



UNIVERSITI PUTRA MALAYSIA

**TRANSFORMATION OF EXTRACTED KNOWLEDGE IN MALAY
UNSTRUCTURED DOCUMENTS INTO AN INTERROGATIVE
STRUCTURED FORM**

FATIMAH SIDI.

FSKTM 2007 10



**TRANSFORMATION OF EXTRACTED KNOWLEDGE IN MALAY
UNSTRUCTURED DOCUMENTS INTO AN INTERROGATIVE
STRUCTURED FORM**

By

FATIMAH BINTI SIDI

**Thesis Submitted to the School of Graduate Studies, Universiti Putra
Malaysia, in Fulfilment of the Requirements for the Degree of Doctor of
Philosophy**

SEPTEMBER 2007



** I admonish you, lest you be one of the ignorant * (Qur'an 11:46)*

Knowledge is a light that leads to wisdom.

It is life for one's soul and fuel for one's character.

**And say: "My Lord! Increase me in knowledge * (Qur'an 20:114)*

If you desire happiness,

then seek out knowledge and enlightenment,

and you will find that anxiety, depression, and grief will leave you.

Abstract of thesis presented to the Senate of Universiti Putra Malaysia in fulfillment of the requirement for the degree of Doctor of Philosophy

**TRANSFORMATION OF EXTRACTED KNOWLEDGE IN MALAY
UNSTRUCTURED DOCUMENTS INTO AN INTERROGATIVE
STRUCTURED FORM**

By

FATIMAH BINTI SIDI

September 2007

Chairman: Associate Professor Hj. Mohd. Hasan Selamat

Faculty: Computer Science and Information Technology

The availability of knowledge discovery operation helps to extract valuable information and knowledge in large volumes of data in structured databases. However, a large portion of the available information is not in structured form but rather collections of text documents in unstructured format, which also implies to Malay unstructured documents. Therefore, structuring characteristics must be imposed to unstructured documents in order to transform information available in unstructured documents into knowledge. A new approach has been established to transform extracted knowledge in Malay unstructured document by identifying, organizing, and structuring them into interrogative structured form. Its architecture is developed based on the implementation of (i) interrogative knowledge identification; (ii) interrogative contextual information; and (iii) interrogative knowledge organization and

structuring with Malay knowledge representation by concepts. It utilizes the Malay language corpus; interrogative theory; as well as object-oriented, ontology, and database model. The research involves system development based on architecture of the Malay/K-Ontology, which is being measured by quantitative retrieval performance using the recall and precision metrics. The development of the Retrieval Interrogative Ontology Analysis Application is used to verify fitness of task for the functionalities and usefulness on the utilization of interrogative contextual information with color coding supplement, additional information annotation, and Malay knowledge representation by concepts. A number of experiments are carried out to quantify the accuracy of knowledge extracted. The Malay/K-Ontology is tested by using stratified random sampling drawn from various sources of Malay unstructured documents such as news, e-mails, articles, magazines, and texts from children story books. The results of the experiments have proved that the approach of Malay/K-Ontology performed well as compared to knowledge extracted manually done by an expert. The results of questionnaires evaluation on the Retrieval Interrogative Ontology Analysis Application have shown good achievement in understanding the main point of the unstructured document easily and clearly. This is to improve better understanding the process of making sense of information into knowledge, maintaining the meaning of the information and gaining the interpretation of the identical knowledge in unstructured document which facilitate identical knowledge perceived by different people.

Abstrak tesis yang dikemukakan kepada Senat Universiti Putra Malaysia
sebagai memenuhi keperluan untuk ijazah Doktor Falsafah

**TRANSFORMATION OF EXTRACTED KNOWLEDGE IN MALAY
UNSTRUCTURED DOCUMENTS INTO AN INTERROGATIVE
STRUCTURED FORM**

Oleh

FATIMAH BINTI SIDI

September 2007

Pengerusi: Profesor Madya Hj. Mohd. Hasan Selamat

Fakulti: Sains Komputer dan Teknologi Maklumat

Dengan adanya operasi penemuan pengetahuan, ianya telah membantu perolehan informasi dan pengetahuan yang berharga dalam saiz data yang besar dalam pangkalan data yang berstruktur. Walau bagaimanapun, sebahagian besar daripada informasi yang ada adalah berbentuk tidak berstruktur tetapi lebih kepada kumpulan dokumen-dokumen teks yang bentuknya tidak berstruktur, begitu juga dengan dokumen Melayu. Oleh itu, ciri-ciri pengstrukturkan perlu dilakukan ke atas dokumen yang tidak berstruktur bagi membolehkan informasi yang terdapat dalam dokumen yang tidak berstruktur tersebut ditukar kepada pengetahuan. Satu pendekatan baharu telah diwujudkan untuk menukar pengetahuan yang diperolehi dalam dokumen Melayu yang tidak berstruktur dengan mengenal pasti, mengorganisasi dan menstruktur pengetahuan yang diperolehi ke dalam

