

**A WRAPPER SELF-MANAGEMENT VIABLE SYSTEM
FRAMEWORK FOR INFORMATION EXTRACTION
IN HEALTH CARE**

SAMEER QASIM AMEEN

UNIVERSITI KEBANGSAAN MALAYSIA

A WRAPPER SELF-MANAGEMENT VIABLE SYSTEM FRAMEWORK FOR
INFORMATION EXTRACTION IN HEALTH CARE

SAMEER QASIM AMEEN

THESIS SUBMITTED IN FULFILMENT FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY

FACULTY OF INFORMATION TECHNOLOGY AND COMPUTER SCIENCES
(FTSM)
UNIVERSITI KEBANGSAAN MALAYSIA
BANGI

2013

KERANGKA KERJA PENYALUT SISTEM KEBOLEHTAHANAN
PENGURUSAN-KEKENDIRIAN DALAM PENGEKSTRAKAN
MAKLUMAT KESIHATAN

SAMEER QASIM AMEEN

TESIS YANG DIKEMUKAKAN UNTUK MEMPEROLEH IJAZAH
DOKTOR FALSAFAH

FAKULTI TEKNOLOGI DAN SAINS MAKLUMAT
UNIVERSITI KEBANGSAAN MALAYSIA
BANGI

2013

DECLARATION

I hereby declare that the work in this thesis is my own except for quotations and summaries which have been duly acknowledged.

03 October 2013

SAMEER QASIM AMIN

P34977

ACKNOWLEDGMENT

First and foremost, I'd like to thank my supervisor and mentor, Prof. Dr. Khairuddin Bin Omar for his technical advice, moral support and the friendship that I have developed over five years in UKM. He has left a lasting impression on me as a fine advisor who is full of patience and tolerance yet having his high expectations met.

I would also like to thank Prof. Dr. Shahrul Azman Mohd. Noah , Prof. Dr. Azurallza Abu Bakar and Dr. Mohd Zakree bin Ahmad Nazri for their support and valuable advice throughout the course of my study.

Last but not least, I would like to dedicate this thesis to my family, my wife Maram, my daughter Dalia and my sons Yazan and Zaid for their love, patience, and understanding that allowed me to spend most of my time on this thesis.

ABSTRACT

This study presents a novel infrastructure of an information extraction system underpinned by two existing models, namely; Distributed Object Oriented Programming and Viable System Model (VSM). Autonomic systems infrastructure has been used successfully in many global computing systems like grid and cloud. Such models for computing environments are required to exhibit self-managing capabilities to reduce cost and complexities and to improve dependability. This research aims to investigate interaction scenarios which enable humans and machines to collaboratively solve problems, applying the interaction results into an autonomic system infrastructure in medical domains by using wrappers and Meta knowledge to control the process of information extraction from patients' records. The main approach in this research for solving the problem of information extraction is, using wrapper approach which relies on semi-structured documents. This research also presents the generic requirements for tools, services and frameworks to facilitate the design and development of self-managed systems dealing with the problem of information extraction with the help of wrappers. This study used consistent query as an approach for interaction between the users to provide a collaborative architecture to extract medical care information data for users to share their experience and knowledge. The methodology of this research consists of two main parts, namely; Wrapper Retrieval and Wrapper Generator. Wrapper Retrieval consists of a user interface and a Wrapper Key Retrieval to recognize the users' needs, search and present a suitable wrapper to the user from the pool of wrappers. In case that none of the wrappers in the pool meets the user's needs, the Wrapper Generator starts functioning by creating a new wrapper which suites the user's inquiries. An autonomic infrastructure system is used to build an information extraction system with the help of wrapper automation for health information where the output from the system is a text document (report) extracted from patients' medical records using automatic extracting wrappers. To implements this methodology, a working system was build, tested and installed in two hospitals for test running. A monitoring model has been designed in this research to provide the situated autonomic computing with feedback and context information from their environment. The implication of this approach represent a significant step forward and may serve as a road sign indicating that this path is an efficient yet easy to implement and control. In support of the monitoring and awareness services, the researcher also conducted a survey questionnaire to collect responses from the users about the performance of the system. The findings which were evaluated using quantitative and qualitative method indicate maximum users' satisfaction.

