) implementation will be found in [18; 19]. A SOM-based classifier was implemented formed by two layers: an input layer with 4 neurons and an output layer with 13 neurons. Each neuron from output layer represent one (benign or malignant) pathological lesion class (calcifications, well-defined/circumscribed masses, spiculated masses, ill-defined masses, architectural distortions and asymmetries) including a normal image.

# 5 Results and discussion

Based on the mammograms information file attached to Mammographic Image Analysis Society (MIAS) database we created an experimental digital repository that we named MIAS-DRI. MIAS-DRI is formed by 322 mammogram digital images, where each one has associated the following information: the mammogram reference number, background tissue, class of abnormality present, severity of abnormality; image-coordinates of centre of abnormality and approximate radius (in pixels) of a circle enclosing the abnormality. Mammograms are gray levels images with a resolution (size) of 1024 x 1024 pixels and 8 bits per pixel.

**Table 1.** Classification results

| Classes | # of ROIs | Neural Networks | | | | | |
| | | FFBP | | | SOM | | |
| | | TP | FP | FN | TP | FP | FN |
|---|---|---|---|---|---|---|---|
| Calcifications | 14 | 11 | 1 | 2 | 11 | 1 | 2 |
| Well-defined/ circumscribed Masses | 12 | 11 | - | 1 | 10 | 1 | 1 |
| Spiculated Masses | 13 | 13 | - | - | 10 | 1 | 2 |
| Other Ill-defined Masses | 14 | 12 | 2 | - | 10 | 2 | 2 |
| Architectural Distortions | 15 | 14 | 1 | - | 11 | 3 | 1 |
| Asymmetry Masses | 12 | 11 | - | 1 | 8 | 1 | 3 |
| Normal | 20 | 18 | 1 | 1 | 17 | 1 | 2 |
| **Total** | **100** | **90** | **5** | **5** | **77** | **10** | **13** |

A representative MIAS-DRI dataset formed by 100 images ROIs was randomly selected, including examples of all pathological lesions classes: calcifications, well defined circumscribed masses, spiculated masses, ill-defined masses, architectural distortions and asymmetries; besides of normal images (see table 1). This dataset was divided into ten folds, using the statistical cross-validation (rotation estimation) [20] method for training and evaluating the accuracy of two ANN formed classifiers. These classifiers were trained ten-times, with a matrix formed by 90 vectors representing appropriated (tested) features (area, brightness, shape and elongation) extracted from already classified PL or normal tissue ROIs.

Classification results were expressed in terms of three parameters: True Positive (TP), False Positive (FP) and False Negative (FN). A TP is obtained when a mammogram PL is classified into the correct (benign or malignant) class. When a benign mammogram pathological lesion is incorrectly classified into another benign

class, into a malignant class or into a normal image, it is defined as a FP. A FN is obtained when a malignant mammogram pathological lesion is incorrectly classified into another malignant class, into a benign class or into a normal image.

Table 1 shows the results achieved after training and simulation. It can be observed that FFBP-based classifier obtained better results 90% of TP classified ROIs compare with the SOM-based classifier with only 77% of TP classified ROIs.

## 6 Conclusions

The usage of the proposed DRI software platform to effectively host a digital repository of mammograms images and related information was successfully validated. This also includes a general framework to test image analysis algorithms that allows the creation of domain specific medical analysis applications. Based on this, an experimental mammography CAD system was developed, to offer a second opinion (diagnosis) about PLs, which demonstrated a satisfactory classification performance. This CAD provides complete functionality to support the working lifecycle of mammograms through different stages (archiving, retrieving, training and analysis). With this, Grid infrastructures can effectively be used to build large collections of mammograms and use them to build and train CAD systems.

This result constitutes the foundation to build production-quality systems to exploit Grid infrastructures for medical image communities. The future work is now focused on optimizing the graphical user interfaces, improving system's scalability and manageability to reach production quality levels and building more complete and specialized CAD systems to address specific needs (such as for certain patient populations, geographical regions, pathologies, etc.)

## REFERENCES

1. Foster, I., Kesselman, C: The Grid: blueprint for a new computing infrastructure. Morgan Kaufmann Publishers, (1999).
2. Arms W.: Key Concepts in the Architecture of the Digital Library. DLib Magazine (1995)
3. Amendolia, S.R. et al.: Managing Pan-European mammography images and data using a service oriented architecture, In: Proceedings of the IDEAS'04 Workshop  pp. 99-108 (2004).
4. Machiraju, V. et.al: Web Services: Concepts, Architectures and Applications. Springer, (2004).
5. TUDOR, The TUDOR DICOM Viewervhttp://santec.tudor.lu/project/optimage/dicom/start
6. Alur,D., Crupi, J., Malks, D.: Core J2EE patterns: best practices and design strategies Prentice Hall, (2003)
7. EGEE, The gLite middleware, http://glite.web.cern.ch/

8. The-Globus-Alliance, The Globus middleware, http://www.globus.org/
9. UNICORE, The UNICORE middleware, http://www.unicore.eu/
10. NEMA, The DICOM Standard, http://medical.nema.org/
11. D'Orsi, C.J., Bassett, L.W., Berg, W.A. et.al.: Breast Imaging Reporting and Data System: ACR BI-RADS-Mammography, American College of Radiology (2003).
12. López, Y., Novoa, A., Guevara, M., Silva, A.: Breast Cancer Diagnosis Based on a Suitable Combination of Deformable Models and Artificial Neural Networks Techniques, Progress in Pattern Recognition, Image Analysis and Applications, pp. 803-811 (2008)
13. Ghassan Hamarneh, G.H.: DTMRI Segmentation using DT-Snakes and DT-Livewire, IEEE International Symposium on Signal Processing and Information Technology, pp. 513 - 518 (2006)
14. Jianming Liang, T.J., Terzopoulos, D.: United Snakes, Medical Image Analysis, pp. 215-233 (2006)
15. Chenyang, X., Prince, J.L.: Snakes, shapes, and gradient vector flow. Image Processing, IEEE Transactions, vol 7, 359-369 (1998)
16. Rumelhart, D.E, Hinton, G.E., Williams, R.J.: Learning representations by back-propagating errors. Nature 323 pp. 533-536 (1986)
17. Widrow, B, Hoff, M.E: Adaptive switching circuits, Neurocomputing: foundations of research, MIT Press, pp. 123-134 (1988)
18. Kohonen, T.: The self-organizing map. In: Proceedings of the IEEE 78 pp 1464-1480 (1990)
19. Vijayakumar, C., Damayanti,G., Pant, R., Sreedhar, C.M.: Segmentation and grading of brain tumors on apparent diffusion coefficient images using self-organizing maps. Computerized Medical Imaging and Graphics vol 31 473-484 (2007)
20. Kuncheva, L.I.: Combining Pattern Classifiers, Wiley-Interscience (2004).

# Ball Detection and Tracking Using Color Features

Catarina B. Santiago[1,2], Armando Sousa[1,2], Luís Paulo Reis[1,3], M. Luísa Estriga[4,5]

[1] Faculty of Engineering of the University of Porto, R. Dr. Roberto Frias s/n, Porto, Portugal
[2] Institute for Systems and Computer Eng. of Porto, R. Dr. Roberto Frias s/n, Porto, Portugal
[3] Artifficial Intelligence and Computer Science Lab, R. Campo Alegre 1021, Porto, Portugal
{catarina.santiago, asousa, lpreis}@fe.up.pt
[4] Faculty of Sports of the University of Porto, R. Dr. Plácido Costa, nº91, Porto, Portugal
[5] Centre of Research, Education, Innovation and Intervention in Sport, R. Dr. Plácido Costa, nº91, Porto, Portugal
lestriga@fade.up.pt

**Abstract.** Recent years brought an increased interest on efficient analysis of the performance of players during a sports match. A player's performance is intimately related with the way he/she interacts with the ball. The focus of this paper is exactly on detecting and tracking the ball. We believe it is possible to construct an automatic system for detecting and tracking a ball based only on visual information. So far our research is based on a GigEthernet camera mounted in the laboratory. Initially the user has to define the color ball subspace. Afterwards a background subtraction helps highlighting the regions of interest which are then scanned to detect the ball color and form color blobs. The tracking algorithm uses past information and ball velocity as inputs in order to define future probable regions. Results are still preliminary but show that with a careful color calibration, around 98% of detection can be achieved.

**Keywords:** Image Processing, Tracking, Sports, Features Detection.

## 1 Introduction

In recent years there has been a growing interest by the sports experts (coaches and teachers) in using automatic techniques in order to record and analyze game sequences or training sessions.

Actually, the simple fact of providing a recorded video of the entire field where all the players and ball movements can be visualized in a single image is already very useful for coaches/teachers. This way, after the game the sports experts can identify weak points, define measures and new training directions to improve the team global behavior.

During this analysis the sport's expert usually focuses on identifying the players, their positions related to the other players (of the same team and of the opponent team) and also related to the ball. In fact, one of the most important elements inside the field area is the ball, because all the players' actions are conditioned by the position, the direction and the velocity it has.

However, this kind of game analysis (or game annotation) is very time consuming when performed by a person, therefore it appears the necessity to develop automatic

systems that can perform these same tasks with the advantages of being able to handle huge amounts of data and execute a systematic evaluation.

This paper focuses on a single part of the referred game analysis that is the tracking of the ball. The main objectives are to design a vision system that is able to detect and track the ball during a handball game providing metrics such as the position and direction. The system must be non intrusive therefore no special tags or colors should be placed on the ball, since most of the championships do not allow it and also supply real time information.

The work developed so far has been tested under more or less controlled conditions on a laboratorial environment.

The remainder of this paper is organized as follows. Section 2 presents some background information about ball tracking research. Section 3 gives a detailed explanation of the image processing system which includes the ball detection and tracking. Section 4 shows the results achieved so far and the last section refers to the conclusions and some future work.

## 2   State of the Art

There are two main types of technologies usually applied to the tracking problem: intrusive where special tags or sensors are placed in the targets and non intrusive where there are no strange objects in the game environment.

Intrusive systems are based on sensors placed in the targets that must be detected and tracked. These sensors broadcast a signal that is picked up by a few receptors placed at strategic points. The position of the targets is afterwards determined by means of triangulation. Witrack and Ubisense [1], [2], [3] are two of these systems that provide small sensors which can be easily placed on the ball.

However, intrusive systems are not allowed in most of the competitions, therefore most researches use images captured by TV broadcast systems or by dedicated cameras placed in very specific locations of the field.

Of course the tracking problem using a vision system becomes a very complex task because the ball can achieve very high speeds or even be occluded by the players. Nevertheless a vision system has the great advantage of not being intrusive and since it provides so much information at the end it pays off to use it.

In order to deal with this kind of complexity several color image segmentation techniques have been developed. A very detailed survey on this subject can be found on [4]. The authors provide a detailed discussion about the five main categories of image segmentation: histogram thresholding, clustering, region growing, edge detection, fuzzy approaches and neural networks as well as several existing color spaces.

For our specific problem, ball detection and tracking, we can find several works applied to tennis [5], [6], baseball [7] and soccer, however for indoor sports like the one we are studying (handball) we could not find any.

The philosophies behind tennis and baseball are much different from that of handball. So the focus of our research is on what has been done for soccer applications. Yet there are still huge differences between soccer and handball.

Handball is played on an indoor sports hall, while football is in an outdoor stadium, the handball ball is much smaller than a football ball and is more often occluded since it is transported by the players' hands instead of the feet.

Due to the geometry of the ball it seems appropriate to use the circle Hough Transform (HT) in order to perform its detection, in fact this technique is widely used in applications where the objects to be detected present a round shape which frequently happens in industrial and medical applications, among others [8], [9], [10]. However the HT requires a high computational effort and is very time consuming, which is not affordable in real time applications. Therefore several authors use it as basis but perform a few adjustments in order to fasten it. This is the case of D'Orazio et al [11] who developed a system which is able to recognize the ball in soccer images from a dedicated camera and aid referees to verify a goal event.

They are truly focused on providing a real time system that is able to recognize the ball in real images where light conditions can vary and the background is constantly changing.

Two different techniques are combined, a fast circle detection based on directional circle HT to identify regions of interest (ROI) and a three layer neural classifier to evaluate these ROIs in order to accept or discard the hypothesis. The usage of constrains based on *a priori* knowledge of the ball dimensions and the influence on the shape of the ball due to light conditions helps fastening and improving the algorithm.

Although their results seem very promising they consider that the ball is not visible if more than 50% of its surface is occluded, which seems a very low value.

More recently their work has evolved from a single camera system to a four high frame rate camera system [12] placed on the sides of the goals. Having two cameras pointed to the ball allows them to detect its 3D position through homography.

The usage of cameras with such a high frame rate provides a more reliable picture of the ball movement and minimizes the effect of motion blur. However these cameras are highly expensive and such high frame rates leave little time to perform the image processing.

Motion blur occurs when objects travel at high velocities (see Fig.1).



**Fig. 1.** Ball at low speed (left) and high speed (right)

In their new approach, first a moving object segmentation based on background subtraction is performed in order to detect the ROIs and only after they apply their fast circle HT  and the neural classifier. This results in a much faster algorithm since the HT is not applied to the entire image but only to the ROIs. Nevertheless they only detect and track the ball near the goals areas.

Liang et al [13] use images provided by TV broadcast. The lower frame rate of these cameras compels them to take into consideration the motion blur caused by the ball speed and make use of the image color properties.

Under the assumption that the ball is nearly white in long view shots, white pixels are first segmented and candidate regions are defined using physical restrictions. They use candidate regions from a set of consecutive frames to construct a weighted graph, where each node represents a ball candidate. The optimal path (which represents the true location of the ball) is extracted using the Viterbi algorithm. They also use a Kalman filter to predict the next ball location.

Results show that occlusion and motion blur problems are not well handled and the ball can be confused with the players' socks.

Images from broadcast soccer are also used by Pallavi et al [14]. Initially each frame is categorize into long, medium or close shot based on the ratio between the grass pixels (green color) and the field region using the YIQ color space which minimizes the effects caused by light changes. Close shot images are not used because they usually represent the face of a player.

Like in [11], the first ball candidates are determined using the circle HT with some radius restrictions. Afterwards, if it is a medium shot view the velocity of each of the remainder ROIs is calculated using the optical flow velocity method of Horn and Schunck. The ROI having the highest velocity is identified as being the ball.

For long shots the non-ball candidates are removed instead of detecting the best ball candidate. Filters based on previous defined heuristics are used to eliminate the weakest candidates and the remnant ones will be used to build a directed weighted graph where the longest path represents the ball trajectory.

The usage of the HT and of predefined heuristics although providing good results can lead to miss detections if the ball is partially occluded or if by some chance the ball behaves contrarily to the heuristics. The time it takes to process a single frame is too long which makes impossible to have real time processing.

From what has been said it is possible to see that in most of the cases researchers focus their study on specific areas of the field (goal area) or use the images provided from TV broadcast, but for us it seems more appropriate and useful to be able to see the entire field all the time so that after determining the ball position it is also possible to evaluate the positions of the players regarding to the ball. Ren et al [15] share this same concern and use a dedicated 8 camera system that allows seeing the entire field.

Ground plane velocity, longevity, normalized size and color features are the inputs of a Kalman filter that assigns a likelihood measure to candidate regions. They further refine their detection using occlusion reasoning, backtracking and the 3D ball position. Results show that backtracking is a key feature of their algorithm.

Our approach is based on the notion of color blob and uses color features extracted from the images in order to create a color ball subspace as it was proposed in [16]. We envision that the usage of simple primitives such as color features will provide a robust and fast algorithm to tackle our problem.

## 3  Ball Detection and Tracking

Our system is composed of two main sub-systems: the ball detection, responsible for identifying in the image the ball position, and the ball tracking, that must be able to keep the track of the ball.

### 3.1 Ball Detection

The ball detection sub-system is the most important part of the overall system because it must be robust enough so that the tracking can have solid basis to work with.

Our image processing is based on color blob definition and the ball is seen as a color blob.

**Color Definition.** Before starting the ball detection it is necessary to define the color of the ball. This is a very important task since a good color calibration will influence the success of the subsequent steps.

The color calibration is performed using the mouse. The pixels selected by the mouse are used to fill in a special lookup table. This lookup table corresponds to a three dimensional vector where each dimension represents a color component of the RGB color space. Whenever a pixel is selected its RGB component is used to address the position on the lookup table and that position is marked with a special value indicating it belongs to the ball color subspace.

The selected pixel is also used as a seed to perform a physical flood and a color growth in the image.

The physical flood process consists in including in the ball color subspace not only the pixel selected by the mouse (one single pixel that correspond to the seed pixel) but also the pixels in the surrounding area as long as their color is similar (using the Euclidean distance) to the seed pixel.

The color expansion is performed on the HSL space in order to minimize the effects of shadows and light variations. Only two dimensions are expanded, saturation and luminance. The expansion is not performed in the hue direction since in this case the intention is not to expand the team color to other colors but expand the same color in terms of white dilution (saturation) and brightness (luminance). Fig.2 illustrates the physical flood process.



**Fig. 2.** Physical flood color expansion (HSL color space)

Color growth also uses the selected pixel as seed but the expansion is performed on the HSL color space without the restrictions of the physical neighborhood.

The conjunction of these two methods allowed achieving very good results in a short time.

**Background Subtraction.** Most of the image area had non-useful information. In fact, the regions of interest were the ones that included the ball therefore a background subtraction was used to highlight those regions.

A pure background subtraction would not have been the best choice since it would fail due to the usual brightness of the handball field floor. Therefore we developed a background subtraction performed on a per pixel basis: whenever a pixel of the image under analysis has a color similar to the color of the same pixel in the background image (measured by the Euclidean distance), the pixel of the image under analysis is marked with the white color as can be seen on Fig.3. The threshold used to perform the pixel subtraction was determined experimentally.



**Fig. 3.** Background subtraction in a slow ball (left) and in a fast ball (right)

The background image is built using initially an image without the ball that is subsequently updated. Each time a new frame is analyzed its content is used to perform an update into the background image. This way a dynamic background subtraction can be achieved.

For the update, the previous background and the image under analysis are divided into a grid and each zone of the two images is subtracted according to equation 1, where A corresponds to the image under analysis, BK to the background image, R, G, B are each of the components of the RGB color space and x and y correspond to the pixel coordinates in the image.

$$\sum_{ys_{\min}}^{ys_{\max}} \sum_{xs_{\min}}^{xs_{\max}} \left( \left( A_{R(x,y)} - BK_{R(x,y)} \right)^2 + \left( A_{G(x,y)} - BK_{G(x,y)} \right)^2 + \left( A_{B(x,y)} - BK_{B(x,y)} \right)^2 \right) \qquad \textbf{(1)}$$

If the result is lower than a given value (determined experimentally) an average between the two images is performed and that section of the background image is updated. Otherwise the section will remain untouched.

**Color Detection.** Once the ball color subspace is defined and the regions of interest are highlighted it is possible to start the color detection. This color detection corresponds to spot, in the image, pixels that belong to the color ball subspace. For that, the entire image is scanned and the color value of each pixel is tested.

This test consists in verifying if the entry on the color lookup table is filled with the special value that indicates it belongs to the ball color subspace or if it is empty. If it corresponds to the ball color subspace then the color of the pixel is replaced with a specific color identifier.

Fig.4 exemplifies this concept where the color identifier is blue. On the right image it is possible to see (blue dots) the color ball subspace defined as a subspace of the RGB color space.
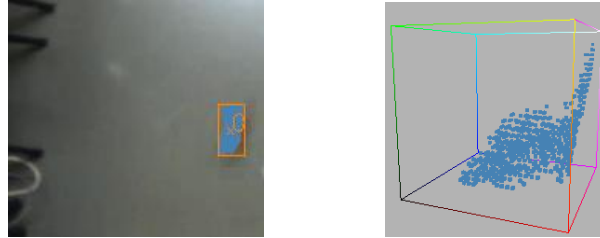
**Fig. 4.** Ball with pixels marked (left) and color ball subspace (right)

Once the pixels are classified a mode filter is applied in order to remove noise. This filter is based on a 3x3 window and consists in analyzing the 8 neighbors of a pixel and detecting if half or more of them belong to the ball color subspace. If this is true then the pixel under analysis is replaced by the color ball identifier.