struktur berbentuk interogatif. Pembangunan seni binanya adalah berdasarkan kepada pelaksanaan pengenalanpastian pengetahuan interogatif, informasi mengikut konteks interogatif, dan pengorganisasian dan pengstrukturian interogatif dengan perwakilan pengetahuan Melayu melalui konsep. Ia menggunakan korpus Bahasa Melayu, teori interogatif, serta model yang berorientasikan objek, ontologi, dan pangkalan data. Penyelidikan ini melibatkan pembangunan sistem berasaskan seni bina Malay/K-Ontology, yang penghasilannya diukur melalui prestasi dapatan semula kuantitatif dengan menggunakan metrik perolehan kembali dan metrik ketepatan. Pembangunan Aplikasi Analisis Dapatan Semula Interogatif Ontologi digunakan bagi menilai ketahanan tugas. Ianya dinilai melalui fungsi dan manfaat berkaitan dengan penggunaan informasi mengikut konteks interogatif dengan penambahan pengekodan warna, tambahan anotasi informasi, dan dengan perwakilan pengetahuan Melayu melalui konsep. Beberapa uji kaji telah dijalankan untuk mengenal pasti jumlah ketepatan pengetahuan yang diperolehi. Malay/K-Ontology diuji dengan menggunakan pensampelan rawak strata yang sumbernya adalah daripada pelbagai dokumen Melayu yang tidak berstruktur seperti surat khabar, e-mel, artikel, majalah, dan teks daripada buku cerita kanak-kanak. Keputusan daripada uji kaji telah membuktikan bahawa pendekatan Malay/K-Ontology telah menunjukkan prestasi yang baik jika dibandingkan dengan pengetahuan yang diperolehi secara manual oleh pakar. Keputusan daripada penilaian soal selidik melalui Aplikasi Analisis Dapatan Semula Interogatif Ontologi telah menunjukkan pencapaian yang baik dalam

memahami perkara utama dokumen yang tidak berstruktur dengan mudah dan jelas. Ini adalah untuk menambah lagi pemahaman proses boleh terima informasi menjadi pengetahuan, pengekalan makna informasi dan mendapat tafsiran pengetahuan yang sama dalam dokumen yang tidak berstruktur bagi memudahkan pengetahuan yang sama diamati oleh orang yang berbeza.



ACKNOWLEDGEMENTS

In the name of *Allah*, the most merciful and most compassionate. Praise to *Allah s.w.t* for granting me strength, courage, patience and inspiration in completing this work. First and foremost, I owe deep gratitude to Assoc. Prof. Hj. Mohd. Hasan Selamat, Chairman of the supervisory committee for his constant guidance and continuous encouragement without which this work would not materialized. I am also thankful to co-supervisors, Assoc. Prof. Dr. Abdul Azim Abd. Ghani, Assoc. Prof. Dr. Hj. Md. Nasir Sulaiman for sharing their intellectual experience. I am also fortunate to have the support and guidance from Dr. Mokhtar Mohd. Yusof of Malaysian Ministry of Health for his support in the early chapters of my work. To Prof. Dr. Hj. Awang Sariyan from Academy of Malay Studies, Universiti Malaya, I am very grateful for his expertise and intellectual experience in verifying Malay language for this research. Last but not least, thanks to Assoc. Prof. Dr. Hamidah Ibrahim for her valuable advice.

I am also indebted to Universiti Putra Malaysia for the study leave, scholarship and allowances which enable me to pursue this research. My gratitude is also extended to my parents, family, colleagues and friends, thank you for sharing all the pains and gains throughout the years. My special thanks and gratitude to my beloved husband Abdul Rahman Hassan and my five children Ibrahiem Luqman, Aqilah, Mahirah, Aniqah and Aqidah for their understanding, caring, support and patience.