ABSTRAK

Kajian ini mencadangkan infrastruktur baru untuk sistem pengekstrak maklumat yang berasaskan dua model sedia ada iaitu; Pengaturcaraan Berorientasikan Objek Teragih dan Model Sistem Berdaya Maju (VSM). Infrastruktur sistem autonomi telah berjaya digunakan dalam banyak model komputeran global seperti grid dan awan. Penggunaan model seumpamanya sebagai persekitaran komputeran adalah diperlukan untuk mempamerkan keupayaan pengurusan sendiri dalam menurunkan kos dan kompleksiti serta meningkatkan kebolehharian. Penyelidikan ini bertujuan untuk menyiasat senario interaksi yang membolehkan manusia dan mesin menyelesaikan masalah secara kolaboratif, mengaplikasi keputusan interaksi ke dalam suatu infrastruktur sistem autonomi untuk domain perubatan menggunakan pembalut dan pengetahuan meta bagi mengawal proses pengekstrakan maklumat rekod pesakit. Kaedah kajian terdiri daripada dua komponen utama, iaitu: bahagian capaian pembalut dan bahagian penjana pembalut. Kajian ini turut menunjukkan keperluan generik alatan, khidmat dan kerangka kerja untuk memudahkan reka bentuk dan pembangunan sistem pengurusan sendiri yang berurusan dengan masalah pengekstrak maklumat menggunakan pembalut. Kajian ini juga menggunakan pertanyaan konsisten sebagai pendekatan interaksi antara pengguna sebagai suatu seni bina kolaboratif bagi pengekstrak maklumat penjagaan kesihatan dengan pengguna dapat berkongsi pengetahuan dan pengalaman. Pendekatan dalam penyelidikan untuk menyelesaikan masalah pengekstrak maklumat adalah menggunakan pendekatan pembalut yang berkait separa-terstruktur dokumen. Kaedah kajian ini terdiri daripada dua bahagian iaitu capaian pembalut dan penjana pembalut. Capaian pembalut terdiri daripada antara muka pengguna dan capaian kunci pembalut untuk mengecam keperluan pengguna, mencari dan membentangkan pembalut yang sesuai bagi pengguna daripada lubang pembalut. Sekiranya tiada pembalut dalam lubang yang sesuai dengan jangkaan pengguna, bahagian penjana pembalut bertindak menghasilkan pembalut baru yang sesuai dengan keperluan pengguna. Infrastruktur sistem autonomi telah digunakan untuk membina sistem pengekstrak maklumat menggunakan automasi pembalut bagi maklumat kesihatan. Output bagi sistem adalah dokumen teks (laporan) yang diekstrak daripada rekod perubatan pesakit menggunakan pembalut pengekstrak binaan automatik. Bagi melaksanakan kaedah ini, sistem sebenar telah dibangunkan, diuji dan dipasang pada dua hospital untuk ujian larian. Model pemantauan telah direka bentuk dalam penyelidikan ini bagi menyediakan komputeran autonomi terkondisi dengan maklum balas dan konteks maklumat daripada persekitaran. Implikasi pendekatan ini menggambarkan satu langkah ke hadapan yang ketara dan dapat melayan sebagai tanda jalan yang menunjukkan bahawa langkah ini adalah efisien, mudah diimplementasi dan dikawal. Untuk menyokong perkhidmatan pemantauan dan kesedaran, penyelidik turut menjalankan kajian soal selidik untuk mengumpul respon dari pengguna mengenai prestasi sistem. Penemuan yang dilakukan secara kuantitatif dan kualitatif menunjukkan kepuasan pengguna yang tinggi.

TABLE OF CONTENTS

	Page
DECLARATION	iii
ACKNOWLEDGEMENTS	iv
ABSTRACT	v
ABSTRAK	vi
TABLE OF CONTENTS	vii
LIST OF FIGURES	xi
LIST OF TABLES	xiii
LIST OF ABBREVIATIONS	xiv
 CHAPTER I INTRODUCTION	
1.1 Overview	1
1.2 Research Background	1
1.3 Problem Statement	4
1.4 Research Objectives	5
1.5 Research Questions	6
1.6 Scope of Research	6
1.7 Research Methodology	7
1.7.1 Theoretical Study	8
1.7.2 Designing	9
1.7.3 Implementation & Observation	10
1.7.4 Performance Evaluation	10
1.8 Significance and Importance of Research	10
1.9 Thesis Outline	11
 CHAPTER II LITERATURE REVIEW	
2.1 Introduction	12
2.2 Autonomic computing	13
2.3 Definitions	13