**Blob Aggregation and Characterization.** So far, the only information collected from the image, is if a pixel belongs to the ball color subspace or not. It is still necessary to establish a relationship between pixels belonging to the same color blob.

The algorithm responsible for establishing this relationship comprises two steps. The first step is based on a per line scan detection and the information is stored in a way similar to that of a run-length encoding using three parameters to define a blob of pixels: $y$, $x_{min}$ and $x_{max}$.

Whenever a pixel belonging to the ball color subspace is reached its x ($x_{min}$) and y values are stored and the subsequent pixels of the same line are checked to see if they also belong to the ball color subspace. Once reached the last pixel belonging to the color ball subspace or the end of the line the $x_{max}$ value is stored.

If the end of the image width was not reached the scan continues in that line and if meanwhile another pixel of the ball is found and its distance to the $x_{max}$ pixel is inferior to a given value these new pixels are considered belonging to the first part.

As a result from the previous procedure a series of single lines are identified as belonging to the specific color ball subspace and they still need to be joined to form a single blob, which represents the second step.

If the distance between two of these lines is small then they are considered as being part of the same blob and are connected together as can be seen by Fig. 5.

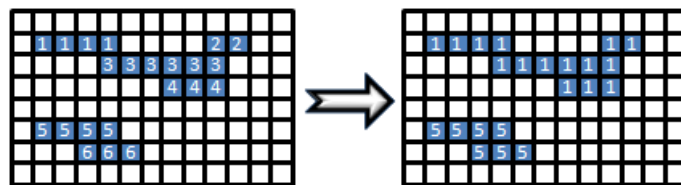With this step the most probable candidates areas were determined and correspond to each blob identified in the image.



**Fig. 5.** Single line aggregation to form a color blob

Once the color blobs are grouped it is possible to perform an estimate of its centre of mass.

So far each blob has been completely identified through the color blob process and characterized in a basic way ($x_{min,max,cm}$, $y_{min,max,cm}$).

False candidates are eliminated using physical constraints such as the area of the rectangle that best fits the blob and the density of pixels belonging to the color ball subspace inside the bounding box.

Until now all the processing was performed on the pixel's world coordinates but it is still necessary to perform the conversion into the real world coordinates.

Since there was little barrel distortion, because the camera was too near of the ground, this conversion consisted on multiplying both the x as the y components by scalar factors that were determined by placing a static object with known size on the ground. However it is predicted that when using this same system under a real game situation where the camera can be 8 meters above the ground the barrel distortion effects will become more severe and will have to be taken into account.

### 3.2 Ball Tracking

Once the ball is detected it is possible to perform its tracking. As stated before the characteristics that define the ball are its area and centre of mass. The tracking method is based on this information and on defining a probable area around the ball that defines its next position.

This probable area takes into account the position and the maximum velocity the ball can achieve.

Once a new frame is picked, the detection phase starts by identifying all the color blobs and those that are identified as being a ball are characterized. The tracking algorithm compares these characteristics with the ones from the previously identified blobs and if they fit inside the probable area of one of those blobs then the new blob is assumed as being the sequence.

At every instant the tracking algorithm can have access to a limited number of previous positions and ball characteristics that are kept on memory. With this approach it is possible to perform the tracking of the ball even if on the previous frame the ball was not detected.

## 4   Results

Results presented in this paper are based on 2 sample footages of around 7 seconds each where fast ball and slow ball movements were recorded (the short duration of the sample footages is due to the limited space on the laboratory). The sample footages were recorded in MJPEG with resolution of 412x708 (in order to catch only the area of interest the ROI property of the camera was used) and 30 frames per seconds.

The left side of Fig.6 shows the original image and the right side the final result after applying the image processing and tracking algorithms. The orange square represents the rectangle that best fits the color blob (ball).
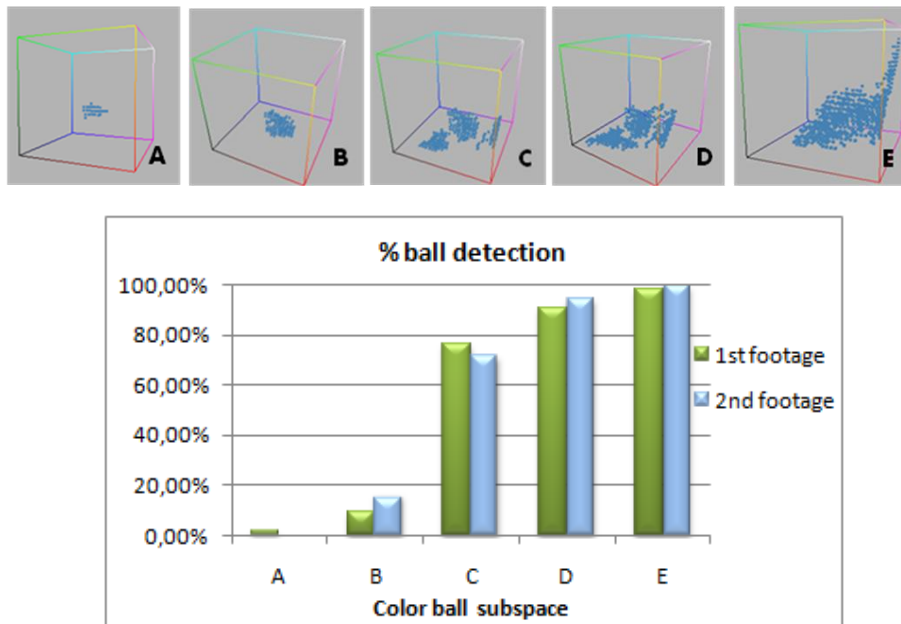
**Fig. 6.** Image before (left) and after (right) processing

Experiments made also showed that the color calibration step is very important and a key factor to have a suitable ball detection as can be seen on Fig. 7.

On images A to E it is possible to see the several tested color ball subspaces represented by the blue dots inside the RGB cube.

A poor representative color ball subspace like A produced very low detection rates for both footages, while a good representative color subspace like E, that took into consideration the different orange shades the ball can have (depending on the light conditions throughout the entire test zone), produced higher detection rates.



**Fig. 7.** Color calibration: color ball subspaces from A to E (top) and corresponding detection rates in both footages (down)

A good color calibration aided by the usage of an adaptative background subtraction was able to overcome severe light conditions as deep shadows or specular reflection as can be denoted on Fig.8. On the processed images it is possible to see the

rectangle that best fits the ball (in orange) which indicates the algorithm was able to correctly identify the ball.



**Fig. 8.** Ball detection under severe light conditions

It is important to notice that specular reflection of the lighting system also appears intensely on the sports hall where we intent to mount our system.

On the first sample footage the ball was detected in 98% of the frames and on the second in 100%. The non-detections on the first sample footage were related with the strong point of light (only one out of the three frames that captured the ball in that zone was correctly detected –Fig.8) and the fast velocity the ball achieved in that area.

Fig. 9 shows the performance of the tracking algorithm. On both sample footages the ball was thrown on the right side of the figure and when it reached the work bench on the left side it bounced back.

Each color section represents a period of time where the tracker was able to correctly follow the ball. Sections red, green and blue correspond to the first sample footage, sections rose and purple to the second.

From section green to blue and red to purple it was impossible for the tracker to behave correctly since the ball was completely hidden below the laboratory work bench.



**Fig. 9.** Performance of the tracking algorithm

The average processing time per frame was 34ms in a laptop computer with 1MB L2 cache and powered by an Intel T2130 processor running at 1.86GHz, under Windows Vista operating system.

## 5 Conclusions

This paper presented a system for detecting and tracking a ball. The main objectives were to develop a system that could handle such a complex task (ball is a small object that travels at high velocity).

Our system is composed of two main blocks, one for the image processing where color blob techniques are used in order to detect the ball and the other for tracking it.

The tracker is based on very simple primitives, previous ball position and maximum velocity, to define an estimative of a future probable area.

Tests conducted on the laboratory using a single GigEthernet camera with a ROI of 412x708, showed that the ball can be detected with high accuracy (between 98% and 100%) and tracked in sample footages of 7 seconds where its velocity goes from fast to slow. Nevertheless and since vision systems are especially influenced by light conditions, some more tests using a regulated projector are already envisaged in order to have a better assessment of the robustness of the algorithms.

Although it was not possible to test the system in a real environment (handball match), the information gathered was very useful to understand the dynamics of a ball and also to have a better perception of its characteristics on a video. Therefore the basis for a ball detection and tracking system were built and in future we believe it can be adapted to a real game situation using a multi-camera system.

It was also possible to confirm our initial hypothesis that a visual system can be a very powerful tool and correctly detect and track a ball.

Future work will concern on making both the ball detection as well as the tracking algorithms more robust so that they can handle other "objects" that appear in the game area (players, referees, marks on the ground) as well as different light conditions which often occur in real and non controlled environments. Ball characteristics like shape (with and without motion blur) and uniform color will be taken into consideration for the detection and Kalman filters for the tracking.

We also intend to perform tests in a real game situation in order to evaluate the performance of the algorithms and assess the potential of a stereo vision system that could provide the 3D position of the ball.

## References

1. Beetz, M., Kirchlechner, B., Lames, M.: Computerized Real-Time Analysis of Football Games. IEEE Pervasive Computing, vol. 4, pp. 33 – 39. IEEE Computer Society, New York (2005)
2. Ubisense.: New ID technology to track patients and records. Card Technology Today, vol.15, pp. 9 (2003)
3. Ubisense.: "Ubisense Series 7000 IP Rated Sensor" (2009)

4. Cheng, H.D., Jiang, X.H., Sun, Y., Wang, J.: Color image segmentation: advances and prospects. In. Pattern Recognition, vol. 34, pp 2259 – 2281. Elsevier Science Inc., New York (2001)

5. Yu, X., Simh, C.-H., Wang, J.R., Cheong, L. F.: A trajectory-based ball detection and tracking algorithm in broadcast tennis video. In: 2004 International Conference on Image Processing, vol.2, pp. 1049 – 1052. IEEE Press, New York (2004)

6. Yan, F., Christmas, W., Kittler, J.: Layered data association using graph-theoretic formulation with applications to tennis ball tracking in monocular sequence. In: IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 30, pp. 1814 – 1830. IEEE Computer Society, New York (2008)

7. Chu, W.T., Wang, C.W., Wu, J.L.: Extraction of baseball trajectory and physics-based validation for single-view baseball video sequences. In: 2006 IEEE International Conference on Multimedia and Expo, pp. 1813 – 1816. IEEE Press, New York (2006)

8. Huang, H., Chen, C., Jia, Y., Tang, S.: Automatic Detection and Recognition of Circular Road Sign. In: Mesa 2008, pp. 626 – 630. IEEE Press, New York (2008)

9. Chaichana, T., Yoowattana, S., Sun, Z., Tangjitkusolmun, S., Sookpotharom, S., Sangworasil, M.: Edge Detection of the Optic Disc in Retinal Images Based on Identification of a Round Shape. In: International Symposium on Communications and Information Technologies, 2008, pp. 670 – 674. IEEE Press, New York (2008)

10. Tomislay, S., Dubrayko, M., Toma, U., Viktor, M.: Machine Vision System for Seam Weld Detection in Longitudinally Welded Pipes. In: 19th International DAAAM Symposium, pp. 1073 – 1074. DAAAM International Vienna, Vienna (2008)

11. D'Orazio, T., Guaragnellab, C., Leo, M., Distante, A.: A new algorithm for ball recognition using circle Hough transform and neural classifier. In: Pattern Recognition, vol. 37, pp. 393 – 408. Elsevier Science Inc., New York (2004)

12. D'Orazio, T., Leo, M., Spagnolo, P., Nitti, N., Mosca, N., Distante, A.: A visual system for real time detection of goal events during soccer matches. In: Computer Vision and Image Understanding, vol. 113, pp. 622 – 632. Elsevier Science Inc., New York (2009)

13. Liang, D., Liu, Y., Huang, Q., Gao, W.: A Scheme for Ball Detection and Tracking in Broadcast Soccer Video. In: Advances in Multimedia Information Processing, LNCS, vol. 3767, pp. 864 – 875. Springer Berlin, Heidelberg (2005)

14. Pallavi, V., Mukherjee, J., Majumdar, A. K., Sural, S.: Ball detection from broadcast soccer videos using static and dynamic features. In: Journal of Visual Communication and Image Representation, vol. 19, pp. 426 – 436. Academic Press Inc., Orlando (2008)

15. Ren, J., Orwell, J., Jones, G. A., Xu, M.: Tracking the soccer ball using multiple fixed cameras. In: Computer Vision and Image Understanding, vol. 113, pp. 633 – 642. Elsevier Science Inc., New York (2009)

16. Sousa, A., Santiago, C., Reis, L.P., Estriga, M.L.,:Automatic Detection and Tracking of Handball Players. In: VipImage 2009 – II ECCOMAS Thematic Conference on Computational Vision and Medical Image Processing, pp. 213 – 219. Porto, Portugal (2009)

# Information Technologies

# Performance Evaluation of LINQ

Paulo Proença [1,2],

[1] ISEP – Institute of Engineering, Polytechnic of Porto
Porto, Portugal
prp@isep.ipp.pt
[2] FEUP – Faculty of Enginnering, University of Porto
Porto, Portugal
paulo.proenca@fe.up.pt

**Abstract.** With LINQ project, Microsoft adds general purpose query facilities to .Net Framework, but there will be productivity gains with the use of LINQ? Specially, when the world is comfortable with ADO.Net there will be an improved performance with LINQ to ADO.Net? In this paper we try to answer this question using a test scenario that pretend to evaluate the performance of the four basic database commands using three competing technologies: ADO.Net, LINQ to ADO.Net and Stored Procedures. This evaluation intends to prove that adding a new middleware layer to the application will not benefit the performance of the software. The measurements values resulting from the tests prove clearly that the purpose related with the development of LINQ were not issues of increased efficiency.

**Keywords:** LINQ, ADO.Net, .Net Framework, SQL Server

## 1 Introduction

Software is simple. It boils down to two things: code and data. Writing software is not so simple, and one of the major activities it involves is writing code that deals with data.

To write code, we can choose from a variety of programming languages. The selected language for an application may depend on the business context, on developer preferences, on the development team's skills, on the operating system, or on company policy.

With this in mind, Microsoft adds new features to .Net Framework that adds query capabilities to the .Net Languages like C# and VB.Net [1].

In the particular case of query relational databases, the ADO.Net is a stable and capable technology that can now be replaced by the LINQ to ADO.Net. However the use of LINQ in relational contexts implies the creation of a new middleware layer between the data access layer and the database, to convert the LINQ commands to the T-SQL commands.

This issue suggests a lack of performance in the database queries and was the motivation to begin this research. This paper proposes an architecture to determine the

performance levels of three alternative technologies (ADO.Net, LINQ to ADO.Net and Stored Procedures use) and attempts to explain the results.

This paper is outlined as follows. The Section 2 briefly describes the state of the art of LINQ. The implementation of the test scenario and the results display are presented in Section 3. Finally, the Section 4 concludes and discusses the results and gives some clues for future improvements of this work.

## 2     LINQ - The State of the Art

Language Integrated Query (LINQ) is a set of Microsoft .NET technologies that provide built-in language querying functionality similar to SQL, not only for database access, but for accessing data from any source. LINQ works as a middle tier between data store and the language environment as displayed in Figure 1.



**Fig. 1.** LINQ workflow (from language to data store)

LINQ aims to solve the problem of manipulating and selecting data from several data sources (databases, XML, object collections) providing a coordinated, consistent and efficient syntax from development environment and by using one chosen programming language, rather than switching between programming languages.

The current LINQ family of technologies and concepts allows an extensible set of operators that work over objects, SQL data and XML data sources. The generalized architecture of the technology also allows the LINQ concepts to be expanded to almost any data domain or technology.

LINQ comes with providers for in-memory object collections, XML documents, SQL Server databases and ADO.Net datasets. These providers define different LINQ technologies. LINQ architecture is displayed in Figure 2.
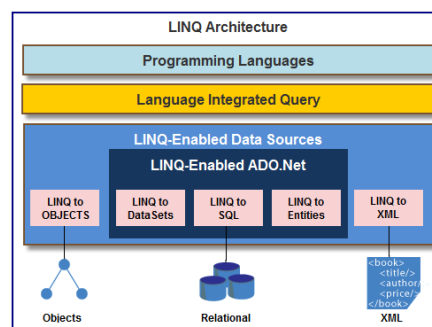


**Fig. 2.** LINQ Architecture

## 2.1 LINQ to Objects

This refers to the use of LINQ queries with any IEnumerable<T> collection directly

```
//Define a enumeration only containing child TextBoxes.
var myChildTextBoxes = from myControls in this.Controls
                       where myChildControl is TextBox
                       select myControls;
// Loop through TextBoxes.
foreach (TextBox myChildTextBox in myChildTextBoxes)
   myChildTextBox.Enabled = false;
```

In a common sense, LINQ to Objects represents a new approach to collections. In the old way, it was necessary to write complex *foreach* loops that specified how to retrieve data from a collection. In the LINQ approach, the developer writes declarative code that describes what he wants to retrieve.

LINQ defines a set of methods (called *standard query operators*, or, more recently, *standard sequencing operators*) and use lambda expressions [3] to language enhancement. This can be used including a reference to *System.Linq* namespace.

The previous example can be re-written in these terms using the *OfType* method provided by LINQ to any enumeration.

```
// Loop through my TextBoxes.
foreach (TextBox myChildTextBox in
this.Controls.OfType<TextBox>())
   myChildTextBox.Enabled = false;
```

LINQ code can be written using query expressions or method expressions (or a combination of both). The difference between the two is only syntactic, as they ultimately compile into the same thing.

## 2.2 LINQ to XML

The LINQ to XML provider converts an XML document to a collection of *XElement* objects, which may be queried using the standard LINQ engine [4][5].

```
XDocument loaded = XDocument.Load(@"contacts.xml");
var q = from c in loaded.Descendants("contact")
        where (int)c.Attribute("contactId") < 4
        select (string)c.Element("firstName") + " " +
               (string)c.Element("lastName");

foreach (string name in q)
  Console.WriteLine("Person name = {0}", name);
```

In the previous example all access to elements was casted to the correct type, and it was necessary to use the *Attribute* and *Element* methods.

To avoid these inconveniences Microsoft is developing LINQ to XSD [6][7] (last version is alpha 2), that allows to code against a strong-typed object model which makes the typing easier, and has access methods to retrieve elements and attributes.

The previous example can be rewritten according LINQ to XSD (as long as the XML has a schema).

```
// contacts is a predefined strong datatype
var xmlSource = contacts.Load(@"Contacts.xml");
var q = from c in xmlSource.contact
        where c.contactId < 4
        select c.firstName + " " + c.lastName;
foreach(string name in q)
  Console.WriteLine("Person name = {0}", name);
```

### 2.3    LINQ to SQL

This provider allows LINQ to be used to query SQL Server databases [8]. Since data is stored in a remote server, LINQ to SQL convert the local query (LINQ) to a SQL query which is sent to the SQL server to be processed. However, since the data is stored as relational data and LINQ works with data encapsulated in objects, it must be defined a mapping between the two representations. For this reason, LINQ to SQL also defines the mapping framework [1].

The mapping is done by defining classes that correspond to tables in the database and the database relations specified by the primary keys are defined using LINQ to SQL attributes.

The Visual Studio IDE includes a mapping designer that automatically creates the corresponding classes from a database schema.

The mapping is implemented by the *DataContext* class [9] (inherit from the core of the LINQ to SQL) that manages the conversion of the LINQ query into T-SQL and retrieves the result set from the database server.