TABLE OF CONTENTS

		Page
DEDICATION		ii
ABSTRACT		iii
ABSTRAK		v
ACKNOWLEDGEMENTS		viii
APPROVAL		ix
DECLARATION		xi
LIST OF TABLES		xv
LIST OF FIGURES		xvii
LIST OF ABBREVIATIONS		xix
CHAPTER		
1	INTRODUCTION	1
	1.1 Background	1
	1.2 Problem Statement	3
	1.3 Research Objectives	6
	1.4 Scope of the Research	6
	1.5 Contributions of the Research	8
	1.6 Organization of the Thesis	9
2	LITERATURE REVIEW	12
	2.1 Introduction	12
	2.2 Spectrum of Data, Information and Knowledge	12
	2.3 Organization of Knowledge	17
	2.4 Extraction and Mapping of Knowledge	21
	2.5 Structured Schema	25
	2.6 Related Work Review Summary and Research Direction	27
	2.6.1 Document Collection	29
	2.6.2 Document Structuring	30
	2.6.3 Discovering Relationships	30
	2.6.4 Limitation and Issues Raised in Related Research	31
	2.7 Summary	36
3	RESEARCH METHODOLOGY	39
	3.1 Introduction	39
	3.2 Research Design	39
	3.2.1 Establishing Research Questions	42
	3.2.2 Concept Development Program	43
	3.2.3 Concept Implementation and Evaluation	43
	3.3 Summary	56

4	MALAY/K-ONTOLOGY APPROACH FOR TRANSFORMATION OF EXTRACTED KNOWLEDGE	57
4.1	Introduction	57
4.2	Theoretical Foundation of the Malay/K-Ontology Approach	58
	4.2.1 Interrogative Knowledge Identification	59
	4.2.2 Interrogative Knowledge Organization and Structuring	63
	4.2.3 Knowledge Structure	66
	4.2.4 Development of Malay/K-Corpus	71
	4.2.5 Development of a Stop Word List	75
4.3	System Design	78
4.4	System Implementation	81
	4.4.1 IKL-Identifier	82
	4.4.2 IKO-Recognizer	88
	4.4.3 IKS-OntologyDB	98
4.5	Summary	103
5	RESULTS AND DISCUSSIONS FOR MALAY/K-ONTOLOGY APPROACH	105
5.1	Introduction	105
5.2	Evaluation Approach	106
	5.2.1 Data Analysis Methods	111
	5.2.2 Data Used	113
	5.2.3 Significance Test	114
5.3	Results of Ontos and Malay/K-Ontology	116
	5.3.1 Analysis of Results of Ontos on English and Malay Obituaries	124
	5.3.2 Analysis of Results of Malay/K-Ontology on Malay Obituaries	128
5.4	Results of Malay/K-Ontology on Knowledge Extraction	131
	5.4.1 Experiment on the Usage of Lexicons	133
	5.4.2 Experiment on the Usage of Phrases	137
	5.4.3 Analysis of Results of Interrogative Knowledge Identification	143
	5.4.4 Analysis of Results of Interrogative Knowledge Organization and Structuring	155
5.5	Summary	164
6	RESULTS AND DISCUSSIONS FOR RETRIEVAL INTERROGATIVE ONTOLOGY ANALYSIS APPLICATION	168
6.1	Introduction	168
6.2	System Architecture	169
6.3	Effects of Interrogative Contextual Information	178
6.4	Effects of Malay Knowledge Representation	182
6.5	Summary	184

7	CONCLUSIONS AND FUTURE WORKS	188
7.1	Research Conclusion	188
7.2	Limitations	194
7.3	Future Works	195
	REFERENCES	197
	APPENDICES	205
	BIODATA OF THE AUTHOR	229
	LIST OF PUBLICATIONS	230



LIST OF TABLES

Table		Page
2.1	Comparison of data, information and knowledge	13
2.2	Epistemology and Methodology Summarization	28
2.3	Summary of Strengths and Weaknesses	31
3.1	The Malay unstructured documents total number of words	48
4.1	Examples of Malay/IK-Corpus	74
4.2	A list of the 35 most frequently occurring words	76
4.3	Example of <i>IKL</i> -Identifier Output	88
4.4	Objects populated by <i>Struktur Kata Nama Am</i>	93
4.5	Objects populated by <i>Struktur Leksikon Semantik</i>	96
4.6	Ontology versus SQL declaration properties	99
5.1	List of data extracted for DeceasedPerson	120
5.2	List of data extracted for DeceasedPersonRelationshipRelativeName	121
5.3	List of data extracted for Viewing	123
5.4	Results of Ontos on English and Malay obituaries	124
5.5	Results of Ontos and Malay/IK-Ontology	128
5.6	Results of lexicons for interrogative lexical constructs	133
5.7	Average of lexicons for interrogative lexical constructs	134
5.8	Results of lexicons for ontological constructs	136
5.9	Average of lexicons for ontological constructs	136
5.10	Results of phrases for interrogative lexical constructs	138
5.11	Average of phrases for interrogative lexical constructs	139
5.12	Summary of the evaluation results of phrases	139