2.3.1	Characteristics of Autonomic Computing	14
2.3.2	Autonomic Computing Capabilities	15
2.3.3	Autonomic Computing Standards	16
2.3.4	Autonomic Computing Interoperability Standards	18
2.4	Monitoring and Self-Awareness	20
2.4.1	Software Instrumentations	21
2.4.2	Intelligent Monitoring System	22
2.5	Self-Management System Requirements	23
2.5.1	Description languages	23
2.5.2	Frameworks	23
2.5.3	Services and Utilities	25
2.6	Self-Management Middleware for Planetary Scale System	28
2.6.1	Self-Organizing	28
2.6.2	Self-Configuration	29
2.6.3	Self-Optimizing	30
2.6.4	Self-Protective	30
2.6.5	Self-Healing	31
2.7	Design and Implementation Requirements	31
2.8	Information Extraction (IE)	35
2.9	IE Methodologies	37
2.9.1	Rule Learning Based Extraction Methods	37
2.9.2	Classification Based Extraction Methods	47
2.9.3	Sequential Labelling Based Extraction methods	53
2.10	VSM Model: a Brief Overview	55
2.11	Information Extraction (IE) in Health Care	57
2.12	Summary	59
CHAPTER III PROPOSED FRAMEWORK		
3.1	Introduction	63
3.2	Building Wrappers for IE Cooperatively	64
3.3	Possible Scenarios and Approaches	67

3.3.1	Approach 1	67
3.3.2	Approach 2	67
3.3.3	Approach 3	68
3.3.4	Approach 4	68
3.3.5	Approach 1'	69
3.3.6	Approach 2'	69
3.4	The IIEHC Approach: Consistency Queries	70
3.4.1	IIEHC Interaction Scenario	71
3.4.2	System Framework (Wrapper Generator Part)	74
3.4.3	System Framework (Wrapper Retrieval Part)	89
3.4.4	Measurements for extracting patterns	97
3.5	Summary	98
CHAPTER IV IMPLEMENTATION		
4.1	Introduction	100
4.2	Wrapper Self-Management System (WRSMS)	100
4.2.1	SM-VSM Pattern For Self-Management System	100
4.2.2	GoF and SM-VSM	103
4.3	Wrapper Self-Management Viable System Scenario	105
4.4	Illustrative Examples	111
4.4.1	Abstract Factory	111
4.4.2	Builder	114
4.4.3	Factory Method	115
4.4.4	Prototype	116
4.5	Summary	116
CHAPTER V QUESTIONNAIRE DESIGN		
5.1	Introduction	118
5.2	Information Collection Tools	119
5.3	Quantitative Descriptive Research Methods	121
5.3.1	Advantages of Survey Research	121
5.3.2	Disadvantages of Survey Research	122
5.4	Classification of Survey Methods	122

5.5	Internet Surveys Methods	124
5.6	Evaluation Methodology	128
5.7	Questionnaire Design	130
5.8	Quantitative and Qualitative Evaluation	137
5.9	Data Collection	138
5.10	Summary	140
 CHAPTER VI WRAPPER SELF-MANAGEMENT SYSTEM (WRSMS) EVALUATION		
6.1	Introduction	141
6.2	Analyzing the Data	141
6.3	Interpretation and Reflection	149
6.4	Recommendations	150
6.5	Summary	150
 CHAPTER VII CONCLUSION AND FUTURE WORK		
7.1	Motivations and Approach	151
7.2	Achievements and Contributions	154
7.3	Future Work	155
REFERENCES		157
 APPENDICES		
A	Test data	165
B	Representation, and Output Text Files	168
C	Programs Sources	170
D	List of Publications	229
E	Gang of Four-Structural Patterns	230
F	Medical Report – Sample	236

LIST OF FIGURES

Figure No.		Page
Figure 1.1	The basic approach behind IIEHC	3
Figure 1.2	Research Stages and Steps	9
Figure 2.1	Functional Details of the Autonomic Manager	18
Figure 2.2	Autonomic computing interoperability standards	18
Figure 2.3	Services and Tools of the Self-Management General System	27
Figure 2.4	Self-configuration scenario	30
Figure 2.5	Automatic information extraction methods	36
Figure 2.6	Example of autoslog concept node	38
Figure 2.7	Example of wrapper induction	44
Figure 2.8	The learn LR Algorithm	45
Figure 2.9	The BWI algorithm	47
Figure 2.10	General approach for building a classification model.	48
Figure 2.11	Example of information extraction as classification	50
Figure 2.12	Extracting Processing Flow In The Two-Level Classification Approach	53
Figure 3.1	The proposed system intelligent information extraction for health care(iiehc) overview	64
Figure 3.2	The proposed IIEHC Interaction Scenario	72
Figure 3.3	Overview of IIEHC Proposed System	73
Figure 3.4	Flowchart shows the learning process	74
Figure 3.5	Tree representation for medical record information	78
Figure 3.6	XML representation for the information inherent in medical records.	79
Figure 3.7	A Multiple Alignment for Three Token Strings	86
Figure 3.8	Key to String Encoder Output	87
Figure 3.9	Patterns Table Construction Process	88
Figure 3.10	Sample patterns table	90
Figure 3.11	Wrapper Code Generation	91
Figure 3.12	User interface (pattern viewer)	91