Since the processing work happens at the database server, local methods, which are not defined as a part of the lambda expressions representing the predicates, cannot be used. However, stored procedures on the server can be used.

```
contactsDataContext db = new contactsDataContext();
var q = from c in db.Contact
          where c.contactId < 4
          orderby c.firstName descending
          select c;
foreach(var c in q)
  Console.WriteLine("{0} {1}",c.FirstName, c.LastName);
```

LINQ to SQL provides a nice and clean way to model the data layer of an application. Once defined the data model it is easy to perform queries, inserts, updates and deletes against it.

## 2.4 LINQ to DataSet

The DataSet is one of the most widely used components in ADO.Net and is a key element of the disconnected programming model that ADO.Net is built on.

The LINQ to DataSet extend the limited query capabilities of the DataSets using the same query functionality that is available for the other providers. It is built on and uses the existing ADO.Net architecture, without attempting to replace it [10].

To query a DataSet it is needed that it contains data. There are several ways to fill a DataSet like using a *DataAdapter* or use LINQ to SQL.

### 2.4.1 Querying Typed DataSets

In order to enable LINQ queries over typed DataSets, the tables contained in the DataSets are represented using classes inherited from a new class named *System.Data.TypedTableBase<T>*.

The new TypedTableBase<T> class still inherits from DataTable, so it is an enhancement that does not break backward compatibility.

As is displayed in Figure 3, this class adds to DataTable implementations of *System.Collections.Generic.IEnumerable<T>* and *System.Collections.IEnumerable*.

This converts DataTables into collections, and thus allows LINQ to query the tables contained in a typed DataSet.



**Fig. 3.** TypedTableBase<T> class extends the DataTable

```
var dataSet = new DataSet();
FillDataSetUsingLinqToSql(dataSet);
var filteredBooks = from book in dataSet.Book
                    where book.Title.StartsWith("L")
                    select new { book.Title,
                                 book.Price };
dataGridView1.DataSource = filteredBooks.ToList();
```

### 2.4.2 Querying untyped DataSets

With a populated the DataSet, LINQ can be used to query information using the same features and syntax used for other data stores.

The DataSet's *DataTable* and the DataTable's *DataRowCollection* do not implement *IEnumerable<DataRow>* so these types cannot be used as sources for LINQ query expressions. In order to solve this problem, an extension method called *AsEnumerable* was added to the *DataTable* type. This extension method takes a

source *DataTable* and wraps it in an object of the type *EnumerableRowCollection* which implements *IEnumerable<DataRow>*. This allows *DataTables* to be a source for LINQ query expressions.

```
var dataSet = new DataSet();
FillDataSetUsingLinqToSql(dataSet);
DataTable bookTable = dataSet.Tables[1];
var filteredBooks =
  from book in bookTable.AsEnumerable()
  where book.Field<String>("Title").StartsWith("L")
  select new {
    Title = book.Field<String>("Title"),
    Price = book.Field<Decimal?>("Price")
  };
dataGridView1.DataSource = filteredBooks.ToList();
```

### 2.4.3  LINQ to Entities

LINQ to Entities is a LINQ implementation that queries objects created by the *ADO.NET Entity Framework* which provides an *Entity Data Model* (EDM) and services that help programmers define and interact with data at a more conceptual level [11]. The EDM can be used to model the data of a particular domain so that applications can interact with data as entities or objects. This makes the object layer an ideal target for LINQ support.

LINQ to Entities enables developers to write queries against the database from the same language used to build the business logic.

The Table 1 provides a comparison of features between LINQ to SQL and LINQ to Entities.

**Table 1.** Features comparison between LINQ to SQL and LINQ to Entities

| Feature | LINQ to SQL | LINQ to Entities |
|---|---|---|
| Language Extensions Support | Yes | Yes |
| Language Integrated Database Queries | Yes | Yes |
| Many-to-Many (3way Join/Payload relationship) | no | No |
| Many-to-Many (No payload) | no | Yes |
| Stored Procedures | Yes | No |
| Entity Inheritance | No | Yes |
| Single Entity From Multiple Tables | No | Yes |
| Identity Management / CRUD features | Yes | Yes |

In October 2009, Microsoft made a statement where it is said that Microsoft is making significant investments in the Entity framework, and therefore this will be the recommended data access solution for LINQ to relational scenarios after .Net Framework 4.0 launch. The developers community has seen this bulletin as the announcement of LINQ to SQL death [12].
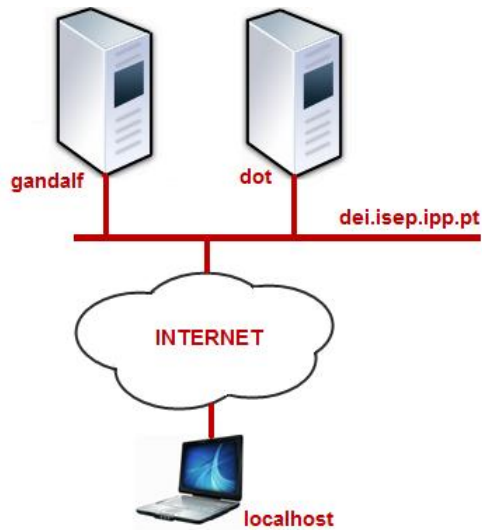
# 3     Performance evaluation test

The goal of this article is to compare the performance of the ADO.net, LINQ to SQL and Stored Procedures using the four basic database operations (SELECT, INSERT, UPDATE and DELETE).

To achieve this, it was necessary to build a test scenario that would allow reliable results. The architecture of this scenario consists of a database management system and two web servers with the features referred in Table 2.

**Table 2.** System features

| Feature | DBMS | Web Server 1 | Web Server2 |
|---------|------|--------------|-------------|
| Address | gandalf.dei.isep.ipp.pt | dot.dei.isep.ipp.pt | localhost |
| Server | Ms. Sql Server 2008 | IIS 6.0 | IIS 6.0 |
| OS | Win. 2003 Server | Win. 2003 Server | Win. Vista Business |
| Memory | 2 GB | 1GB | 4GB |

The Figure 4 shows the scenario used to perform the tests.



**Fig. 4.** The performance test system architecture

The tests will evaluate the execution time of the chosen commands against the database using each technology. That command will be repeated one thousand times and the average of the evaluated times will be chosen to represent the execution time of that command.

These bulky tests were repeated thrice at different stages of the day in order to test different loading condition of servers and network.

The client application used to query the database is an ASP.net web application developed over .Net Framework v3.5. This application was deployed into two web servers, one of them in the same private network as the DBMS (I will refer to this as

the local server) and the other directly connected to the network (hereafter will be identified as remote server). The main objective is to verify if the behavior identified in the remote server also applies to the local server.

The database commands used in the test will be the same (in functionality) despite of the technology used (ADO.Net, LINQ to .Net or Stored Procedures).

### 3.1 The Database

The database used in this test has a single table with three fields as is show in Figure 5.



**Fig. 5.** Database table PaperMIC

### 3.2 The web application

This application was developed in ASP.net and is intended to evaluate and display the commands execution time in the database using each of the three technologies.

```
DataSet ds=null;
double sum = 0;
for (int x = 0; x < 1000; x++){
  SQLDal dalObj = new SQLDal();
  DateTime tInic = System.DateTime.Now;
  ds = dalObj.getCmdSQL();
  TimeSpan tExec = System.DateTime.Now.Subtract(tInic);
  sum += tExec.TotalMilliseconds;
}
sum = sum / 1000;
```

### 3.3 The SELECT Command

The task we pretend to execute is to select from the database the first five records with lowest value in *IdField*. All the tests of this command were made with five thousands records in the database.

```
Select top 5 * from PaperMic order by IdField;
```

The LINQ solution for this task was implemented with LINQ to Datasets, and the initial dataset was filled with LINQ to SQL.

```
public DataSet getCmdLINQ(){
  DataTable tableDs = dsDb.Tables[0];

  var PaperRecs =(from rec in tableDs.AsEnumerable()
                  orderby rec.Field<int>("IdField")
                  select rec).Take(5);
  DataSet ds = new DataSet();
  ds.Tables.Add(PaperRecs.CopyToDataTable());
  return ds;
}
```
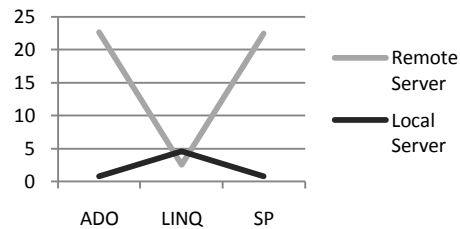
The results are displayed in Table 3, where the values are averages of bulky tests results.

**Table 3.** Results of the SELECT command test

| Server | ADO | LINQ | SP |
|---|---|---|---|
| Remote Server | 22,6517 | 2,5590 | 22,4548 |
| Local Server | 0,7656 | 4,5844 | 0,7604 |

values in milliseconds



### 3.4   The INSERT Command

This task pretends to insert a new record in the database with a *date/time* stamp and information of the used technology.
The LINQ version uses LINQ to SQL to perform the task.

```
public bool InsertCmdLINQ(){
  var newPaperMIC = new PaperMIC{
    DateField = System.DateTime.Now,
    TextField = "LINQ"
    };
  paper.PaperMICs.InsertOnSubmit(newPaperMIC);
  paper.SubmitChanges();
  return true;
}
```

The results are displayed in Table 4.

**Table 4.** Results of the INSERT command test

| Server | ADO | LINQ | SP |
|--------|-----|------|-----|
| Remote Server | 26,7032 | 79,2745 | 26,7750 |
| Local Server | 6,4896 | 6,7448 | 6,4479 |

values in milliseconds



## 3.5 The UPDATE Command

The update task will change the date/time stamp and the used technology information of the lowest *IdField*.

The SQL command to use in this task is:

```
Update PaperMIC set
  DateField=@DateField,TextField=@TextField
  where IdField =(select Min(IdField) from PaperMIC);
```

The LINQ to SQL code used to perform this task is:

```
public bool UpdateCmdLINQ(){
  var upPaper = (from r in paper.PaperMICs
                 orderby r.IdField
                 select r).First();
  upPaper.DateField = System.DateTime.Now;
  upPaper.TextField = "LINQ";
  paper.SubmitChanges();
  return true;
}
```

The test results for the UPDATE command are in Table 5.

**Table 5.** Results of the UPDATE command test

| Server | ADO | LINQ | SP |
|--------|-----|------|-----|
| Remote Server | 28,5183 | 111,1727 | 27,7373 |
| Local Server | 1,5052 | 5,5469 | 1,5572 |

values in milliseconds

### 3.6    The DELETE Command

This task pretends to delete the lowest *IdField* record.  To accomplish such task, the LINQ version of the code uses a lambda expression
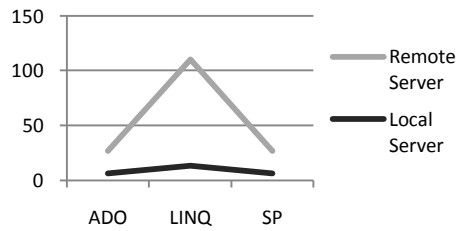
```
public bool DeleteCmdLINQ(){
  var oldPaper = paper.PaperMICs.Single(
       p => p.IdField == (
           (from r in paper.PaperMICs
            select r.IdField).Min())
       );
  paper.PaperMICs.DeleteOnSubmit(oldPaper);
  paper.SubmitChanges();
  return true;
}
```

The measured values to this command are expressed in the Table 6.

**Table 6.** Results of the DELETE command test

| Server | ADO | LINQ | SP |
|--------|------|--------|--------|
| Remote Server | 27,0038 | 110,4019 | 26,5611 |
| Local Server | 6,3906 | 13,4375 | 6,4062 |

values in milliseconds



## 4    Conclusions and Results Discussion

LINQ provides a coordinated, consistent and efficient syntax, allowing the manipulation and selection of data in the same programming language used for the development of computer applications. Therefore the objective of full integration of the data access layer in the development environment of the application is fully achieved.

However, the performance of LINQ to ADO.Net is very disappointing. Excepting the results of the SELECT command test in the remote server, where the measured times for the LINQ version were almost 9 times faster than the ADO.Net version and the Stored Procedures version, the other results of LINQ were worse than the other two technologies.

Even the best result of the LINQ version on the remote server SELECT command test cannot be completely attributed to LINQ. Actually the SELECT command was implemented with LINQ for Datasets and therefore the required query was made in a Dataset object in the web server and not directly on the Database. This also explains the worst result of the LINQ version of SELECT command in the local server. The

features of the local server (namely the physical memory) are substantially poor then the remote server.

Figure 6 shows a comparison of results between the two test scenarios.



**Fig 6.** Comparison between the two test scenarios

On the other hand, the results of the ADO.Net version and Stored Procedure version of the application were very similar. The queries were not very heavy, but the UPDATE and the DELETE commands require a search for the record with lower *IdField*, and the database always had several thousands of records.

For further improvement, other commands can be tested (like join selections and transacted operations) and the results can be analyzed against different kind of information read (like integer read, string read, block data read and substring reads).

# References

1. Box, D., Hejlsberg, A.: The LINQ Project - .NET Language Integrated Query. Microsoft Corporation (2005)
2. Marguerie, F., Eichert, S., Wooley, J.: LINQ in Action. Manning Publications Co.(2008)
3. Microsoft Corporation: Lambda Expressions (C# Programming Guide), http://msdn.microsoft.com/en-us/library/bb397687.aspx; [Cited: Jan. 2010]
4. Microsoft Corporation: Language-Integrated Query (LINQ), http://msdn.microsoft.com/en-us/library/bb397926.aspx; [Cited: Jan. 2010]
5. Meijer, E., Beckman, B., Bierman, G.: LINQ: reconciling object, relations and XML in the .NET framework. In: SIGMOD '06: Proceedings of the 2006 ACM SIGMOD international conference on Management of data. New York, NY, USA: ACM; 2006. p. 706.
6. LINQ to XSD – Overview – An incubation project on typed XML programming, Microsoft Corporation (2006)
7. Kumar, N.:LINQ Quickly, pp 155 – 169, Packt Publishing Ltd (2007)
8. Microsoft Corporation, LINQ to SQL:.NET Language-Integrated Query for Relational Data, http://msdn.microsoft.com/en-us/library/bb425822.aspx, [Cited: Jan. 2010]
9. Microsoft Corporation, DataContext Class (System.Data.Linq), http://msdn.microsoft.com/en-us/library/system.data.linq.datacontext.aspx, [Cited: Jan.2010]
10. Microsoft Corporation: LINQ over DataSet; May 2006
11. Microsoft Corporation: LINQ to Entities, http://msdn.microsoft.com/en-us/library/bb386964.aspx, [Cited: Jan. 2010]
12. Update on LINQ to SQL and LINQ to Entities Roadmap, http://blogs.msdn.com/adonet/archive/2008/10/29/update-on-linq-to-sql-and-linq-to-entities-roadmap.aspx, [Cited: Jan. 2010]

# An Evaluation of RIA Frameworks

Bruno André A. Loureiro[1]

[1] FEUP – Faculty of Engineering, University of Porto
Rua Dr. Roberto Frias, S/N 4200-465, Porto, Portugal
bruno.loureiro@fe.up.pt

**Abstract.** This paper aims to show the role of Rich Internet Applications (RIA) as a way to allow a Web application to provide the same functionality of a desktop application. For developing RIA Applications it's important to make a selection of the best framework, so, this paper will present some frameworks that support the RIA development. Frameworks will be evaluated according to some parameters. The results of this analysis enable us to perform a comparison and show the strengths and limitations of each framework.

**Keywords:** Web Applications, RIA, Frameworks.

## 1    Introduction

The Internet has emerged as the default platform for application development [1]. Web applications are often used by enterprises due to its simplicity of use, using only a Web browser, avoiding the additional installation in the operating systems. One of greatest problems in the development of Web applications is to provide the user with an environment which approximates better as possible to daily used applications, such as office applications, e-mail applications, and others.

Some time ago the development of a Web application implied problems with page refresh for example. The use of technology like AJAX[1], can solve partially this problem, but was unable to provide the same user experience as expected in a desktop application.

The arrival of RIA frameworks allowed the developer to employ in a Web application all richness of desktop applications avoiding the problems previously encountered, and with all better user experience.

In a traditional Web application the result is frequently a frustrating, confusing or disengaging user experience resulting in unhappy customers, lost sales, increased costs and the disappointment that the Internet has not lived up to its promise [1].

Given the high number of RIA frameworks, the developer is faced with a difficult task to select the best framework that allows the development of a solution, quickly and efficiently.

As each framework has strengths and limitations, the main goal of this work is to make a comparison among several RIA frameworks.

---

[1] AJAX - Asynchronous JavaScript Technology and XML

This section introduced the RIA environment. Section 2 gives a brief overview of Web applications, including a comparison between Web Applications and Desktop Applications, definition and benefits of a RIA and frameworks for RIA development. Section 3 presents the criteria for evaluation and a detailed comparison of several frameworks. Section 4 shows the results of the comparison. Finally, section 5 concludes the paper and discusses future work.

## 2 Web Applications

The number of corporate web applications has grown exponentially and most organizations are continuing to add new applications to their operations [11]. The main feature of the Web Applications is its execution in a Web Browser. Some aspects of Web applications are briefly presented.

### 2.1 Comparing Web Applications with Desktop Applications

A desktop application is software that needs to be installed in computer's operating system while a Web application needs only a Web browser. An example of a desktop application is the Microsoft Word and of a Web application is an e-commerce site.

The distinction between a Web application and a desktop application is made based on following features:

- **Maintenance**: Web based applications need to be installed only once whereas desktop applications are to be installed separately on each computer [4]. Updating a desktop application needs to be carried out in every computer which it was installed. That is not the case for Web applications.
- **Ease of use**: desktop applications are confined to a physical location and hence have usability constraint [4]. Web applications can be used by the users from any location using the Internet.
- **Security**: Web applications are exposed to more security risks than desktop applications [4]. While in a desktop application we can ensure protection from various vulnerabilities, in Web applications this is hard to attain.

### 2.2 Rich Internet Applications

The term RIA was introduced in March 2002 by vendors like Macromedia who were addressing limitations at the time in the "richness of the application interfaces, media and content, and the overall sophistication of the solutions" by introducing proprietary extensions [3].

RIA are Web applications that are considered to have many of the features and functions of traditional desktop applications [2], in desktop application the user interface has components like menus, multi-window, multi-tab, combo boxes, etc.

The distinction between traditional Web applications and RIA's, is its richness and responsive, features that were difficult to incorporate in traditional Web applications.

Among several definitions of RIA, Macromedia defines RIAs as combining the best user interface functionality of desktop software applications with the broad reach

and low-cost deployment of Web applications and the best of interactive, multimedia communication (Figure 1) [1].
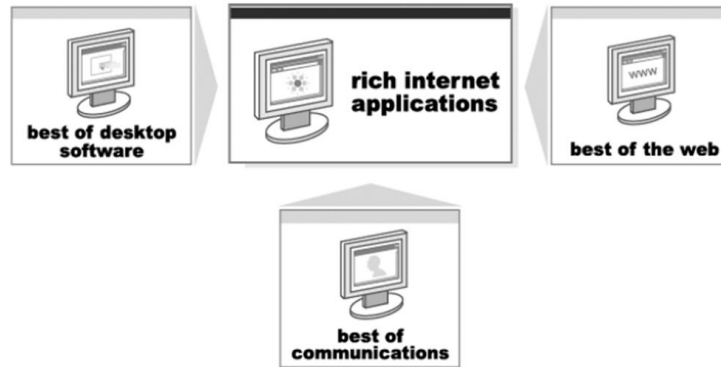


**Fig. 1.** Composition of a RIA [1]

Joining the best of desktop software, communications and the Web promotes an application more intuitive, responsive, and effective user experience.

### 2.3 Benefits of RIA

In a traditional Web application all the activity is wrapped around client/server architecture with a thin client [5]. All processing is done at the server while the client is only used to display static content.