5.13	Results of phrases for ontological constructs	139
5.14	Average of phrases for ontological constructs	140
5.15	Comparison on the usage of lexicons and phrases for interrogative lexical constructs	144
5.16	Number of interrogative lexical constructs extracted by an expert	147
5.17	Number of interrogative lexical constructs generated by <i>IKL-Identifier</i>	147
5.18	The differences between the number of interrogative lexical constructs by <i>IKL-Identifier</i> and expert	148
5.19	Frequency distribution of the direction of differences between <i>IKL-Identifier</i> and an expert (with <i>p</i> values)	152
5.20	Comparison on the usage of lexicons and phrases for ontological constructs	156
5.21	Number of ontological constructs extracted by an expert	158
5.22	Number of ontological constructs generated by <i>IKO-Recognizer</i>	159
5.23	The differences between the number of constructs by <i>IKO-Recognizer</i> and an expert	160
5.24	Frequency distribution of the direction of differences between <i>IKO-Recognizer</i> and expert (with <i>p</i> values)	162

LIST OF FIGURES

Figure	Page
2.1 The Conceptual Transformation of Unstructured Documents Framework	29
3.1 Research Orientation and Operational Framework	41
4.1 Interrogative Knowledge Identification Platform	61
4.2 Protégé-Frame Classes Tab for Building Concepts	67
4.3 Protégé-Frame Forms Tab for Knowledge Acquisition	68
4.4 Protégé-Frame Instances Tab for Building Instances	69
4.5 Protégé-Frame Queries Tab for Querying Instances	70
4.6 The Malay/ <i>I</i> K-Ontology Model	80
4.7 System Architecture of Malay/ <i>I</i> K-Ontology	81
4.8 <i>I</i> KL-Identifier Processes	83
4.9 <i>I</i> KO-Recognizer Processes	90
4.10 Example of ' <i>siapa</i> ' (who) object and its hierarchical classification	92
4.11 Example of Class ' <i>Benda</i> ' in HTML format	99
4.12 Sample of object-relationship model instance	101
5.1 Evaluation approach for Malay/ <i>I</i> K-Ontology accuracy	108
5.2 Precision and recall for a given unstructured document	109
5.3 Ontos High-Level Data-Extraction processes	117
5.4 Average Recall-Precision Chart for interrogative lexical constructs	144
5.5 Average Recall-Precision Chart for ontological constructs	156
6.1 System Architecture of the Retrieval Interrogative Ontology Analysis Application	169
6.2 Malay/ <i>I</i> K-Database collection and integration with Protégé knowledge-base system	171

6.3	Example of additional information annotation	173
6.4	Example of interrogative contextual information on interrogative element of ' <i>apa</i> ' (what)	174
6.5	Example of interrogative contextual information on interrogative element of ' <i>mengapa</i> ' (why)	174
6.6	Participant's satisfaction level on the effects of interrogative contextual information	180
6.7	Participant's satisfaction level on the effects of knowledge representation	183



LIST OF ABBREVIATIONS

CF	Concept Frames
CFG	Concept Frame Graph
DAML+OIL	Darpa Agent Markup Language
DML	Data Manipulation Language
FLogic	Frame Logic
HTML	Hypertext Markup Language
KDD	Knowledge Discovery in Databases
KIF	Knowledge Interchange Format
KRSs	Knowledge Representation Systems
NER	Name Entity Recognition
NLP	Natural Language Processing
OCML	Operational Conceptual Modeling Languages
OIL	Ontology Interface Layer
OKBC	Open Knowledge Base Connectivity
OOP	Object-Oriented Programming
OSM	Object-oriented System Model
RDF	Resources Description Framework
RDF(S)	RDF Schema
SHOE	Single HTML Ontology Extension
SQL	Structured Query Language
XML	eXtended Markup Language
XOL	XML-based Ontology Languages

CHAPTER 1

INTRODUCTION

1.1 Background

The difficulty of defining knowledge in unstructured documents is due to the paradox that knowledge resides in a person's mind and at the same time, it has to be captured, stored, and reported. For that, philosophers classify knowledge into knowing-that and knowing-how. Knowing-that is factual where data are stored in databases and facts can be recalled, processed, and disseminated. While knowing-how is actionable to do something, turning data into information and in turn into knowledge (Spiegler, 2003).