Figure 3.13	Numbering the medical abbreviation using the Tree	92
Figure 3.14	Meta Data base table	93
Figure 3.15	Altova process for producing a meta-Data report.	94
Figure 3.16	The meta-data Output Report	94
Figure 3.17	A reporting wrapper	95
Figure 3.18	Sample output from wrapper of Figure 3.17 above	95
Figure 3.19	Mapping the input fields to the output report	96
Figure 3.20	First page of Medical Evaluation – Physician’s Report	97
Figure 4.1	Self-management viable system model	101
Figure 4.2	Wrapper Self-Management system UML Use Case Diagram	106
Figure 4.3	Wrapper self-management system - sequence diagram	107
Figure 4.4	Wrapper self-management system -activity diagram	108
Figure 4.5	Wrapper self-management system -UML class diagram	110
Figure 4.6	Gof Abstract Factory Pattern	112
Figure 4.7	C# Code for Generating an Abstract Factory for Self-Tuning Capability	112
Figure 4.8	Gof Builder Pattern	114
Figure 4.9	Gof Factory Method Pattern	115
Figure 4.10	Gof Prototype Pattern	116
Figure 5.1	Types of quantitative descriptive research methods	121
Figure 5.2	Classifications of Survey Methods	122
Figure 5.3	Improving response rates methods	127
Figure 5.4	Activity diagram of the questionnaire process	129
Figure 5.5	Opening screen	131
Figure 5.6	Medical specialty	131
Figure 5.7	Years the participant been practicing medicine and where.	133
Figure 5.8	Opinion about the old system	135
Figure 5.9	Level of satisfaction	137
Figure 5.10	List of all specialties	139
Figure 5.11	Application of Crosstab for the first 5 specialties	140

LIST OF TABLES

Table No.		Page
Table 2.1	Requirements for the Survival of the Self-Management Service for the Planetary-Scale System	34
Table 2.2	Example of initial tagging rule	41
Table 2.3	Example of generalized tagging rule	42
Table 2.4	Example of Correction Rule	43
Table 2.5	Major Systems of Viable Systems Model	56
Table 3.1	Subset of named entities for medical records' domain.	66
Table 3.2	List of meaning of some entities from medical field.	77
Table 3.3	The encoding hierarchy for text document.	80
Table 5.1	Information Collection Tools	119
Table 5.2	Major Criteria for Selecting a Survey Method	126
Table 5.3	Number of participants after applying Crosstab for the first 5 specialties	139
Table 6.1	Years of practicing medicine for each medical specialty	142
Table 6.2	Practicing medicine in this community	142
Table 6.3	The type of practice for participants	143
Table 6.4	Participants Opinion regarding the old system	144
Table 6.5	Most important factors in participant's decision	146
Table 6.6	Participants' satisfaction with IIEHC	147
Table 6.7	Comparing IIEHC with similar product.	148
Table 6.8	The extent to which the participants are likely to continue using IIEHC	149

LIST OF ABBREVIATIONS

Abbreviation	Description
ASIDL	Assembly Services and Infrastructures Description Language
ASIF	Assembly Services and Infrastructures Framework
BMU	Best Matching Unit
CIM	Common Information Model
DISCO	Discovery
EHMS	E-Health Monitoring System
GoF	Gang of Four
JSU	Job Schedule Unit
KNN	K-Nearest Neighbourhood
MOWS	WSDM Management of Web Services
MSDL	Monitor Session Description Language
MUWS	WSDM Management Using Web Services
ODS	On-Demand Service
OGSA	Open Grid service Architecture
QoS	Quality of Service
SADL	Sensor and Actuator Description Language
SAF	Sensors and Actuators Framework
SLA	Service Level of Agreement
SM-VSM	Self-Management Viable System Model
SOAP	Simple Access Description Language
SOM	Self Organising Map
SRU	Service Reservation Unit
SSM	Soft Systems Methodology
SVM	Support Vector Machine
UDDI	Universal Discovery, Description and Integration
UML	Unified Modelling Language
VServer	Virtual Servers
VSM	Viable System Model
WS	Web Service
WSDL	Web Service Description Language
WSDM	Web Services Distributed Management
XML	Extensible Markup Language

CHAPTER I

INTRODUCTION

1.1 OVERVIEW

This chapter contains an explanation of the various components of this research. This introductory chapter seeks to set the context by which the study is framed and encompasses an explanation of the background to the problem, the problem statement, research objectives, research scope, research methodology, research significance and the contribution of the research to the existing body of knowledge and to existing environments and its importance thereof. In general, this chapter not only gives a preliminary depiction of the research, but also provides an executive summary of the entire work.