The biggest drawback as experienced over the years with this system is that all the interaction with the application must pass through the server [5], any request is send to server and the server sends the response to client, and the page is reloaded.

RIAs on the other hand have an enhanced client side technology which can execute instructions on the client's computer [5]. This processing technology in client side is generally termed as client engine.

The client engine of a RIA makes possible the following characteristics:

- **Richer:** RIA applications provide improved UI[2] behaviors normally not obtainable with standard browser based Web-applications;
- **More responsive:** The interface behaviors are typically much more responsive than those of a traditional Web application;
- **Client/ Server Balanced:** The demand for client and server computing resources is better balanced, it allows the shared processing between client and server, while in traditional Web application a server is seen like a workhorse;
- **Asynchronous:** The client engine can interact with the server asynchronously, without requiring a user action of pressing a button or link. This option allows RIA exchange data between the client and the server, without waiting;

---

[2] UI **-** User Interface

**- Network Efficient:** The network traffic is significantly reduced because an application-specific client engine is more intelligent than a traditional Web application and has the capability to decide what/which data needs to be exchanged with servers.

### 2.4 Frameworks for RIA Development

To begin the development of a RIA's application it's necessary to select a framework among the several frameworks available on the market. The software framework is typically responsible for downloading, updating, verifying and executing the RIA [5]. Examples of RIA frameworks are: OpenLaszlo[3], Thinwire[4], Adobe Flex[5]/AIR[6], JavaFX[7] and Microsoft Silverlight[8].

The RIA frameworks can be grouped in tree major classes: plug-in based, browser and Java based frameworks. Table 1 shows the RIA frameworks classification.

**Table 1.** RIA frameworks classification

| RIA Frameworks | |
| --- | --- |
| **Type** | **Examples** |
| Plug-in based frameworks | Adobe Flash (with flex and AIR), Microsoft Silverlight, Sun JavaFX, OpenLaszlo, Thinwire |
| Browser based frameworks | Dojo, echo3, GWT, Yahoo UI to BackBase |
| Java based frameworks | Applet and java Webstart technologies |

## 3   RIA Frameworks Evaluation

To make an evaluation of a RIA Framework the first step was to select the evaluation criteria, so the selected parameters were: Maturity, IDE[9] Integration, UI Code, Community Involvement and Software requirements. These parameters seem to be most relevant to a beginner in this matter.

In this study, we only will evaluate the frameworks: Adobe Flex, Microsoft Silverlight, Sun JavaFX and OpenLaszlo. Was selected Adobe Flex and Microsoft Silverlight due to their high popularity, sun JavaFx due to its novelty and OpenLaszlo because was the first RIA framework that appeared.

---

[3] http://www.openlaszlo.org/

[4] http://www.thinwire.com/

[5] http://www.adobe.com/products/flex

[6] http:// www.adobe.com/products/air/

[7] http://javafx.com/

[8] http://silverlight.net/

[9] Integrated Development Environment

## 3.1 Maturity

Maturity is measured by a combination of the current version of the framework [7].

**Adobe Flex.** Development of the Flex framework was initiated by Macromedia in 2002. Adobe inherited Flex when the company bought Macromedia in 2004. The Flex SDK[10] was released to the open source community in 2006 with Flex V2. The current version is Flex V4 SDK beta 2 [11].

**Microsoft Silverlight.** The first version of Silverlight was launched in April, 2007 and enables the development of the next generation of Microsoft .NET - based media experiences and rich interactive applications (RIAs) for the Web. Silverlight is delivered as a cross-platform and cross-browser plug-in that exposes a programming framework and features that are a subset of the .NET Framework[12]. The current version is Silverlight 4.0 beta 1 [13].

**JavaFX.** JavaFX is a software platform that allows the creation of RIA that can run across a wide variety of connected devices, launched in 2008. The current release (JavaFX 1.2, June 2009) enables building applications for desktop, browser and mobile phones.

**OpenLaszlo.** OpenLaszlo is an open source platform created by Laszlo Systems Inc for the development and delivery of rich Internet applications, started in 2001 with name Laszlo Presentation Server (LPS). In the Version 3 in 2005 the name of LPS was changed to OpenLaszlo. It is released under the Open Source Initiative-certified Common Public License. The current stable version of OpenLaszlo is V4.6.1.

**The leader: OpenLaszlo.** The first RIA framework on market was OpenLaszlo, development technically began before that of Flex, even though it was under a different name, and Laszlo Systems has consistently stayed ahead of Adobe with its releases. OpenLaszlo has been able to maintain a consistent release schedule-currently in V 4.6.1 while the Flex SDK remains in the beta of its fourth version.

## 3.2 IDE Integration

IDE integration is measured according to the existence of a plug-in for Eclipse[14], one of best IDE for application development.

**Adobe Flex.** Adobe® Flex® Builder 3.0.2 Professional Eclipse Plug-in is a professional Eclipse™ based developer tool enabling intelligent coding, interactive step - through debugging, and visual design of user interface layout and behavior for Flex applications. Its greatest drawback is to be a commercial product.

---

[10] Software Development Kit

[11] http://labs.adobe.com/technologies/flex4sdk/

[12] http://www.microsoft.com/net/

[13] http://silverlight.net/getstarted/silverlight-4-beta

[14] http://www.eclipse.org

**Microsoft Silverlight.** Eclipse4SL[15] is an open source, feature-rich and professional RIA application development environment for Microsoft Silverlight in Eclipse.

**JavaFX.** JavaFX 1.2 Plugin[16] for Eclipse provides features to allow you to edit, build, and debug JavaFX applications in the Eclipse IDE. The editor features a set of drag-and-drop snippets to quickly add JavaFX objects with transformations, effects, and animation, bringing a set of examples.

**OpenLaszlo.** IDE4Laszlo[17] is a IDE for OpenLaszlo, is an Eclipse-based development environment for creating, editing, debugging, and testing applications based on the LZX declarative mark-up language.

**The leader: OpenLaszlo, JavaFX, Microsoft Silverlight.** All frameworks offer integration with IDE Eclipse, in the Adobe Flex case is required one commercial plug-in so it is a drawback in comparison with the others.

### 3.3 UI Code

The Adobe Flex, OpenLaszlo, Microsoft Silverlight frameworks implement their own declarative XML-based object declaration language. More specifically, Flex uses MXML, OpenLaszlo uses LZX and Microsoft Silverlight uses XAML. JavaFX uses JavaFx script is a new script language introduced as part of the JavaFX technology.

**Adobe Flex.** Adobe Flex uses MXML (Magic eXtensible Markup Language) as its declarative markup language, which follows the conventions of basic XML [7]. Figure 2 shows an example of MXML code.
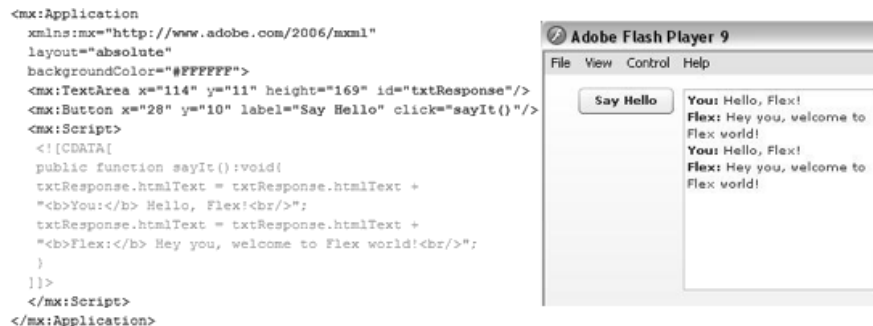


**Fig. 2.** Sample of MXML Code [8]

**Microsoft SilverLight.** Extensible Application Markup Language (XAML) was originally created for the Windows Presentation Foundation (WPF) technology released with .NET 3.0. WPF and XAML provide a way to integrate designers into the application development process [9]. Silverlight 1.0 brought XAML to the world

---

[15] http://www.eclipse4sl.org/download/

[16] http://javafx.com/docs/gettingstarted/eclipse-plugin

[17] http://www.syte.ch/en/laszlo.xml

of RIA development that can run directly in the browser once the Silverlight plug-in has been installed. An example of XAML code can be seen in Figure 3.

```
<Canvas x:Name="parentCanvas"
  xmlns="http://schemas.microsoft.com/client/2007"
  xmlns:x="http://schemas.microsoft.com/winfx/2006/xaml"
  Loaded="Page_Loaded"
| (...)>
<TextBlock Canvas.Left="130" Canvas.Top="11"
 Height="169" x:Name="txtResponse"/>
 <UIControls:Button  Canvas.Left="10"
 Canvas.Top="10" Text="Say Hello"
 Click="SayIt"/>
</Canvas>
```

Say Hello  You: Hello, Silverlight!

Silverlight: Hey You, Welcome to the Silverlight world!

You: Hello, Silverlight!

Silverlight: Hey You, Welcome to the Silverlight world!

**Fig. 3.** Example of XAML Code [8]

**JavaFX.** JavaFX Script is a scripting language is part of the JavaFX family of technologies on the Java Platform for Rich Internet Application development being quite different from others frameworks. Figure 4 illustrates an example of JavaFx code.

```
package hellojavafx;
import javafx.ui.*;
class HelloModel {
    attribute saying: String;
}
var model = HelloModel {
    saying: ""
};
Frame {
    title: "Hello, JavaFX"
    width: 500
    height: 300
    content:  Panel {
        content:
    [Button {
     x: 10 y: 10 text: "Say Hello"
    action: operation() {
    model.saying = model.saying.concat("You: Hello, JavaFX!\n");
    model.saying =
    model.saying.concat("JavaFX: Hey you, Welcome to JavaFX World!\n");
    }},
    (...)
  ] }
 visible: true}
```

Hello, JavaFX

Say Hello  You: Hello, JavaFX!

JavaFX: Hey you, Welcome to JavaFX World!

You: Hello, JavaFX!

JavaFX: Hey you, Welcome to JavaFX World!

You: Hello, JavaFX!

JavaFX: Hey you, Welcome to JavaFX World!

**Fig. 4.** Example of JavaFX Code [8]

**OpenLaszlo.** Remember that OpenLaszlo uses LZX as its declarative markup language, which also follows the same basic principles as any well-formed XML.

With LZX, you can place a JavaScript function inside a <method /> tag and assign that method its own ID [7]. Figure 5 shows an example of LZX Code.
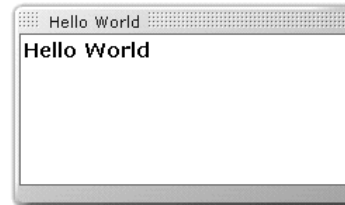
```
<canvas width="800" height="600">
<window name="mainwindow" width="250" height="150"
title="Hello World!">
<simplelayout axis="y" spacing="10"/>
<text text="Hello World" fontsize="14"
fontstyle="bold" />
</window>
</canvas>
```

**Fig. 5.** Example of LZX Code [10]

**The leader: Adobe Flex, Microsoft Silverlight and OpenLaszlo.** The basic principle of continued development for each framework remains the same: to do more with less. The XML based languages, like MXML, XAML and LZX are better intuitive for promote the RIA development.

JavaFX uses the JavaFx Script, language that is not based in XML, so the developer needs to have some knowledge about Java programming.

### 3.4 Community involvement

The level of community involvement is generally a key factor to its adoption, as well as its rate of evolution [7].

Community involvement is more easily recognized by the number of technical documentation produced, like tutorials, hints, and bug's database with always some quality-assurance given by experts in matter.

The measurement of community involvement was made based on the number of results obtained by Google[18] containing the name of the framework and the term tutorial. For example: flex tutorial. Table 2 shows the results of the community involvement.

**Table 2.** Measurement of the community involvement

| Framework | Search Terms | Results |
|---|---|---|
| Adobe Flex | flex tutorial | 3.620.000 |
| Microsoft Silverlight | silverlight tutorial | 1.600.000 |
| JavaFX | javafx  tutorial | 334.000 |
| OpenLaszlo | openlaszlo  tutorial | 118.000 |

**Adobe Flex.** Prior to V2, Flex had already gained a substantial amount of attention, largely because of strong marketing efforts after Adobe bought Macromedia [7]. As a result, the Flex community seemed to scale exponentially as soon as the Flex SDK was released as open source in 2006 [7].  The number of results founded was the best with 3.620.000. The Flex community continues to grow.

Some interesting sites are: http://flex.org and http://www.flexdeveloper.eu.

---

[18] http://www.google.com

**Microsoft Silverlight.** The community involvement of Silverlight developers presents one greater dimension like Adobe Flex, but a bit smallest. The number of results was 1.600.000. An example is: http://xamlpt.com.

**JavaFx.** JavaFX has a small community due to its recent appearance on the RIA market. The number of results was 334.000. An example is: http://jfx.wikia.com/wiki/Planet_JFX_Wiki.

**OpenLaszlo.** The size of the community around OpenLaszlo remains relatively small [7]. One reason for this smallest community can be associated with marketing budget of Laszlo Systems, comparatively with Adobe, Microsoft and Sun. The number of results was 118.000. An example is: http://forum.openlaszlo.org/.

**The leader: Adobe Flex.** The community involvement is important to provide resources for developers. After analysis, the conclusion arrived was that the RIA framework with the greatest community involvement is Adobe Flex followed Silverlight, both JavaFx and OpenLaszlo need some time to grow.

In addition of specific RIA Frameworks communities referenced in analysis, exist some generalist communities of RIA development that can be useful since the beginners at experts. Examples of it are: RiaPT[19], InsideRIA[20], among others.

### 3.5 Software requirements

This category refers to necessary plug-ins or software that you must install before you can run an application built with the respective framework [7]. This is huge entry barrier because users sometimes don't trust in software that need to be installed, so when is required an additional installation some users leave the application.

**Adobe Flex.** Flex applications run from Adobe Flash® Player, which currently has the highest penetration rate of all software applications and plug-ins [7]. In June 2009, worldwide penetration of the Flash Player had reached a record 98.8 percent, with 97.1-percent penetration in emerging markets [7].

**Microsoft Silverlight.** Web applications developed using Microsoft Silverlight needs a Silverlight plug-in[21] powered by the .NET framework and compatible with multiple browsers, devices and operating systems, bringing a new level of interactivity.

**JavaFX.** The framework JAVAFX needs the Java Plug-in with JavaFX[22] extension to run on a Web browser. This plug-in has the largest size relatively of the others.

**OpenLaszlo.** OpenLaszlo compiles to both SWF[23] files which play from Flash Player and DHTML[24]. The timeline of OpenLaszlo development shows that the intention of

---

[19] http://www.riapt.org

[20] http://www.insideria.com

[21] http://www.microsoft.com/silverlight/

[22] http://java.sun.com/javase/downloads

[23] ShockWave Flash

[24] Dymanic Hypertext Markup Language

this framework is to eventually be capable of compiling to any format which runs in a Web browser, including Java applets.

**The leader: OpenLaszlo and Adobe Flex.** The fact that OpenLaszlo produce more than one file format, makes a clear distinction in relation of others RIA frameworks, assuring in this way a greater compatibility. However the strong presence of the Adobe Flash Player in the computers, causes Adobe Flex is a good choice, too.

## 4    Results

After the comparison according several parameters, this enables us to identify the strengths and limitations of each framework. Table 2 summarizes our evaluation.

**Table 3.** Summary of RIA frameworks evaluation

| Evaluation Parameters | RIA Frameworks | | | | |
| | Adobe Flex 1 | Microsoft Silverlight 2 | JavaFx 3 | OpenLaszlo 4 | Lider |
|---|---|---|---|---|---|
| **Maturity** | 4.0 beta 2 | 4.0 | 1.2 | 4.6.1 | 4 |
| **UI Code** | MXML | XAML | JavaFX Script | LZX | 1,2,4 |
| **IDE Integration (Eclipse)** | Adobe® Flex® Builder 3.0.2 Prof. | Eclipse4SL | JavaFX 1.2 Plugin | IDE4Laszlo | 2,3,4 |
| **Community involvement** | Very Good | Good | Sufficient | Weak | 1 |
| **Software Requirements** | Adobe Flash Player | Silverlight Plug-in | Java Plugin with JavaFX extension | Adobe Flash Player or DHTML | 1,4 |

As can be seen in the Table 3, the OpenLaszlo appears as lider in several parameters lacking the parameter Community Involvement.

One aspect important of OpenLaszlo is allowing the output of a SWF file or using a DHTML, it can give some contribution for a greater compatibility in systems, but losses due to its weak community, so we select the Adobe Flex and Microsoft Silverlight as better frameworks for RIA development.

Relatively about JavaFX, this framework is immature and need some substantial growth for arrive the top and achieve a position near of Adobe Flex and Microsoft Silverlight.

## 5    Conclusion

During this work was made a comparison among the several frameworks available for the RIA development, Adobe Flex, Microsoft Silverlight, JavaFX and OpenLaszlo, in

terms of Maturity, IDE integration, UI Code, Community Involvement and Software Requirements.

It can be concluded that, in the continuous struggle between Microsoft and Adobe to get atop, Adobe has some advantage, because of its maturity. The software required for running an application is Adobe Flash Player; it's available in most systems and has a very good community that will be most valuable to developers who opt for their choice.

Although the comparison among RIA frameworks could be included others parameters, it is thought that has been given a contribution in this area, allowing knowledge of the technologies available for the development of RIA.

### 5.1   Future Work

An important issue in this comparison of RIA frameworks is certainly the parameters used in realizing the comparison, other parameters can be incorporated like: code complexity management, Web services support, user experience, performance, accessibility, some of them were not included in this work due to space restrictions.

Another important aspect is the nature of this work, these technologies are always evolving, and this work can become outdated, arising the need a new comparison in future.

## References

1. Duhl, J.: Rich Internet Applications, IDC, November (2003)
2. Rhodes, J., Gonzalez, J.: Generate and publish RIA applications directly from CA Plex, ADC Austin (2008)
3. Allaire, J.: Macromedia Flash MX – A next-generation rich client, Adobe (2002)
4. Smit, J.: Desktop Applications Vs Web Applications, ArticlesBase (2008) http://www.streetdirectory.com/travel_guide/114448/programming/desktop_applications_vs_web_applications.html
5. Rich Internet Application - Development Platform, CURL, FLEX and JAVA Comparison, Sonata Software, http://www.curl.com/knowledge-center/docs/kc_1190311376.pdf
6. Grosso, W.: Laszlo - An Open Source Framework for Rich Internet Applications (2005) http://today.java.net/pub/a/today/2005/03/22/laszlo.html
7. Orlando, D.: Use the Best Open Source Client-Side Framework for Cloud Computing, IBM (2009)
8. Zhang, M.: Hello, Flex, Silverlight and JavaFX, O'Reilly (2008) http://www.insideria.com/2008/02/hello-flex-silverlight-and-jav.html
9. Wahlin, D.: Getting Started with XAML in Silverlight, M-DEV (2008)
10. Pillai, S.: Introducing OpenLaszlo ,O'Reilly (2006) http://www.xml.com/pub/a/2006/10/11/introducing-open-laszlo.html?page=3
11. Web Application Security: GamaSec Solution, GamaSec (2006)

# Developing Dynamic Reports using Rich Internet Applications

Luís Matias

Faculdade de Engenharia da Universidade do Porto
Rua Dr. Roberto Frias, s/n 4200-465 Porto – Portugal
luis.matias@fe.up.pt

**Abstract**. In the business world, there are several software tools to generate reports automatically. We can do it with software like Crystal Reports and we can use different options to configure the reports. However, we observed that it is not flexible in several points like the structure of the data base or the changing of report's parameters after its generation. To solve this problem, we propose the using of RIA – Rich Internet Applications in order to build dynamic reports. These dynamic reports will allow us to change the report's contents and parameters after its first generation. With this, it will be possible to ask for different data in runtime, with no need for generating new reports. To implement this solution using RIAs, we tested and compared two different technologies: Microsoft Silverlight and Adobe Flash. We expected that the Silverlight should be more reliable than Flash in this context because it was integrated in the .NET framework. After our experiments, we concluded that the dynamic reports can be generated using RIAs. For that solution, we stated that Silverlight was better than Flash because it's easier to use and to develop with.