Recent advancements in computer technologies and databases have given many approaches to generate and extract knowledge. Among them is data mining and many people treat data mining as a synonym for another popularly used term, Knowledge Discovery in Databases (KDD). KDD is the field that is evolving to provide automated analysis solutions in extracting and generating knowledge. It is the process of identifying valid, novel, potentially useful, and ultimately understandable patterns/models in data (Fayyad *et al.*, 1996). Moreover, it has emerged as a rapidly growing interdisciplinary field that merges together databases, statistics, machine learning and related areas in order to extract valuable information and knowledge in large

volumes of data. KDD has deeply transformed the methods to interrogate traditional databases, where data are in structured form, by automatically finding new and unknown patterns in huge quantity of data. Most previous work in knowledge discovery is concerned with structured, numerical, heterogeneous databases and data warehouses, that will be extracted from the operational (day-to-day processing) systems (Poe, 1996; Iritano and Ruffolo, 2001; Kroeze *et al.*, 2003). Much work has been done in developing and building knowledge discovery systems by using KDD process where the source data are from operational database (Buchheit *et al.*, 2000; Chen *et al.*, 2001; Ho *et al.*, 2001; Tsai and Chen, 2001; Valafar and Valafar, 2002; Leung *et al.*, 2003).

However, structured data represent only a little part of the overall organization of knowledge; in fact, the major part of this knowledge is incorporated in textual documents. For example, available business data are captured in text files that are not structured, e.g. memoranda and journal articles that are available electronically (Fayyad *et al.*, 1996; Iritano and Ruffolo, 2001; Kroeze *et al.*, 2003). A large portion of the available information does not appear in structured databases but rather in collections of text articles drawn from various sources (Feldman, 1999). Thus, the main concern here is to dig knowledge from the available vast amount of textual documents.

1.2 Problem Statement

It is estimated that 90% of electronically available material is unstructured and the amount of unstructured textual documents, accessible through the web, intranets, news groups, etc. is enormously increased every year (Iiritano and Ruffolo, 2001). Hence, huge amount of unstructured documents are available on the web and intranets. The amount of information available to us is constantly increasing and our ability to absorb and process this information remains constant. "We are being drowned in information while being starved of knowledge and distracted from wisdom", taken from Dr. Norman Myers cited in Feldman (1999). Knowledge exists and is found everywhere (ubiquitous) in unstructured documents, so extracting knowledge in unstructured documents is essential.

Unstructured documents cannot be queried in simple ways. Therefore, knowledge contained in unstructured documents can neither be used by automatic systems nor could be understood easily and clearly by humans. Hence, identifying knowledge from unstructured documents to be easily realized and understood by humans is one of the most exciting areas to be explored.

Most work on knowledge discovery is concerned with structured databases. It is clear that this paradigm requires handling huge amount of information that is available in unstructured documents. To apply traditional knowledge

discovery or query operation on unstructured documents, it is necessary to impose some structures that will be rich enough to allow knowledge discovery operations techniques such as data mining to play their roles.

A structured database is a collection of related pieces of information stored electronically with structural description of the type of facts held in it. The most common model for a structured database is a relational model. It is used to represent any relationship between any collections of data that can be represented. It composes of tuples or records, and attributes or fields which unstructured document lacks. Embley *et al.* (1998a; 1998b; 1999a; 1999b) established and developed an approach of extracting information from unstructured documents and reformulating the information as relations in a database. The purpose of their work is to impose structure by establishing relations over the information contents of an unstructured document. The approach is based on data extraction ontology, a conceptual model instance that describes the data interest, including relationships, lexical appearance and context keywords. Later, Embley (2004) has extended the use of information extraction ontologies as an approach that leads to semantic understanding based on a foundation of Medows's definitions for data, information, knowledge and meaning. It is being reported that their approach generates precision ratios near 98% in extracting data on unstructured documents that are data rich. However, the work done does not attempt to extract "deep-level understanding" and does not depend upon complete sentences. The approach used is more appropriate for web pages

that publish information (like classified ads) which rarely contains complete sentences. Moreover, their approach is in extracting data on information of unstructured document; i.e., transforming information into data but does not focus on transforming information into knowledge. Unfortunately, nowadays a large part of knowledge is stored in an unstructured document or textual format which is usually written in complete sentences.

The growth in the number of unstructured documents written in Malay language is enormously available on the web and intranets. This triggers the need to investigate the availability of knowledge in Malay unstructured documents. The investigation looks on transforming information into knowledge. Hence, there is a need to identify the information in unstructured documents that have knowledge. The identified knowledge need to be extracted and then to be transformed into structured form by imposing structuring characteristic over the contents of the unstructured document. This is to enable the purpose of querying by using database standard Data Manipulation Language (DML), and increase understanding by humans on the main point of the unstructured document. At present, there are researches done in information retrieval in Malay language (Abu Ata, 1994; Ahmad, 1995; Tan, 1998; Abdullah, 2006). Unfortunately, these researches are focused on its information retrieval effectiveness but not on knowledge extraction.