1.2 RESEARCH BACKGROUND

Information Extraction (IE) represents a rich research area within the domain of Natural Language Processing (NLP). NLP possesses the potential for addressing many real-world applications such as database access, question answering, and information extraction. However, in-depth natural language processing is an expensive endeavour that can strain computational resources. As an alternative to full-blown natural language processing, some researchers in the NLP community have turned their attention to adopt (IE) techniques.

The goal of an IE system is to extract specific types of information from text. For example, an IE system might extract the names of all persons, their age distributions, physical addresses etc. This simplifies the job of the NLP system

considerably. IE is less computationally expensive than full-blown natural language processing because many phrases and even entire sentences can be ignored if they are not relevant to the domain. Since the system is only concerned with the domain-specific portions of the text, some of the most difficult problems in NLP are simplified (e.g., part-of-speech tagging, ambiguity resolution, etc.). IE is a more practical and robust technology than in-depth NLP and has met with notable success in recent years.

IE systems are capable (at least partially) of understanding the content in the input. This capability is realised by cutting out the uninteresting portions of the text and extracting information from the remainder. This information is then structured and provided as output (Vincent Breton 2005). Different techniques are proposed for information extraction methodologies. These techniques are summarised by the following approaches (Hercules 2007) : Named Entity Recognition, Co-reference Resolution, Template Element Construction, Template Relation Construction and Scenario Template Production. These approaches are chosen according to the domain characteristics. Due to the variety of information to be extracted, one single method may not be suitable for all cases.

Integrated Information Extraction for Health Care (IIEHC) as depicted in Figure 1.1 uses NLP techniques to combine information from multiple sources represented as text data and automatically extracts patient-information from them (databases, text), then the system allow the user to use these information to produce different types of outputs. This is done by using programs called wrappers; these are programs that extract contents of a particular information source and put it into structured data form.

There are two main approaches to wrapper generation: wrapper induction and automated data extraction. Wrapper induction uses supervised learning to learn data extraction rules from manually labelled training examples. Automated wrapper generation using unsupervised pattern mining. Automated extraction is possible if the system follow fixed templates. Discovering such templates or patterns enables the system to perform extraction automatically (Liu 2011).

Wrapper generation is an important problem with a wide range of applications. Extraction of such data enables one to integrate data/information from multiple sources to provide value-added services.

For any health care system, the most important factor is the efficiency of extracting accurate patient information this information are unique to each patient and constitutes as the backbone of the health care system upon which all subsequent medical actions depend.

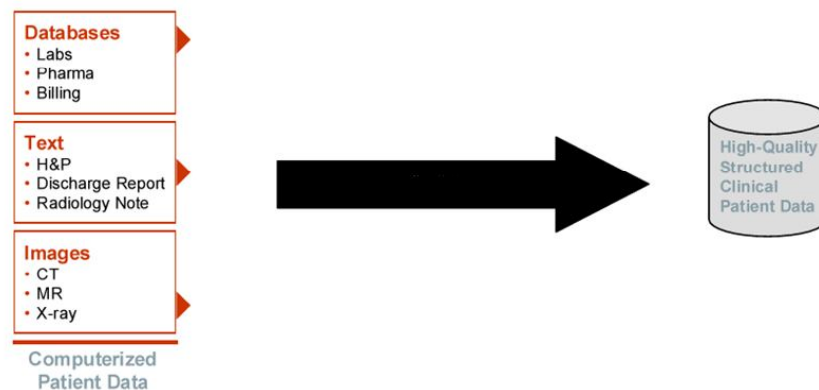


Figure 1.1 The basic approach behind IIEHC

Complex tasks in health care systems are often solved in teams as no one individual has the collective expertise, information, or resources required for effective performance. Among the more notable benefits of working in teams is the combination of a collective genius and its application in problem solving, which by nature involves a multitude of activities such as gathering, interpreting and exchanging information, creating and identifying alternative courses of action, by arriving at a common resolution amongst the differing and often conflicting preferences of team members, and ultimately implementing a choice, determining how incremental progress will be measured, and monitoring its evolution. Such efforts and initiatives rest on the solid assumption that the collective genius is preferable to that of any given individual.