## 1    Introduction

In the business world, over the last years, we witnessed the creation of several tools destined to support the automatic generation of reports from different sources and data types. These reports usually have rigid models that are independent from its data source if it has a previously defined structure.

One of the more capable and reliable software tool to do this type of interaction is Crystal Reports [1]. However, this type of applications has, due to its nature, high level restrictions to the information retrieval process based in the generated reports. It is possible to resume big data set in small lines, like annual invoices. In the opposite way, it is not allowed to do the inverse and see what lines fulfill that sum or change, in a runtime environment, the timeline that we want to report and analyze.

To fulfill this type of requests, it was proposed to develop prototypes of a new type of reports: dynamic reports. The dynamic reports are reports that are automatically generated and that allow analysis and changing with no need to do a new generation of it. We will be able to do it in runtime. To solve this problem, we used RIA's – Rich Internet Applications.

In this context, we tested and compared the two dominant technologies that implement RIA's nowadays: the Microsoft Silverlight and the Adobe Flash [8, 9, 10, 11], in a way to understand which one of these technologies would be the best solution to our problem.

We hope to conclude that it is possible to develop applications to support automatic generation of dynamic reports with the two chosen technologies, especially with Microsoft Silverlight.

This paper is structured as follows: Section no.2, titled "Automatic Report's Generation", states a general perspective about this problem. Section 3 has the title "Dynamic Reports with RIA's" and it states, in a general way, the RIA's basic characteristics and the main requirements to build a prototype created with this type of technology, based in Management Control. The Management Control basic functions and characteristics are also mentioned. In the section 4, the results obtained by the developing of the two prototypes are lightly described, considering the requirements previously defined. Section 5, titled "Comparative Study between RIA's technologies", exhaustively describes the results obtained through the analysis of the developing and the using of the prototypes. The sixth section of this paper presents a small discussion about the results obtained in the previous section. Finally, the section 7, "Conclusions", presents the main conclusions of the work.

## 2 Automatic Report's Generation

The automatic generated reports are commonly used since the early 90s, with the reports associated to the starters' software-based database management system for non-professional use: the Office and the Lotus Suite [2, 3, 4]. However, it only became a standard feature when the Crystal Reports was introduced in the market.

The Crystal Reports is a Business Intelligence application. This type of applications is commonly used in the collection, integration, analysis and presentation of information to support the decision. Originally developed for Business Objects, this application is used to develop and generate reports with several types of data sources. A large number of applications use OEM (Original Equipment Manufacturer) versions of Crystal Reports as its own tool of reporting generation.

The reports generated using Crystal Reports allow multi-applications and multi-database. They can also be used in the conception of those applications.

However, the flexibility ends here: the Crystal Reports needs the data structure to be exactly the same as the first data base used to generate it. Despite the fact that we can configure the reports with different data sources, the report can never be changed after its generation. If you consider that the time needed to generate is as big as the data source, this must be considered as a serious problem to solve when we want to, for instance, compare reports related to different timelines.

# 3    Dynamic Reports with RIA

The dynamic reports are reports that can be changed after its generation, like a dynamic query with some parameters (i.e. the beginning and the end date of a period). These reports are generated from a data set and they can be manipulated after its generation to demonstrate different perspectives. They must be portable and easily accessed, like an web portal, to be used in different situations. These characteristics present one clear requirement: we need to implement a web application to implement a dynamic report as it is described above.

However, it is not the best solution to use a traditional web application. As it is stated by Christodoulou Stygaras in [20], these applications present several problems with the processes, the configuration and the data. This last one is particularly important, because they do not support interactive explorations of the data, compelling the user to navigate the hypertext to see the desired data [21].

As stated in [21], the traditional web applications were extended in several directions to improve interactivity and ease of use. The RIA followed one of these directions. It is the mix between the interactive and the multimedia user interface functionality of the Desktop applications and between the portability offered by the Web tools, allowing to build Rich applications with data and also with multimedia contents, like dynamic charts and media clips.

Next, it is explained how the RIA can be a solution to implement Dynamic Reports and how we will use it.

## 3.1    RIA as a solution

The problem that we have is RIA related. There are several technologies available to implement a RIA. In this case, we compared two different technologies that were able to develop the same type of applications to use in this context.

The RIAs are web applications that are usually executed under browser's plug-ins. We described, in some bullet points, the main characteristics of RIAs:

- They support the graphical render and the inclusion of media clips – video and audio;
- They are easily installed because they are executed in a plug-in that adapts itself to the technological reality it finds (browser, operative system, etc.). This characteristic allows it to maintain the user experience with every platform used to execute it. The maintenance of these applications is simple because their plug-ins are automatically updated;
- They appear as more secure applications because they are usually pre-compiled and they sometimes use the sandbox's  concept to the platform they are executed over. This will limit their access to the client.
- They have a better performance because they shrink the existing latency, comparing to the traditional web applications that need to constantly connect themselves to the server to allow it to process the data.

In order to develop an application to automatically generate dynamic reports, we chose to develop a prototype focused on Management Control.

## 3.2 Management Control Requirements

The Management Control is the discipline that studies the impact of the Strategic Management on the organizations based in metrics that evaluate their performance [6].

Using tangible values to measure intangible characteristics of the organization, like its branding or its CRM (client relationship management), it is possible to evaluate the current and future impact of a strategy on the operational results. We can easily verify, for instance, if the organization profit has grown like the expected and if the budget defined to certain activities is enough to achieve the defined goals.

We find some methodologies in the market which implement a more efficient Management Control over the organizations. One of the most important is the Balanced ScoreCard [5, 7, 19]. With this, it is simple to justify our choice - because it is highly suitable to our software market: the top managers need - usually as support to their decision processes – of large amounts of data.

An example of a map needed by the Management Control, in its financial perspective, is the balance sheet's temporary map. A balance sheet is a summary of the financial balances of a company or organization. Assets, liabilities and ownership equity are listed from a specific date, such as the end of its financial year, as you can see in the following figure.

| | | $'000 | $'000 |
|---|---|---|---|
| **ASSETS** | | | |
| **Financial Assets** | | | |
| Cash | | 49,088 | 28,713 |
| Receivables | | 498,365 | 377,026 |
| *Total Financial Assets* | | 547,453 | 405,739 |
| | | | |
| **Non-Financial Assets** | | | |
| Land and buildings | | 1,706,653 | 1,517,193 |
| Infrastructure, plant and equipment | | 94,549 | 75,184 |
| Assets held for sale | | 14,485 | 1,409 |
| Intangibles | | 8,892 | 12,706 |
| Inventories | | 16,396 | 20,467 |
| Other non-financial assets | | 21,174 | 21,320 |
| *Total Non-Financial Assets* | | 1,862,149 | 1,648,279 |
| | | | |
| **Total Assets** | | | |

**Fig1**. Example of a Balance Sheet.

## 3.3 Prototypes Requirements

After the analysis of the RIAs characteristics and of the Management Control requirements, we defined the following list of requirements to the prototypes [12]:
- To develop one or more Management Control maps in a application using RIAs;
- To develop an application with dynamic contents and dynamic configurations, in runtime environment, to the developed maps;

- Allow to edit the values existing in the reports to be used later;
- To develop mechanisms to obtain, in runtime, different perspectives of the maps;
- To develop an easy-use application with strong visual components and highly intuitive and appellative aspect.
- To allow the automatic generation of bar charts or other chart types, in runtime, with the data contained in the maps;
- To allow the generation of several maps simultaneously using different configurations and visually compare themselves using the automatic generation of charts;

These requirements were used to develop two prototypes of applications to automatically generate dynamic reports to support the Management Control. One of the prototypes was developed using Microsoft Silverlight and other with the Adobe Flash. The results are presented in the following section.

## 4   Results of the Developed Prototypes

Two prototypes of the desired application were created: one using Silverlight, other using Flash, as it is stated in the following sections.

### 4.1   Flash Prototype

This prototype was easily developed, using the Adobe Software Development Kit for Flash. It is, similarly to the other prototype, connected to a database running locally in a Microsoft SQL Server using a web service. It was not too hard to develop some animations or a user friendly interface, but it was really hard to find a way to generate charts in runtime. It is really hard to develop your own Flash controls and you can't use anything from the OS. So, the last two defined requirements were not accomplished.

### 4.2   Silverlight Prototype

Like the previous one, the prototype was easily developed. It is possible to change some parameters of the map, as shown in the Figure 2.

**Fig 2**. Screenshot of the Silverlight prototype

Using the Silverlight Toolkit, it is possible to include chart components that are easily fed using LINQ – Language Integrated Query [22]. You can see the buttons on the left to navigate through the lines used in the represented sums, you can change the budget used to feed the map or the dates used to generate it. You are able to create a new map and you can save it in a file using Silverlight Isolated Storage, as you will be able to see in the further sections. Both prototypes can be accessed by a browser.

### 4.3 Prototype Developing Results

The prototypes and the maps above prove that there are many things to improve from the old Crystal Reports maps. These prototypes are more flexible, more portable and more user-friendly than the older ones. Like that, they fill the gap found in Crystal Reports transforming the Dynamic Reports into a reality. Now, we need to know which of the technologies is best to continue this work on Dynamic Reports.

In the next section, the comparative study between the two technologies is presented, based in the prototypes previously developed and in the defined metrics used to measure them.

## 5   Comparative Study between RIA's Technologies

In the previous sections, the basic characteristics of a RIA and the prototypes developed using two different RIA technologies were explained. Considering these

characteristics and the context we want to use them, we defined the following metrics to measure and compare the RIA's technologies:

– User Experience – this metric pretends to measure the capacities of the technologies to develop a good user experience like the interactivity's levels, the graphic render, the quality of the media clips and the global performance;
– Security, Communication and Client's Access - this metric pretends to measure the security level of the technology and how can we access and change the data in the Client and to communicate with other applications;
– Installation and Execution – these metric measures defines the prerequisites needed to install and execute an application of this kind and which are the requirements of the client's application;
– Maintenance and Reliability – we measure the effort to maintain this applications working as the reliability proportioned by the technology;
– Developing Easiness – it is defined if it is – or if it is not – easy to develop applications with this technology;

These metrics were used to measure and compare the RIA's technologies: Silverlight and Flash, like you can verify in the following lines. A ranking from zero to five stars ( ) was defined for each one of the metrics to support a tangible evaluation of the technologies.

## 5.1    Adobe Flash

Using the previous section and the RIA's market analysis done in [12], the following analysis to Adobe Flash was concluded:

**User Experience.** The Adobe Flash allowed a good user experience because it has several possibilities to include different rendered animations and audio/video streaming. It has a large amount of possibilities to support the interaction and a high level of styling. The Flash assures the quality of the reproduction independently from the client's monitor resolution (vector-based).

**Security, Communication and Client's Access.** In order to communicate, Flash didn't demonstrate any functionality to allow the connection to Database Management Systems. However, it can use MSMQ – Microsoft Message Queuing [18] to communicate with other local applications.

Due to its communication limitations, Flash is really secured to execute a presentation of this type. Flash is in the market from many years till now and, with that, it developed a real trustful relationship with their users: programmers and consumers.

**Installation and Execution.** Flash is in the market from many years till now, so, it supports the majority of existing operative systems and browsers like Windows 9x, Linux, Opera, IE, Mozilla Firefox, MacOS, Safari, etc... It can be executed in a browser's plug-in or in an isolated player as a standalone application. The applications developed using Flash are short-sized and the plug-in startup is really fast.

**Maintenance and Reliability.** Flash's reliability is high due to its long years of existence. The Flash's plug-in is easily maintained because it is automatically updated. Despite the advantages of Flash to generate isolated presentations, they require a new compilation of all components every time the presentation is changed. This fact difficult the maintenance of the applications built in Flash.

**Developing Easiness.** It is not simple to develop in Flash. It only supports one programming language: the Action Script. This language requires a high learning curve because it is not a standard programming language and it is not used in other applications besides those who are developed with Flash. Therefore, the source code can't be used again in other context.

It is not possible to use any OS controls and the animations are built using the frames that are shown. Flash only uses the matrical transformations to animate its presentations - we must somehow assure that the client will maintain the application's frame rate. If this doesn't happen, our animation can last for 1 or for 5 seconds (!), for instance. Flash's developing environment was designed, as well as its programming language, merely to develop Flash applications and it requires, once more, a long time to learn. This developing environment is graphical and, due to that fact, it is more suitable to use for designers than for programmers.

### 5.2 Microsoft Silverlight

Now, using the previous section and the RIA's market analysis done in [12], the following analysis to Microsoft Silverlight was concluded:

**User Experience.** The Silverlight has a set of capabilities similar to Flash. Rendered animations, streaming of audio and video, etc...

However, it's global performance it greater than Flash's. The compilation of its applications is made using CLR - Common Language Runtime [17]. Every language that uses this type of environment is managed code's languages: languages that are compiled first in a virtual machine, then, executed by the CPU.

On the other hand, Silverlight presents a technology that distinguishes itself from other technologies: the Deep Zoom. The Deep Zoom allows us to randomly zoom in large images with an enormous performance [16]. A good example of the use of this technology is the Hard Rock Memorabilia [14]. The Silverlight possesses one more unique characteristic: graphic acceleration supported by hardware, with the technology Direct3D. Because of these characteristics, we realized that Silverlight is

better than Flash in the perspective of user experience. Silverlight must probably be RIA's technology with the best user experience actually in the market.

**Security, Communication and Client's Access.** Like Flash, Silverlight hasn't got any functionality to allow the direct connection to Database management systems. However, it is not allowed to access the client file system. The simplest way to turn over this problem is to use a web service. The .NET framework and CLR allow the programmer to use technologies like WCF, ADO.Net and link to access to this kind of systems.

As far as security is concerned, we can define Silverlight like a pretty secure platform. This fact occurs because Silverlight is executed in a sandbox. A sandbox typically restricts the access to the platform's native API, controlling the resources that the application can and cannot use like the disk space and the memory space, the access to system's information, the read of input devices, etc... To minimize those effects, the Silverlight uses the Isolated Storage [15]. The Isolated Storage is a technology that stores data from web applications outside of web browser's cache. This maintains the data outside of the client even when the browser's cache is "clean". One of the Silverlight's biggest limitations is its need to communicate asynchronously [13]. Every Silverlight's communication process with web services is done through asynchronous calls.

**Installation and Execution.** The Silverlight is also executed in a browser's plug-in but it does not support many standard platforms like the Linux operative system or the Opera browser. It does not support any execution outside of a browser like a standalone application. Like Flash, the plug-in startup is fast but their files have a bigger size than those of Flash. As we observed, they are, in average, 10 to 20 times bigger as they don't compress any of its source files. Beside the fact that it needs several source code files to its execution, it's dependencies are not built-in in the application (like the bitmap and jpeg files, icons or audio/video files).

**Maintenance and Reliability.** Like Flash, the plug-in has an easy maintenance because it is automatically updated. Despite its short time in the market, Silverlight uses technologies largely used to build business oriented applications (.NET,C#,VB .NET, CLR, etc.). Consequently, the programmers state a high reliability to Silverlight. Silverlight has different source files and because of that the files that are responsible for the communications protocols, data queries, design, etc. are easily identified. With that, it's easy to create new applications and functionalities using tested and existing source code (C# code, for instance).

**Developing Easiness.** Silverlight mainly uses technologies from .NET framework in its applications like CSharp, VB, WPF, WCF, ASP.NET and LINQ. These languages are commercially used to develop applications with all type of business plans. Therefore, its source code can be used over and over again.

The used IDE (Integrated Development Environment) is also well known – Visual Studio – and the most programmers have a great familiarity with its use. Silverlight has the Expression Blend as IDE too. This tool allows us to edit our applications

design and it is pointed to the designers themselves. The creation of Silverlight's animations is, again, very easy: it is possible to define time-based animations (otherwise, it is not in Flash, where they are frame-based). We can merely define the first and the last state and the render software generates the remaining states throughout the time. The Silverlight use the Windows Presentation Foundation framework as platform to use several Windows's controls without any integration efforts. Another interesting factor is the existence of the Silverlight Toolkit. The Silverlight Toolkit is a skin and controls are ready to use in Silverlight applications. One of those controls is especially relevant: it's a graphic generator with several presentation forms which is highly suitable to use in business oriented applications like our own.

However, this technology has negative points too. The projects debugging is simple because it is possible to edit just a small component of the project without compiling it all. But we need to consider the slow startup of the application because it needs an ASP.NET server to be running, located locally or abroad.

### 5.3 Summary

Using the described work, the prototypes and the analysis previously done, we present in Table 1 a summary of the whole study.

**Table 1.** Comparative study between Flash and Silverlight

|  | Adobe Flash | MS Silverlight |
|---|---|---|
| User Experience | **** | ***** |
| Security, Communication and Client's Access | *** | *** |
| Installation and Execution | ***** | *** |
| Maintenance and Reliability | **** | **** |
| Developing Easiness | ** | **** |

## 6    Discussion

As we can see through the table analysis, in 25 possible points, the technologies could not achieve grades above 80%.

The Silverlight had one point more than Flash. This is relevant but, in our opinion, is far from being decisive. In a technological perspective, the Flash's market share is high and this makes Flash the best option. However, observing the unique characteristics of both technologies, we can say that Silverlight has bigger potential to grow comparing with Flash and it will became a really serious competitor in a short time. The Flash is a tested and commercial format with great portability e adaptability to different platforms and its executables are short sized. These characteristics should maintain it as market leader in the near future.

However, in a market with a constant expansion as the RIA's one, there are characteristics that will define how will the software applications be in the future – the use of those applications in mobile devices. In this aspect, the Flash will have more constraints to move on with the natural market's evolution. In our opinion, those characteristics are:

- 3D Hardware graphic acceleration;
- Time-based animations;
- Programming languages easy to understand and usable in different contexts.

Observing these characteristics, we can say that Silverlight presents itself in the front line to succeed Flash in the RIA's market. If Flash does not evolve itself in this way, it could become obsolete in a medium term period.

In the specific context of developing a tool to automatically generate dynamic reports, both technologies prove to be sufficient to accomplish it. However, the Silverlight presents a better performance due to its developing easiness and user experience, allowing to build complex applications to automatically generate dynamic reports, in a richer and faster way.

## 7    Conclusions

Nowadays, the reports can no longer be represented by blank and white paper reports, strictly static in the way they represent information and hard to search and to extract relevant information.

The automatic generation of dynamic reports is, in our opinion, the future of this kind of information representation. The developing of this kind of tools using RIA is possible and highly profitable in the products' quality and programming perspective.

Comparing the technologies used to implement these solutions, the Microsoft Silverlight demonstrated to be the best solution for, essentially, two reasons: on one hand, its user experience that transforms the views and the editions of the reports in simple and pleasant operations. On the other hand, the developing easiness allowed by the using of the .Net framework is high because the Silverlight uses technologies belonging to that framework, well known and largely used in the whole world like CSharp.

As future work, we will develop a multi-database architecture that will allow us to develop, in a medium term future, a commercial product to automatically generate the dynamic reports supporting one, or more, platforms.

## References

1. Peck, George. Crystal Reports 2008: The Complete Reference. McGraw-Hill Professional, 2008.
2. IBM. liveSite - Website Content Management Software. Obtained in 1st January of 2010. http://www.lotusmuseum.com/

3. Microsoft. microsoft's timeline from 1991 - 2008. http://www.thocp.net/companies/microsoft/microsoft_company_part2.htm (Obtained in 1st January of 2010).
4. The Microsoft Office Fluent user interface overview. http://office.microsoft.com/en-us/products/HA101679411033.aspx (Obtained in 1st January of 2010).
5. Kaplan, R. S., & Norton, D. P. (1996). The Balanced Scorecard. Harvard Business School Press.
6. Rodrigues, J. A. (s.d.). Controlo de Gestão e Performance. Obtained in February 16th of 2009, in Controlo de Gestão e Performance: http://www.indeg.org/cursos/mestradosexecutivos/contabilidade/controlo/
7. Kaplan, R. S., & Norton, D. P. (2004). Strategy Maps - Converting Intangible Assets Into Tangible Outcomes. Harvard Business School Press.
8. Ezell, J. (3rd May of 2007). Silverlight vs. Flash: The Developer Story. Obtained in 19th February of 2009, in Silverlight vs. Flash: The Developer Story: http://weblogs.asp.net/jezell/archive/2007/05/03/silverlight-vs-flash-the-developer-story.aspx
9. Microsoft. (11th October of 2008). The Official Microsoft Silverlight Site. Obtained in 24th October of 2009, de The Official Microsoft Silverlight Site: http://silverlight.net/
10. Advance Flash. (2009). Stock Flash. Obtained in 25th February of 2009, in Advance Flash: http://www.advanceflash.com/
11. Adobe. (7th August of 2008). Adobe - Flash Lite: Architecture. Obtained in 16th February of 2009, in Adobe - Flash Lite: Architecture: http://www.adobe.com/products/flashlite/architecture/
12. Moritz, F. (January of 2009). Rich Internet Applications (RIA). Obtained in 5th January of 2009, in Rich Internet Applications (RIA): http://www.flomedia.de/diploma/documents/DiplomaThesisFlorianMoritz.pdf
13. Microsoft. (2008). Data Points: Service Driven Apps With Silverlight 2 and WCF. Obtained in 17th February of 2009, in MSDN: http://msdn.microsoft.com/en-us/magazine/cc794260.aspx
14. Hard Rock. (2008). Hard Rock Memorabilia. Obtained in 5th January of 2009, in Hard Rock: http://memorabilia.hardrock.com/
15. Microsoft. (2008). Isolated Storage. Obtained in 17th February of 2009, in MSDN: http://msdn.microsoft.com/en-us/library/3ak841sy(VS.80).aspx
16. Microsoft. (November 2008). Deep Zoom. Obtained in 17th February of 2009, in MSDN: http://msdn.microsoft.com/en-us/library/cc645050(VS.95).aspx
17. Microsoft. (2008). Common Language Runtime Overview. Obtained in 22th January of 2009, from MSDN: http://msdn.microsoft.com/en-us/library/ddk909ch(vs.71).aspx
18. Microsoft. (s.d.). Microsoft Message Queuing. Obtained in 16th February of 2009, from: http://www.microsoft.com/windowsserver2003/technologies/msmq/default.mspx
19. Kaplan, R. S., & Norton, D. P. (2008). The Execution Premium - Linking Strategy to Operations for Competitive Advantage. Harvard Business School Press.
20. Christodoulou S. P., Styliaras G. D. and Papatheodorou T. S., "Evaluation of Hypermedia Application Development and Management Systems". 9th ACM conference on Hypertext and Hypermedia, ACM Press, Pittsburgh, 1998, pp. 1 - 10, ISBN:0-89791-972-6.
21. J.C. Preciado, M.Linaje, F.Sanchez and S.Comai, "Necessity of methodologies to model Rich Internet Applications". WSE - Proceedings of the Seventh IEEE International Symposium on Web Site Evolution, IEEE Computer Society, Washinton DC, 2005, pp. 7 - 13, ISBN ~ ISSN:1550-4441 , 0-7695-2470-2
22. Microsoft. (March 2008). Swiss MSDN Team Blog : Silverlight 2 Beta 1 + WCF + LINQ to SQL. Obtained in 5th January of 2009, de MSDN: http://blogs.msdn.com/swiss_dpe_team/archive/2008/03/17/silverlight-2-beta1-wcf-linq-to-sql-a-powerfull-combination.aspx

# Learning Environments

# A Survey on Serious Games for Rehabilitation

Paula Rego[1,2], Pedro Miguel Moreira[1] and Luís Paulo Reis[2,3]

[1] Escola Superior de Tecnologia e Gestão do Instituto Politécnico de Viana do Castelo
Avenida do Atlântico s/n, 4900-348 Viana do Castelo, Portugal
[2] Faculdade de Engenharia da Universidade do Porto
Rua Dr. Roberto Frias s/n, 4200-465, Porto, Portugal
[3] Laboratório de Inteligência Artificial e Ciência de Computadores
Rua do Campo Alegre 1021, 4169-007, Porto, Portugal
{paularego,pmoreira}@estg.ipvc.pt, lpreis@fe.up.pt

**Abstract.** Serious Games are growing into a significant area spurred by the growth in the use of video games and of new methods for their development. They have important applications in several distinct areas such as: military, health, government, and education. The design of computer games can offer valuable contributions to develop effective games in the rehabilitation area. This paper presents fundamental concepts relating to Serious Games followed by a survey of relevant work and applications on Serious Games for Rehabilitation. We propose a classification designed to properly distinguish and compare Serious Games for Rehabilitation systems in what concerns their fundamental characteristics. We also describe a particular Serious Game for Rehabilitation, RehaCom, as a case study. Finally, the paper presents some challenges and research opportunities in this area.

**Keywords:** Serious Games, Rehabilitation, Game Design, Health Informatics.

## 1  Introduction

Serious games are growing into a significant research area spurred by the advances in game development and in computer graphics hardware, in turn driven by the success of video games. They are becoming so popular that they are acquiring new audience's categories such as casual players. At the same time traditional games are becoming too complex for all but the most hard-core players in the industry. An example of this shift in audience is Nintendo Wii system, which has brought a new audience for gaming among families, women and older people, who had been turned off for the traditional games that before dominated the market.

As the success and proliferation of video games grows, they have the potential to be more than just entertainment, just like books, movies and television. According to Michael & Chen [1], the time has come for video games to become more relevant, more responsible, and more important or, in other words, to get serious. As a result there has been a trend towards the development of more complex games, which have both pedagogical and fun (game-like) elements.

The amount of research in this relatively new field has grown significantly during the past decade. All around the world a growing number of seminars and conferences are arranged. In 2002, was formed the Serious Games Initiative [2] that helps the area of serious games emerge into an organized industry of developers aiming to solve problems in diverse areas and informs us about the potential of games and how to merge innovation and developers from one discipline with those in another . In 2009 was organized the first conference specialized in serious games: VS-GAMES'09 – First IEEE International Conference in Games and Virtual Worlds for Serious Applications, in the World's premier Serious Games Institute in UK.

Although, there are many definitions of the term, it is agreed by the different authors that the tern refers to the use of computer games that have a main purpose that is not pure entertainment. In fact, Serious Games have been applied in many diverse areas: corporate and military training [3], health[4,5,6], education[7,8,9,10,11,12], cultural training[4,10,13,14,15,16]. Many of these areas are related and sometimes overlapped areas like e-learning, edutainment, game-based learning and digital game-based learning, but all these can be included in a broader concept that is "Serious Games".

This paper focuses on Serious Games in the rehabilitation area. Rehabilitation is a dynamic process of planned adaptive change in lifestyle in response to unplanned change imposed on the individual by disease or traumatic incident [18]. High social costs result in a major part from high costs in the rehabilitation of a variety of deficits resulting from diseases or traumatic incidents. The success of a rehabilitation program depends on various factors: appropriate timing, patient selection, choice of rehabilitation program, continued medical management and appropriate discharge planning. This can be achieved in a multidisciplinary way and with a appropriately equipped rehabilitation department, where medical, nursing, physical therapy, occupational therapy, speech and language therapy, clinical psychology and social work personnel work together towards a common goal in a planned and coordinated way[18].

It has been evidenced that games contribute to increase motivation in rehabilitation sessions, which is the major problem in therapy sessions, caused by the repetitive nature of exercises. We believe that the design of computer games can offer valuable contributions about how to develop more effective games for rehabilitation programs. In that sense, one of the fundamental goals of our research is to identify, classify and assess game features that are relevant for the design of computer games in this area. The work herein presented reviews relevant work described in the literature. This paper also presents our developments and proposes of a classification schema towards a taxonomy based on a set of criteria for the design of more effective rehabilitation games. This classification framework enabled us to present a comparison of serious games for rehabilitation. Another important result herein presented concerns with the identification of research opportunities and open problems.

The rest of the paper is structured as follows. Section 2 presents an introduction on Serious games including fundamental terminology and concepts, classifications and examples of application. Following, Section 3, introduces the motivation of Serious Games for Rehabilitation. Section 4 introduces our proposed classification criteria towards a taxonomy for Rehabilitation Serious Games and the following Section 5 presents a review of Serious Games for Rehabilitation. In Section 6, a reference

system - Rehacom System - is described. Finally, in Section 7, major conclusions are drawn and directions for future work are suggested.


## 2  Serious Games

Today, the term serious game is becoming more and more popular, but there is no current single definition of the concept. Zyda [19] defines a serious game as: "a mental contest, played with a computer in accordance with specific rules, which uses entertainment to further government or corporate training, education, health, public policy, and strategic communication objectives." Furthermore, Zyda [19] argues that serious games have more than just story, art, and software. It involves pedagogy (activities that educate or instruct, thereby imparting knowledge or skill) and is this what makes games serious. However, he also stresses that pedagogy must be subordinate to story and that the entertainment component comes first.

Michael and Chen [1] define serious games as "games that do not have entertainment, enjoyment or fun as their primary purpose". This is not to say that serious games are not entertaining, enjoyable, or fun, just that there is another purpose. So, as is recognized by many authors, serious game is not merely the application of games and game technology for non-entertainment purposes in such diverse domains. A game's purpose may be formulated by the user or by the game designer, which means that also a commercial off the shelf (COTS) game, used for non-entertainment purposes, may be considered a serious game.

Serious games can be of any genre, use any game technology, and be developed for any platform. They can be entertaining, but usually they teach the user something. In our work we define Serious Games as games that engage the user, and contribute to the achievement of a defined purpose other than pure entertainment (whether or not the user is consciously aware of it).

Serious games can also be classified in a number of different ways. Zyda[19] states that serious games technology can be applied to domains as diverse as healthcare, public policy, strategic communication, defense, training, and education. Michael and Chen [1] classify serious games into a number of markets: military games, government games, educational games, corporate games, healthcare games, and political, religious and art games. Susi et al. give an overview of Serious Games [20]. Despite such classifications, many games could belong to more than one category. Sawyer and Smith introduced in 2008 a Serious Games taxonomy to serve as a starting point to define serious games moving forward and invited community research to contribute to a next version of their classification.

Serious games can be applied to a broad spectrum of areas. For example, the Military area has a long history of using games for training. The first "serious game" designed and used for military training was Army Battlezone, designed by Atari in 1980. However, one of the most well known was released in 2002 – America's Army – which in contrast to most video games, is free to download.

Examples of serious games application can be found in several other diverse areas as referred before. A major application is in the rehabilitation area which is the main focus of the current research and is presented in the next section.

## 3   Serious Games for Rehabilitation

Evidence from human studies shows that goal-oriented, task-specific training improves function and that increased amounts of training produce better outcomes.

One problem with task-specific treatment approaches, however, is maintaining people's interest in performing repetitive tasks and ensuring that they complete the treatment program [21]. A lack of interest or a short attention span also can impair the potential effectiveness of therapeutic exercises.

Considering that traditional training and exercises are often boring and repetitive, using computer games to augment physical and cognitive rehabilitation offers the potential for significant therapeutic benefit. Games are likely to engage a person's attention because they require cognitive and motor activity [22]. Besides, most games challenge the player to achieve sustained success through progression to increasingly difficult levels. Another very important aspect is that games can be used to aid in the management of pain, as they divert the patient's attention [21, 22].

In the last decade there has been a growing interest in the use of Virtual Reality (VR) technology for the rehabilitation of cognitive and motor deficits. Stroke patients have become one of the main target populations for these new rehabilitative methods [21]. These VR based-methods offer the patients an opportunity to participate in experiences engaging and rewarding.

In this work we want to identify important game characteristics in the rehabilitation area. Several similar works have been reported in the literature in this area. Flores et al.[23] proposed a classification of games for elderly rehabilitation that could serve as a general indicator for the appropriateness or adaptability of each game in that area. They defined a set of criteria for game design for stroke rehabilitation, in a combination with a set of criteria for design entertainment in elderly population based on the fact that the ideal game for use in the rehabilitation of stroke patients would satisfy all the criteria defined in the two sets. They analyzed and compared two sets of games, according to these criteria: existing games currently used in stroke rehabilitation, chosen subjectively as the most entertaining for the elderly users according to their criteria, and other popular games that are not currently used in this area, but could have distinct advantages over existing games for stroke rehabilitation and thus could be adapted for this purpose. They concluded that current rehabilitation games lack entertainment qualities and popular games lack essential components for rehabilitation effectiveness.

Burke et al. [21] identified the game design principles important for upper limb stroke rehabilitation, presented several games which have been designed using these principles, and evaluated them using questionnaires made to a small number of participants. The two principles defined were: meaningful play and challenge. Meaningful play emerges from a game in the relationship between a player's actions and the system's outcome [21].

In the next Section we present a set of criteria relevant to the classification e comparison of serious games for rehabilitation. In Section 5 we present in detail a review of more relevant works found in literature, and a table comparing their characteristics.

## 4 Towards a Taxonomy for Rehabilitation Serious Games

Based on the literature reviewed we identify as important main criteria for the classification of Serious games in the rehabilitation area the following ones:

*Application area*: is the domain application in which a game can be applied; despite this domain can be very vast, we may consider however two main applications: cognitive rehabilitation (Cognitive) and physical/motor rehabilitation (Motor). In cognitive rehabilitation the goal is to achieve the most independent or highest level of functioning. The rehabilitation is based on individualized goals that take into consideration the patient's current strengths and weaknesses. Traumatic brain injury (TBI) occurs when mechanical force causes damage to brain tissue resulting in the disruption of brain functioning. Cognition is frequently damaged after TBI. Motor rehabilitation is related with many deficiencies: stroke rehabilitation (upper and lower limb extremity training, spatial and perceptual-motor training), balance training [24,31], acquired brain injury, wheelchair mobility, Parkinson's disease, orthopedic rehabilitation, functional activities of daily living training, and telerehabilitation [24].

*Interaction Technology*: the technology used by the patient to interact with the system. This can vary from the traditional methods using a mouse or keyboard process to VR based methods. For instance, in VR, patients can have visual interfaces including desktop monitors and head-mounted displays (HMDs), haptic interfaces, and real-time motion tracking devices that are used to create environments that allow users to interact with images and virtual objects in real-time through multiple sensory modalities (vision, haptics, proprioception, audition). In telerehabilitation the most used modalities are webcams, tele-videoconferencing over phone lines, videophones and webpages containing rich internet applications.

*Game interface*: the interface used in the game. It can be two-dimensional (2D) or three-dimensional (3D).

*Number of players*: number of patients playing the game: single player (single) or multi-player (multi).

*Game Genre*: the games genre can vary in relation with the technology used. Examples found include: games to evaluate the movement (catch, reach and grasp) and games that are simulations or strategy.

*Adaptability* (Yes/No): the capacity of the system to adapt dynamically game difficulty or challenge, according to the patient performance and abilities in the game. This requires that player actions have to be captured and analyzed as the game is played and game elements may change dynamically to maintain an appropriate level of challenge, making the game easier or harder derived by the user's performance. Many games use levels to structure difficulty. Usually at the start of a new game, a player generally desires a low level of challenge to meet their correspondingly low level of ability/familiarity with the game. Other games may not have recognizable levels as such, but the challenge might increase as particular points in the game are reached, again, indicating that an appropriate level of understanding and acquisition of skills has been achieved.

*Performance Feedback* (Yes/No): this dimension is related with the capacity of the system to transmit to the patient the results of the interaction. This feature gives the patient a measure of his progress in achieving goals, or in their skills over time. It can

be aural, visual and haptic and can be used to signify correct or incorrect actions or responses. Without feedback the interaction of the patient with the system looses significance due to a lack of visible meaning.

*Progress monitoring*: is the capacity of the system to allow saving the results of patients interaction with the system

*Game portability*: is related with the capacity of the system to be used at home, or at a hospital or clinic.


## 5 Review of Serious Games for Rehabilitation

Several serious games for rehabilitation have been reported in literature. In this section we review the work developed in this area, focusing in the main criteria adopted in the previous Section. We make also a reference in each game to the evaluation test done, in particular, the size of the sample used in the test and the evaluation method chosen.

Betker et al. [25] described a serious game for Balance Rehabilitation. The video game-based tool was controlled by a center–of–pressure (COP) for the maintenance of balance in a short-sitting position caused by spinal cord and head injuries. It was developed to use a pressure mat. The flexibility of the pressure mat allows games to be performed on solid, fixed surfaces and allows progression to compliant surfaces, with the pressure mat being placed between the patient and the surface. The games require movements of the patients in one or in all directions. Difficulty levels can also be configured, helping to ensure patient competitiveness while exercising his full range and speed of voluntary movement. This is important to prevent a player from becoming frustrated and quickly losing interest. The portability of the system affords its use in monitored at-home programs, which makes this therapy approach cost-effective. The games also provided the patients and the therapist with instantaneous feedback about performance and goal attainment. The games were evaluated using a questionnaire administered after the exercises and with stability measurements obtained during a set of tasks performed before and after exercise.

Ma et Bechkoum [26] described a serious-game based movement therapy which aims to encourage stroke patients with upper limb motor disorders to practice physical exercises. Their framework uses functional tasks, such as wrist extension, reaching, grasping and catching, and serious games. The system allows patients to interact with virtual objects in real-time through multiple modalities and to practice specific motor skills. Physiotherapists are necessary for initializing the system and controlling the scripting of tasks. Input devices include the ordinary devices mouse and keyboard for the operator and a range of real-time motion tracking devices—Data gloves to capture finger flex and hand postures; wireless magnetic sensors to track the patient's hand, arm and upper body movements. Output has visual, audio and haptic modalities. The dual output visual interface includes a desktop computer LCD for the operator and a high resolution HMD for patients. The HMD equipment displays an immersive virtual environment, providing a better sense of presence. The software components include a 3D graphic engine and a movement therapy module which creates functional training and non-functional serious games. It has also a dynamic adaptation module

that uses patient profiles and progress data to select tasks and initially configure the difficulty level of the tasks and games. A pilot study was made with 8 participants, concluding that serious games intervention did have an impact on the recovery of movement.

Conconi et al.[27] introduced PlayMancer, a platform for rapid development of serious games, with a special focus on therapeutic support games for behavioral and addictive disorders, i.e. eating disorders and pathological gambling. It is modular and combines techniques from multimodal interaction (speech, touch, biosensors and motion-tracking), 3D engines, virtual and augmented reality, speech recognition and natural language processing. The prototype to be adopted for chronic mental disorders (mainly eating disorders and behavioral addictions) treatment, introduces the player to an interactive scenario which aims to increase his general problem solving strategies, self-control skills and control over general impulsive behaviors. The 3D interactive environment is made up of different islands that will be used as scenario. Each island will permit access to one or several types of resources which will facilitate and improve the game character's, and hence the player's, relaxation techniques and planning skills. The game encourages the player to learn and develop new confrontation strategies.

Caglio et al. [28] assessed the modifications occurring in cognitive functions, in particular spatial and verbal memory in a TBI patient after a 3D video game rehabilitation training. The video game was a driving simulator. During the training the participant was requested to explore a complex virtual town from a ground-level perspective.

Cameirão et al. [29] developed the Rehabilitation Gaming System (RGS), a VR based system for the rehabilitation of patients suffering from stroke and TBI. The system uses a camera based motion capture system with gaming technologies to activate intact neuronal systems that provide direct stimulation to motor areas affected by brain lesions. The RGS is designed to engage the patients in task specific intensive training tuned to the patients needs and with continuous monitoring [29, 30, 31].

Ryan et al.[32] described the Balance Rehabilitation Games project which aims to design a game to older adults while incorporating appropriate balance exercises. The game is a maze-solving problem for one or two players. The goal of the game is to navigate the maze and collect all the treasures. The player's score is the final time through the maze. Players move forward by walking in place on a Wii Fit balance board. Longer `strides' produce more rapid progress, to reward better balance rather over rapid stepping. Cooperative and competitive two-player versions of the game are also being prototyped. In the cooperative version, the players work together to collect all the treasures and finish the maze as quickly as possible. In the competitive version the treasures are omitted and it is simply a race to complete the maze as quickly as possible.

Burke et al. [21] developed several games designed for upper limb stroke rehabilitation, which use low-cost webcams as input technology to capture video data of user's movements. The position of the player hands are tracked, so he has to wear a glove or hold a marker which can be an object of a single color, such as a piece of card. The games use user profiling and an option to adaptability.

RehaCom software is used for enabling cognitive rehabilitation [33, 34, 35, 36]. This system is described in more detail in the next section because we find it to be the reference in rehabilitation serious games.

Table 1 displays the classification proposed used the set of criteria defined in Section 4. The "--" means that this feature is not mentioned in the paper reviewed.

**Table 1-** Classification and Comparison of Rehabilitation Serious Games.

| | Betker et al.[24] | Ma et Bech-koum [25] | Conconi et al.[26] | Caglio et al. [27] | Cameirão et al.[28] | Burke et al.[21] | Ryan et al.[31] | System Re-haCom [32] |
|---|---|---|---|---|---|---|---|---|
| **Application Area** | Motor | Motor | Cognitive | Cognitive | Motor and Cognitive | Motor | Motor | Cognitive |
| **Interaction Tech.:** | Body Weight Movement | Motion Tracking + HMD | Speech + Touch+ Motion Tracking + Biosensors | Keyboard | Motion Tracking | Motion Tracking | WiiMote Wii Balance | Special Keyboard + Joystick |
| **Game Interface** | 2D | 3D | 3D | 3D | 3D | 2D | 2D | 2D |
| **No. Players** | Single | Single | Single | Single | Single | Single | Single/ Multi | Single |
| **Competitive / Collaborative** | None | None | None | None | None | None | None | None |
| **Game Genre** | Memory + Simulation | Simulation | Strategy | Simulation | -- | Simulation | Maze | Assorted |
| **Adaptability** | Yes | Yes | Yes | No | Yes | Yes | -- | Yes |
| **Progress Monitoring** | Yes | Yes | Yes | No | Yes | Yes | -- | Yes |
| **Performance Feedback** | Yes | Yes | Yes | -- | Yes | Yes | -- | Yes |
| **Portability** | Home | Clinic | Clinic | Clinic | Clinic/Home | Home | -- | Clinic |

## 6 RehaCom – A Reference System

In this section we describe the RehaCom system in more detail, a system widely used and tested in the area of cognitive rehabilitation. His effectiveness has been demonstrated in a number of studies all very well referenced (with a description of the study conducted) in the RehaCom Catalogue [33]. Since this system is well established in various hospitals and clinics, with a great number of patients, we can more easily have access to it and to the patients being treated by it in rehabilitation programs. Consequently we can more easily study his efficiency. For these reasons, we find it to be a reference system in these area and we choose it as our case study.

RehaCom is a computer-assisted modular system that requires an experienced therapist. The system concept was developed by Hans Regel in 1986 and since then it has been refined over 20 years in clinics, with input from experts in the area. Since 1996 it has been developed by Hasomed (Inc, Ltd). For a few years, it has been market leader in Europe [33] and is currently available in 15 languages.

The system is composed of training procedures for training different skills: attention, memory, executive, field of view and visuomotors. Each training procedure consists of a specific task that the patient must accomplish. As an example of this is

topological memory where the patient has to memorize the position of cards bearing pictures (e.g. of books, a TV, a camera etc.), as in a memory game, or geometric figures. When the pictures are hidden, he needs to know which picture was in which position.

Table 2 displays the training procedures RehaCom offers for each training program or application area.

**Table 2 -** Training Procedures of System RehaCom by Application Area.

| Attention Training | | Memory Training | Executive Functions | Field of View Training | Visuomotor skills |
|---|---|---|---|---|---|
| Alertness | - Acoustic Reactivity<br>- Reaction Behavior | Topological Memory | Shopping | Saccadic Training | Visuomotor Coordination |
| Vigilance | - Vigilance | Physiognomic Memory | Plan a Day | Exploration | |
| Visuo-spatial Attention<br>Selective Attent. | - 2D<br>- 3D<br>- Attention & Concentration | Memory of Words<br>Figural memory | Logical Reasoning | | |
| Divided Attention | - Divided Attention | Verbal Memory | | | |

RehaCom may be classified according with the classification proposed in the previous section, as follows:

- *Application Area*:   The system assists in cognitive rehabilitation. In this field, applications can be various: clinical psychology, geriatrics, developmental psychology, sport psychology, work psychology and driving.
- *Interaction Technology*:  Training with the system can be carried out using a special panel, the computer keyboard, the mouse or a touch screen.
- *Game Interface*: the interface of the games that compose the computer system are all two-dimensional.
- *Number of Players* (Single/Multi-player): RehaCom games are to be played with only one participant at a time.
- *Competitive/Collaborative*: the system only allows one participant at a time, consequently it doesn't permit an interaction competitive or collaborative with other patients.
- *Game Genre*: the system supports games of varying genres. Games classification can be made by category area in which they are included. These categories are the ones showed in Table 2.
- *Adaptability*: The training programs are adaptive which means that task difficulty increases automatically as learning progresses so the patient is not faced with tasks that are too easy or too difficult for him.
- *Progress monitoring*: the therapist can analyze the patient's progress and identify and influence his performance deficits and reserves, even while the patient is training. It enables to monitor progress and adjust training goals as necessary. The results of all training sessions are saved and thus a new session always starts where the last one finished. When the pre-set training time is up, the session ends.
- *Performance feedback*: The system informs the patient of his progress in motivational ways. If an error is made the patient receives specific feedback and the users readily accept the computer's assessment. Training always begins with the instructions. In many programs the instructions are based on the "learning by doing" principle. The first training task then begins. If he needs to, the patient can view the instructions again; he can also take a break whenever he wishes.  As soon as a session ends a performance

chart appears on the screen where the patient sees his progress made from session to session. In addition, it is available also a more detailed description of the results.

- *Portability*: The training programs are designed to enable the patient to train on his own most of the time. It requires the presence of the therapist in order to discuss the patient's training goal and results at the beginning and end of training. Therefore, it reduces considerably the workload of participants in the therapy.

RehaCom system offers a lot of advantages as described above according to the criteria defined, but it has some aspects that could be improved mainly in the lack of the cooperative and competitive dimension and in the game interface which is very simple and is only two-dimensional.

## 7 Conclusions and Future Work

Serious games design is a recent and active research area. Making use of design methodologies, narrative structure, visual arts, interaction techniques and modalities and technologies commonly available in computer games, is proven to be effective in applications without entertainment as the main purpose. In this paper we reviewed the main concepts of serious games and focused on those used with rehabilitation purposes.

We conducted a survey of the most relevant work available in the literature. As the described systems can be very different we proposed a classification in order to assist the comparison and classification of such systems in respect to the more relevant features identified. Based on this classification, existing games could be modified in order to satisfy a large number of the classification criteria and become more functional tools for the rehabilitation therapy. We adopted as a reference system the RehaCom System, a serious game application widely used for cognitive rehabilitation and well grounded in neuro-scientific theory.

As a conclusion we noticed that when compared to popular games that are designed to be engaging and fun, current therapeutic games lack the qualities required to be entertaining and are thus potentially less motivating for patients. Thus, a direction for future work is to identify and measure the impact of the more relevant aspects that can improve the suitability and effectiveness of a game for rehabilitation. We also noticed that almost all reported work is tested with a small number of patients/users.

As another major research opportunity we identify the study of how the effectiveness of computer games for rehabilitation can be increased by incorporation of a social dimension. There is not reported work for systems where collaboration or competitiveness performs a major role on the rehabilitation process. Collaboration and competition add a new dimension that could allow the patients to enjoy the interaction and found the motivation and encouragement from others playing the same game. How attain this collaboration or competition, namely when patients can be at different stages of the rehabilitation process or have different handicaps is an important research problem.

# References

1. Chen, S., Michael, D.: Serious Games: Games that Educate, Train and Inform. Thomson Course Technology PTR, Boston (2006)
2. Serious Games Initiative, http://www.seriousgames.org/. Accessed 19 Dec 2008
3. Numrich, S.K.: Culture, Models, and Games: Incorporating Warfare's Human Dimension. In: IEEE Intelligent Systems, vol. 23, no. 4, pp. 58--61. IEEE Press (2008)
4. Blackman, S.: Serious Games … and Less!. In: ACM SIGGRAPH Computer Graphics , vol. 39 , no. 1, pp. 12--16 (February 2005)
5. Macedonia, M.: Virtual Worlds: A New Reality for Treating Post-Traumatic Stress Disorder. In: IEEE Computer Graphics and Applications, vol. 29, no.1, pp. 86--88. IEEE Press (2009)
6. Sawyer, B.: From Cells to Cell Processors: The Integration of Health and Video Games. In: IEEE Computer Graphics and Applications, vol. 28, no. 6, pp. 83--85. IEEE Press (2008)
7. Von Wangenheim, C.G., Shull, F.: To Game or Not to Game?. In: IEEE Software, vol. 26, no. 2, pp. 92--94. IEEE Press (2009)
8. Mayo, M.: Games for Science and Engineering Education. In: Communications of the ACM, vol. 50, no. 7, pp. 31--35. ACM Press (2007)
9. Kelly, H., Howell, K., Glinert, E., Holding, L., Swain C., Burrowbridge, A., Roper, M.: How to Build Serious Games. In: Communications of the ACM, vol. 50, no.7, pp. 44--49. ACM Press (2007)
10. Zyda, M.: Creating a Science of Games. In: Communications of the ACM, vol. 50, no. 7. ACM Press (2007)
11. Westera, W., Nadolski, R.J., Hummel, H.G.K., Wopereis, I.G.J.H.: Serious Games for higher education: a framework for reducing design complexity. In: Journal of Computer Assisted Learning, vol. 24, pp. 420--432 (2008)
12. Connolly, T.M., Stansfield, Hainey, T.:An Application of games-based learning within software engineering. In: British Journal of Educational Technology, vol. 38, no.3, pp. 416--428 (2007)
13. Zielke, M.A., Evans, M.J., Dufour, F., Christopher, T.V., Donahue, J.K., Johnson, P., Jennings, E.B., Friedman, B.S., Ounekeo, P.L., Flores, R.: Serious Games for Immersive Cultural Training: Creating a Living World. In: IEEE Computer Graphics and Applications, vol. 29, no. 2, pp. 49--60. IEEE Press (2009)
14. Barnes,T., Encarnação, L.M., Shaw, C.D.: Serious Games. In: IEEE Computer Graphics And Applications, vol. 29, no.2, pp. 18--19. IEEE Press (2009)
15. Ardito, C., Buono, P., Costabile, M.F., Lanzilotti, R., Pederson, T., Piccinno, A.: Experiencing the Past through the Senses: An M-Learning Game at Archaelogical Parks. In: IEEE Multimedia, vol. 15, no. 4. IEEE Press (2008)
16. Freitas, S., Jarvis, S.: Serious Games – engaging training solutions: A research and development project for supporting training needs, British Journal of Educational Technology, vol. 38, No.3, pp. 523--525 (2007)
17. Anderson, E.F., McLoughlin, L., Liarokapis, F., Peters, C., Petridis, P., Freitas, S.: Serious Games in Cultural Heritage. In: 10th Int. Symposium on Virtual Reality, Archeology and Cultural Heritage (VAST'09) Short and Project Proceedings, Eurographics, Malta, pp. 22--25 September, pp. 29--48, (2009)
18. Gunasekera, W.S., Bendall, J.: Rehabilitation of Neurologically Injured Patients. In: Gunasekera, W. S., Bendall, J., Neurosurgery, pp. 407-421. Springer London (2005)
19. Zyda, M.: From Visual Simulation to Virtual Reality to Games. In: IEEE Computer, vol. 38, no. 9, pp. 25--32 (2005)
20. Susi, T., Johannesson, M., Backlund, P.: Serious Games – An Overview. Technical report, School of Humanities and Informatics, University of Skovde, Sweden (2005)

21. Burke, J.W., McNeill, M.D.J., Charles, D.K., Morrow, P.J., Crosbie, J.H., McDonough, S.M.: Optimising engagement for stroke rehabilitation using serious games. In: Visual Computer, vol. 25, pp. 1085--1099. Springer Press (2009)

22. Krichevets, A.N., Sirotkina, E.B., Yevsevicheva, E.B., Zeldin, L.M.: Computer games as a means of movement rehabilitation. In: Disability & Rehabilitation, vol. 17, Number 2, pp. 100--105 (1995)

23. Flores, E., Tobon, G., Cavallaro, E., Cavallaro, F.I., Perry, J.C., Keller, T.: Improving Patient Motivation in Game Development for Motor Deficit Rehabilitation. In: Proc. of the 2008 Int. Conf. on Advances in Computer Entertainment Technology, vol. 352, pp. 381--384. ACM Press, Yokohama (2008)

24. Holden, M. K.: Virtual Environments for Motor Rehabilitation: Review. CyberPsychology & Behavior., vol. 8, no.3, pp. 187--211 (2005)

25. Betker, A.L., Desai A., Nett C., Kapadia, N., Szturm, T.: Game-based Exercises for Dynamic Short-Sitting Balance Rehabilitation of People with Chronic Spinal Cord and Traumatic Brain Injuries. In: Physical Therapy, vol. 87, Number 10, pp. 1389--1398 (2007)

26. Ma, Minhua, Bechkoum, K.: Serious Games for Movement Therapy after Stroke. In: IEEE Int. Conf. on Systems, Man and Cybernetics, pp. 1872--1877 (2008)

27. Conconi, A., Ganchev, T, Kocsis, O., Fernández-Aranda, F., Jiménez-Murcia, S.: PlayMancer: A Serious Gaming 3D Environment. In: Int. Conf. on Automated Solutions for Cross Media Content and Multi-channel Distribution - AXMEDIS'08, pp. 111--117. IEEE Press (2008)

28. Caglio, M., Latini-Corazzini, L., D'Agata, F., Cauda, F., Sacco, K., Monteverdi, S., Zettin, M., Duca, S., Geminiani, G.: Video game play changes spatial and verbal memory: rehabilitation of a single case with traumatic brain injury. Journal of Cognitive Processing, vol. 10, pp. S195--S197 (2009)

29. Cameirão, M.S., Badia, S.B., Zimmerli, L., Oller, E.D., Verschure, P.F.M.J.: The Rehabilitation Gaming System: a Review. Studies in Health Technology and Informatics, vol. 145, pp. 65--83 (2009)

30. Cameirão, M.S., Badia, S.B., Zimmerli, L., Oller, E.D., Verschure, P.F.M.J.: A Virtual Reality System for Motor and Cognitive Neurorehabilitation. In: Proc. 9th European Conf. for the Advancement of Assistive Technology in Europe – AAATE 2007, pp. 3--5 October, San Sebastian, Spain (2007)

31. Cameirão, M.S., Badia, S.B., Zimmerli, L., Oller, E.D., Verschure, P.F.M.J.: The Rehabilitation Gaming System: a Virtual Reality Based System for the Evaluation and Rehabilitation of Motor Deficits. In: Proc. of Virtual Rehabilitation 2007, pp. 29--33. IEEE Press (2007)

32. Ryan, M., Smith, S., Chung, B., Cossell, S., Jackman, N., Kong, J., Mak, O., Ong, V.: Rehabilitation Games: Designing computer games for balance rehabilitation in the elderly. Available at http://vivianong.com/images/wiihab/fdg09.pdf. Accessed in 21 Dec 2009

33. RehaCom Catalogue, Schuhfried, available at http://www.schuhfried.at/fileadmin/pdf_eng/catalog_RehaCom_en.pdf. Accessed in 21 Dec 2009

34. RehaCom Basic Manual, Hasomed GmbH, available at http://www.hasomed.de/fileadmin/user_upload/Rehacom/Manuale/ENG/RehaComEN.pdf. Accessed in 21 Dec 2009

35. Mak, M.: Neurorehabilitation as a form of treatment in mental diseases. Terapia, vol.16, no. 1, pp. 47--49 (2008)

36. Maia, L., Gaspar, C., Azevedo, M., Loureiro, M. J., Silva, C. F.: Reabilitação Cognitiva Assistida por Computador: o programa RehaCom e a sua utilização no GEARNeurop. Psiquiatria Clínica, vol. 25, no.2, pp. 83--105 (2004)

# Information and Communication Technologies as a support for learning in Higher Education

Bertil Marques[1,2,3], Carlos Vaz Carvalho[1,2]

[1] Informatics Engineering Department
[2] GILT- Graphics, Interaction and Learning Technologies
[3] ProDEI – FEUP
Oporto, Portugal
[1,2] {bpm,cmc}@isep.ipp.pt; [3] pro09020@fe.up.pt

**Abstract.** In this article it is intended to do an analysis of current developments on the use of Information and Communication Technology (ICT) in the context of Higher Education (HE) in the light of changes in the information society. In particular, focusing on the adoption of online learning platforms and their impact in the teaching-learning model in HE, looking as an example to the integration of MOODLE platform at the Instituto Superior de Engenharia do Porto (ISEP). This analysis seeks to understand whether there was a change of practices by teachers as a consequence of the Bologna model, making teaching more dynamic, less expository, and more ICT supported. However, It was found that there is a low rate of adherence to the use of online teaching platform, even as a repository of information and there is even a much greater reluctance to use the platform for blended learning settings.

**Keywords:** ICT, e-Learning, b-Learning, Higher Education, MOODLE.

## 1    Introduction

Nowadays, and since the 90's, there have been important changes in society's information and economy, now also known as knowledge society and economy. These changes had much to do with technological change, led by the liberalization of telecommunications; it was found that the World Wide Web (Web) grew daily, projecting worldwide that same society, their information and/or it knowledge. So the changes were so striking and disseminated that quickly branched. This growth came to serve all branches of the aforementioned society and economic knowledge.

### 1.1    Technological evolution

The educational systems evolution has emerged as a major challenge for humanity. The evolution of technologies and it impact on the evolution of teaching, particularly in HE, lead us to explore in more detail the evolution on the use of learning support platforms. These developments have affected different types of influences on human

cognitive functions such as memory; imagination, and perception and reasoning, therefore two features cannot be ignored: new forms of access to information and new ways of thinking and developing knowledge, as they help to understand what has changed in the education landscape [1].

Every day, new Web tools focusing on education come up. However, if the use of these technologies is not reflected in the change in procedures, the value of real influence on the results in learning is to be questioned. It is important to note that change from the traditional education to a mode based on Web technologies implies adaptation and a willingness to implement this change in procedures.

In the new generation of Web applications, there are a few ones where resources are shared on the same platform. It is a new level of interaction that facilitates collaboration and information sharing. They facilitate the provision of resources in different formats like text, video and audio, links to sites, notices to students, student-teacher interaction through communication tools, tools to support collaborative learning and registration activities. Those platforms, in general, have been used in blended-Learning (b-Learning) or as a support to face-to-face education in HE [2].

The pressure to use of Web technologies came together with the implementation of the Bologna Declaration[1] - Declaration of Europe [3] that aims to harmonize the structures of Higher Education, with the goal of increasing the competitiveness of European HE system and the promotion of mobility and employability of graduates. Thus, to HE teachers, the demands go beyond the expertise of knowledgeable content. It is intended that teachers monitor the progress of technology and start to contribute in creating new learning environments. The aim is that teachers take students to a "…scientific development and entrepreneurial (...), trained graduates in different fields of knowledge, suitable for insertion in the professional sectors and assist in the continuous training to encourage work search and research, raise the desire for constant improvement by integrating cultural and professional knowledge that are acquired throughout life, (...) to promote critical thinking, freedom of expression and freedom for research…", [4] following the progress of technologies.

## 1.2   Instituto Superior de Engenharia do Porto

Instituto Superior de Engenharia do Porto (ISEP), or Porto Superior Institute of Engineering, had several names and roles since its foundation in 1852 by the minister Fontes Pereira de Melo as Escola Industrial do Porto (Porto School of Industry), to support the industrialization of Portugal, and some time later upgraded to Instituto Industrial do Porto (Porto Institute of Industry).

In 1989, the Instituto Superior de Engenharia do Porto (ISEP) was integrated into the Polytechnic Higher Education subsystem, with two distinct levels: the Bachelor's degree, with duration of three years, and Specialized Studies, lasting two years, which conferred the Graduate degree (Licenciatura). In 1998, under a new reform of the Polytechnic Higher Education system, ISEP started to minister the two-stage degree: Bachelor, first cycle with three-years - followed by a second two-year cycle - to obtain the Graduation. In 2006, due to the adhesion of Portugal to the Bologna

---

[1] Bologna Declaration - Signed in June 1999 by 29 EU countries, published in Portugal, Decreto-Lei nº 42/2005 de 22 de Fevereiro, Diário da República – I Série – A nº 37, pg.1494.

Declaration, ISEP started to offer a new Study Plan, consisting of Graduation and Master in various fields of Engineering [5]. Besides Bachelor (1º Cycle Pre-Bologna), Graduation (1º Cycle Pre-Bologna), Graduation (1º Cycle Bologna), Master's (2º Cycle Bologna), ISEP also offers Post graduation plans in several areas of Engineering, since 2007/2008.

As an institution of HE, ISEP has always followed the development of ICT. It invested on the infrastructures of each department (7 in total) and its services in order to make available the access and use of ICT by teachers, students and services of the Institute, while developing their professional activity.

### 1.3 Subject

In this article the aim of study is the evolution of the usage of the MOODLE learning platform at ISEP, more specifically in the Department of Informatics Engineering (DEI), between 2006/2007 and the 1st semester of 2009/2010 academic year. Data was collected from MOODLE ISEP management with the prior authorization of the direction of ISEP, involving data for each academic year in operation at DEI.

The aim was to check if there have been changes in behavior by teachers to implement the Bologna model in order to make teaching more dynamic, less expository, supported by ICTs, in this particular case with the MOODLE platform.

This study could only be applied to one of ISEP's departments: the choice of the DEI was made because the authors are part of the staff, and have their courses included in the study that will be presented. It is a starting point to make a global study of ISEP.

In chapter one, it is explaine of Web technology that led to the learning support platforms, then a short history of the Instituto Superior de Engenharia do Porto and its link to the online learning platform is given.

In chapter two, the platform chosen for ISEP is explained. The decision to collect information from only one department and the values are presented.

In chapter three, results will be highlighted, analyzed and discussed.

The conclusions based on the analysis in chapter three are done in chapter four.

## 2 The Platform and its application

The MOODLE[2] platform, an open source virtual learning environment created in 2001 by Martin Dougiamas, was chosen as ISEP's learning support platform.

As usual in the open source software, since it was publicly available, MOODLE has been developed collaboratively by a community of professionals from different areas, with new features and functionality continually being added [6].

MOODLE includes a set of features that can be systematized in four basic dimensions: providing content, exercises; existence of tools and media services, synchronous like chats or asynchronous like forums; protected access and

---

[2] MOODLE [9] stands for Modular Object-Oriented Dynamic Learning and simultaneously acronym Martin Object-Oriented Dynamic Learning and Martin is the name of its creator.

management of user profiles; control systems activities. In addition to the functional characteristics it possesses, the MOODLE platform is currently being translated into over 60 languages and is free to use, factors that explain its use all over the world.

## 2.1 MOODLE platform in DEI

The global availability of MOODLE to all ISEP courses was made in 2006/2007 academic year and it is operational till these days. All curricular units are generated in each semester, being allocated a teacher by default as editor in charge. It is not compulsory that teachers effectively use the platform.

The choice of the Department of Informatics Engineering (DEI) to implement this study was due to the fact that the authors are part of the staff, and have their curricular units this study. It is a starting point for a comprehensive study of ISEP, department by department.

The teacher does not define the role of students, who are automatically enrolled in the platform based on information from the system database applications. If students are registered to the curricular unit, they are automatically enrolled as student in the curricular unit of their MOODLE.

Table 1 shows data concerning the use and revitalization of the virtual learning platform MOODLE of DEI-ISEP. Based on these four academic years and the course of Department of Informatics Engineering, the data are for the academic years:

- 2006/2007 – Computer Engineering Degree – 1st Cycle – Bologna
  Computer Engineering Degree – 2nd Cycle – Pre-Bologna
- 2007/2008 – Computer Engineering Degree – 1st Cycle – Bologna
  Master Computer Engineering – 2nd Cycle – Bologna
  Computer Engineering Degree – 2nd Cycle – Pre-Bologna
  Post graduation Computer Science Applied to Health
- 2008/2009 – Computer Engineering Degree – 1st Cycle – Bologna
  Master Computer Engineering – 2nd Cycle – Bologna
- 2009/2010 – Computer Engineering Degree – 1st Cycle – Bologna
  Master Computer Engineering – 2nd Cycle – Bologna

It should be explained that the 1st cycle of the Computer Engineering Degree (Bologna) on DEI-ISEP, is a 12+4 weeks work. The first 12 weeks are theorical, theorical-practice and practice and laboratory classes of regular curricular units. The 4 final weeks in each semester, are practice and laboratory and tutorial classes. In this period there is a special curricular unit named Laboratory Project (LAPR), in which a prototype is developed that involves the contents of the regular curricular units. So the 3 LAPR units for the 1st semester (LAPR1, LAPR3 and LAPR5) and for the 2nd semester (LAPR2 and LAPR4) are created but were not yet fully operational because they started after the collection of data.

**Table 1.** Data regarding use of MOODLE in DEI-ISEP [7].

| Computer Engineering Course | 2006/2007 | 2007/2008 | 2008/2009 | 2009/2010[3] |
|---|---|---|---|---|
| Total Curric. Units. | 29 | 207 | 120 | 111 |
| B-Learning[4] Curric. Units. | 10 | 24 | 29 | 11 |
| Repository[5] Curric. Units. | 10 | 46 | 34 | 19 |
| Total enrolled students | 4129 | 11154 | 13366 | 12966 |
| Total subscribed teachers | 133 | 574 | 628 | 289 |
| Total subscribed teachers-editors[6] | 90 | 320 | 226 | 144 |
| Nº actions VIEW | 3825 | 12830 | 1103444 | 569741 |
| Nº actions ADD | 0 | 22 | 3676 | 2354 |
| Nº actions UPDATE | 2 | 28 | 3931 | 1685 |
| Nº actions DELETE | 0 | 27 | 558 | 427 |
| Nº EVENTS | 62 | 162 | 265 | 91 |
| Nº GROUPS | 29 | 117 | 211 | 7 |
| Nº QUESTIONS | 465 | 464 | 536 | 560 |
| Nº res. DIRECTORY | 78 | 164 | 132 | 83 |
| Nº res. FILE | 762 | 1995 | 1951 | 1048 |
| Nº res. WEB PAGE | 118 | 241 | 95 | 35 |
| Nº res. TEXT PAGE | 8 | 22 | 33 | 7 |
| Nº res. LABEL | 204 | 833 | 601 | 353 |
| Nº mod. ASSIGNMENTS | 13 | 114 | 213 | 64 |
| Nº mod. CHATS | 3 | 1 | 1 | 0 |
| Nº mod. CHOICES | 0 | 0 | 0 | 0 |
| Nº mod. DATABASES | 4 | 2 | 2 | 0 |
| Nº mod. FORUMS | 35 | 227 | 145 | 121 |
| Nº mod. GLOSSARIES | 3 | 1 | 0 | 0 |
| Nº mod. LESSONS | 0 | 0 | 0 | 0 |
| Nº mod. QUIZZES | 14 | 15 | 14 | 14 |
| Nº mod. WIKIS | 0 | 2 | 0 | 0 |
| Nº mod. WORKSHOPS | 0 | 0 | 0 | 0 |
| Nº mod. SURVEYS | 0 | 0 | 9 | 0 |
| Nº Curric. Unit. Files | 3682 | 6342 | 10348 | 6735 |
| Nº Curric. Unit. Files submitted | 842 | 3557 | 6402 | 2041 |

---

3 The data for this academic year was provided on 14 December 2009.

4 We considered all curricular units that have more than 5 events, more than 5 modules or more than 20 files uploaded.

5 Subjects were considered not entering the parameter in previous note in curricular unit area having more than 20 files.

6 Not all teachers of the curricular units are editors; this decision is taken by the responsible of a curricular unit that assigns roles to other teachers.

# 3   Results and discussion

By observing the data in Table 1 an individual analysis between different aspects can be made and see its development.

## 3.1   General Results

In the first year (2006/2007), only 29 curricular units have been created. It was the first year of operation, and although there were 2 plans operating in the DEI, the 1st cycle of Bologna Computer Science Degree and 2nd cycle degree of two phases Pre-Bologna only the requested units were created. 29 created, 69% used the platform (of which 50% in b-learning and the other 50% in the repository), so 31% of teachers-editors responsible for these curricular units did even not log in, as can be seen on Figure 1.

**2006/2007**

■ % C.U. working     ■ % C.U. not working



**Fig. 1.** Curricular Units working vs. not working in the academic year 2006/2007, the first year of entry into operation of the platform MOODLE in DEI-ISEP.

**Registration Data.** To help us understand the values presented, in Table 2 data by academic year with the number of curricular units vs. total curricular units running on the platform (b-Learning + repository) is show, as well as the number of students' enrollments. In Figure 2 it is shown a visual analysis of the Curricular Units on MOODLE platform total vs. working units.

**Table 2.** Data regarding the use of MOODLE in DEI-ISEP in the academic years indicated, only the totals of curricular units, students enrolled and registered teachers [7].

| Computer Engineering Course | 2006/2007 | 2007/2008 | 2008/2009 | 2009/2010 |
|---|---|---|---|---|
| Total Curric. Units. | 29 | 207 | 120 | 111 |
| Total enrolled students | 4129 | 11154 | 13366 | 12966 |
| Total Curric. Units. Working | 20 | 70 | 63 | 30 |

**Fig. 2.** Curricular Units on MOODLE platform on the academic years 2006/2007, 2007/2008, 2008/2009 and 2009/2010.

**Actions Results.** MOODLE ISEP administration office, informed that data on the numbers of actions VIEW, ADD, UPDATE and DELETE is valid only for the academic years of 2008/2009 and 2009/2010, due to a problem with the platform server. The previous data were erased and turned to zero at the beginning of 2008/2009; therefore in table 3 the data since that reset time will be presented.

Remains to explain the meaning of each action:

VIEW – Any click that runs on any resource available on the curricular unit, do by students or teachers;

ADD – Addition of a resource in a curricular unit, do by teacher-editor;

UPDATE – Change on a resource available on a curricular unit, do by teacher-editor;

DELETE – elimination of a resource on a curricular unit, do by teacher-editor.

**Table 3.** Data regarding the use of MOODLE in DEI-ISEP in the academic years indicated, only the values of the number of actions VIEW, ADD, UPDATE and DELETE [7].

| Computer Engineering Course | 2006/2007 | 2007/2008 | 2008/2009 | 2009/2010 |
|---|---|---|---|---|
| Nº actions VIEW | a) | a) | 1103444 | 569741 |
| Nº actions ADD | a) | a) | 3676 | 2354 |
| Nº actions UPDATE | a) | a) | 3931 | 1685 |
| Nº actions DELETE | a) | a) | 558 | 427 |

a) Results not considered because they are incomplete.

**Fig. 3.** Variation in the number of actions VIEW performed in platform MOODLE in DEI-ISEP.



**Fig. 4.** Variation in the number of actions ADD, UPDATE and DELETE performed in the platform MOODLE in DEI-ISEP.

**Resources and Modules Results.** Here are included all the features to support learning as the activities module ASSIGNMENTS, CHATS, CHOICES, DATABASES, FORUMS, GLOSSARIES, LESSONS, SURVEYS, QUIZZES, WIKIS, and WORKSHOPS. These modules are those that allow the various types of interaction between students/teachers, between students/students and between teachers/teachers. By the numbers in Table 1 only three of these activities are preferred, ASSIGNMENTS, FORUMS and QUIZZES, as can be seen in Table 4 and in Figure 5.

**Table 4.** Data regarding use of activities ASSIGNMENTS, FORUMS and QUIZZES on MOODLE in DEI- ISEP, on the academic years indicated [7].

| Computer Engineering Course | 2006/2007 | 2007/2008 | 2008/2009 | 2009/2010 |
|---|---|---|---|---|
| Nº mod. ASSIGNMENTS | 13 | 114 | 213 | 64 |
| Nº mod. FORUMS | 6 | 20 | 25 | 10 |
| Nº mod. QUIZZES | 14 | 15 | 14 | 14 |



**Fig. 5.** Variation in number of activities ASSIGNMENTS, FORUMS and QUIZZES on MOODLE in DEI-ISEP, in the 4 academic years marked.

The module ASSIGNMENTS is more often used because it allows dated deliveries defined by the teachers and it is reliable. The module FORUMS is being now also more used by the students.

**Files Results.** Figure 6 presents the number of files added in curricular units as well as the number of files submitted through the platform also in curricular units.

**Fig. 6.** Variation of number of files vs. files submitted curricular units on the platform in DEI-ISEP, in the 4 academic years.

One can see that values have increased every academic year that passes.

## 3.2 Discussion of Results

Whereas valid then and now with the curricular units enrolled automatically in the following academic years, notice that the number of curricular units between the academic years 2007/2008 and 2008/2009 is drastically low; this was the conclusion of the transitional period between the bi-step system and Bologna. In the academic year 2008/2009 there was no edition of the Post graduation Computer Science Applied to Health reducing then four curricula for only two.

**Registration Results Discussion.** Looking at the full data of the four academic years, the oscillations mentioned in Table 2 can be seen, as far as he total number of curricular units of plans and the number of curricular units operating on the platform (b-learning + repository). At the time of data received in this academic year 2009/2010, the 111 courses planned and created for the 1st and 2nd semester, yet only 30 were operational on the platform, 11 in b-Learning and 19 in repository, as shown in Figure 2. All previous years have proved low compliance in b-Learning or as a repository.

**Actions Results Discussion.** For the data of 2008/2009, a point of comparison cannot be made as that is only one academic year with this information. However, if the values of 2009/2010 are looked at, while representing only 27% of the disciplinary units of academic year, already has values close to those presented in 2008/2009, as seen in Figures 3 and Figure 4.

**Resources and Modules Results Discussion.** Resources DIRECTORY, FILE, WEB PAGE, TEXT PAGE and LABEL, do not show relevant differences in presented values, only the resources WEB PAGE and LABEL where values decrease slightly.

The number of elements in ASSIGNMENTS represents the number of materials (Word documents, PDF, PowerPoint, or others) that teachers allow students to submit in the platform. In FORUMS shows that although the values seem reasonable, the

reality is that MOODLE automatically creates a default forum for curricular unit set up, so the actual values are for 2006/2007 of 6 forums, for 2007/2008 of 20 forums, 2008/2009 were 25 forums and with the data that is currently available for 2009/2010 of 20 forums. The data have been gathered in Table 4 and Figure 5 shows the distribution. The QUIZZES module is used in all academic grades, had values of percentage use of 7%, 21.4%, 22% and 46% (which cannot be taken into account as mentioned above), and was a far low. The other modules have null or meaningless values.

**Files Results Discussion.** Two parameters that have values always increasing are the numbers of files in the curricular units and the numbers of files sent on curricular units. After analyzing the figures for 2009/2010, as there are values of a magnitude significant considering what has been said about the specificity of the date were acquired from the data being analyzed. By using the academic year 2008/2009 and 2009/2010 platform just to plan courses of Bologna (1st and 2nd cycle) there has been a rise in the number of subjects using the MOODLE, however, with these numbers one could see and anticipate that have been mainly as a repository of information.

# 4    Conclusions

In order to make a good study, more information will be needed, at least 2 more academic years. It will also have to be more subdivided: instead of having information on the academic year, to have the semester/ academic year and eventually be able to provide data by discipline/semester/academic year/curricular unit.

Though from the various tables and figures presented, it is found that there is still have a low rate of adherence to the use of MOODLE as a mere repository of information, but there is a much greater reluctance when the numbers for a real use of the potential that the platform allows are analyzed.

Also the perception resulting of academic activities and close contact with reality with many professional colleagues/teachers, is that the introduction of the MOODLE platform in school functioned as a motivating factor for some who saw in this tool an opportunity to innovate some of their teaching practices. For others, it was just the curiosity to try a new technology, initially with the potential at the virtual online learning or an extra "difficulty" in their learning in ICT. On a second phase, some went ahead motivated, others left for various reasons. Others, still, may have never engaged to the point of finding that extra difficulty that made some arrive to the second stage.

While talking about a department technologically advanced, with teachers who should be receptive to these new technologies and new teaching methodologies, there hasn't been a real support and commitment in using it more effectively.

It appears that, due to lack of knowledge, inertia or simply for convenience, there is some resistance from some teachers who continue to use in their classes the exhibition model, making the students to take a passive attitude of mere recipients of knowledge. It is also found in these cases a lack of vision, when analyzing the records of units in the field for the planning, component of distance learning or more specifically the component of e-Learning [8].

In HE institutions, including ISEP, the term "blended learning" (b-Learning) is often used incorrectly. It hides the fact that the institutions use a few ICT resources, and continue to use mostly the same kind of teaching and learning that used previously.

Baring in mind the analysis of the use of the platform, it seems that a mere introduction of an ICT platform is not enough to change the teaching/learning paradigm: it would also be necessary a change of teaching methodology used by teachers on their classes.

Therefore, with the information were given, we conclude that there were some changes in behavior, some of which have been seen by direct observation. But as for the hold up of ICT, to support the online learning platform available to MOODLE and how this verification is helped by the figures above, the hypothesis placed initially could not be proven.

# References

1. Baptista, R., Carvalho, C. V.: O Ensino através de ambientes de jogo RPG: Uma experiência num contexto específico. In: V CITICE. Braga (2007)
2. Carvalho, A. A. A.: Rentabilizar a Internet no Ensino Básico e Secundário: dos Recursos e Ferramentas Online aos LMS. In: Revista de ciências da educação. May/August, Vol.3, sísifo (2007)
3. Website of Direcção Geral do Ensino Superior - DGES - MTCES, http://www.dges.mtces.pt/DGES/pt/ [Cited: December 8, 2009]
4. Freitas, J. A.G.: Bolonha e a formação universitária e profissional em Ciência da Informação. In: Cadernos de Biblioteconomia, Arquivística e Documentação, Issue 1, pp. 10--15 (2006)
5. Website of Instituto Superior de Engenharia do Porto, http://www.isep.ipp.pt [Cited: December 8, 2009]
6. Lopes, A. M., Gomes, M. J. B.: Ambientes virtuais de aprendizagem no contexto do ensino presencial: Uma abordagem reflexiva. In: V CITICE. Braga (2007)
7. Website access to MOODLE Instituto Superior de Engenharia do Porto, http://moodle.isep.ipp.pt/admin [Cited: December 10, 2009]
8. Almeida, P. F.P.T.: Uma Metodologia para a Integração das tecnologias WEB nas Unidades Curriculares de Sistemas e Sistemas de Tecnologias de Informação no Ensino Superior. Minho University - Engineering School. Doctorate Thesis (2008)
9. Website of MOODLE platform, http://moodle.org [Cited: December 20, 2007]

# PAPERS IN ALPHABETICAL ORDER

# AUTHORS IN ALPHABETICAL ORDER