



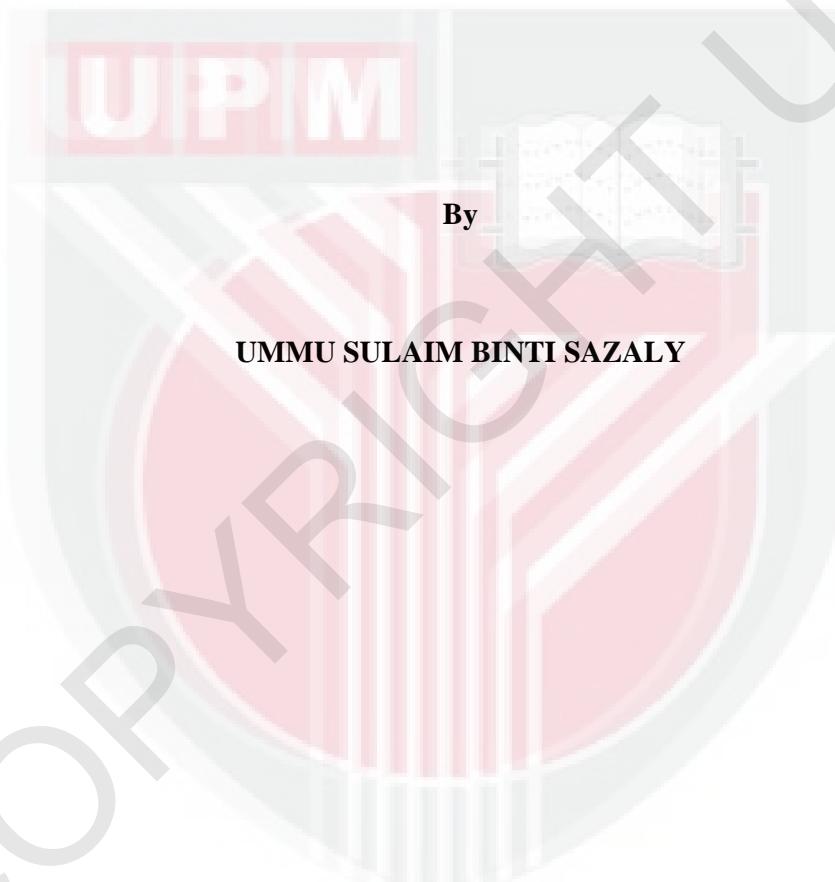
UNIVERSITI PUTRA MALAYSIA

***ENHANCE EFFICIENCY OF ANSWERING XML KEYWORD QUERY
USING INCOMPACT STRUCTURE OF MCCTREE***

UMMU SULAIM BINTI SAZALY

FSKTM 2013 3

**ENHANCE EFFICIENCY OF ANSWERING XML KEYWORD QUERY
USING INCOMPACT STRUCTURE OF MCCTREE**



UMMU SULAIM BINTI SAZALY



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

November 2012

Abstract of thesis presented to the Senate of Universiti Putra Malaysia in fulfilment
of the requirement for the degree of Master of Science

**ENHANCE EFFICIENCY OF ANSWERING XML KEYWORD QUERY
USING INCOMPACT STRUCTURE OF MCCTREE**

By

UMMU SULAIM BINTI SAZALY

November 2012

Chair : Associate Professor Mohd Hasan bin Selamat

Faculty : Computer Science and Information Technology

People nowadays live in cyber life where everything can be done by just typing through keyboard and system will complete the process. As the interaction is done through online, data sharing is the most important service to send and deliver information. Extended Markup Language (XML) has been chosen as the most important data sharing medium as it is very friendly for human and machine to interpret. Due to the importance of it, many studies have been done to increase the effectiveness of retrieving information from XML file. Many notions and techniques have been introduced especially to process query of information.

Compact Lowest Common Ancestor (CLCA) and Maximal Compact Lowest Common Ancestor (MCLCA) implemented in algorithms named CGTreeGenerator and MCCTreeGenerator has been proven in returning an accurate result in answering XML keyword query. CGTreeGenerator compacted the XML tree by eliminating irrelevant nodes based on CLCA notion, which produced Compact Global Tree (CGTree). MCCTreeGenerator used CGTree to select subtree called Maximal

Compact Connected Tree (MCCTree) as query result based on MCLCA notion. However, the MCCTree cannot be used directly in its ranking method because calculation in ranking method used the structure of subtree as before it has been compacted. If the result cannot be used directly by the ranking method, the algorithm has an ineffective process. Moreover, if the ineffective process requires re-examining the original tree, the efficiency of the process of the algorithm will be reduced. This study is a response to these weaknesses. This study proposes a new algorithm, namely XMCCTreeGenerator, to enhance the efficiency of the CGTree-MCCTreeGenerator.

This study identifies the effective processes needed in producing XML query result using MCLCA notion and without compacting it. Those processes constructed XMCCTreeGenerator algorithm which will produce the same subtree as MCCTree but difference in its structure. This new returned subtree called Extended MCCTree (XMCCTree) can be used directly by the ranking method because it is in an incompact structure. An experiment is run using XML datasets available in XML Data Repository from University of Washington's website. Two files are selected which consist of different data structure and divided into three ranges of size. Keywords are manually randomly selected from the files and executed between three to five numbers of keyword.

Two prototypes are developed which implement CGTree-MCCTreeGenerator and XMCCTreeGenerator. Since this study focuses on efficiency of the algorithm, elapsed time for each execution is collected from the experiment. In conclusion, the

proposed XMCCCTreeGenerator is more efficient than the previous CGTree-MCCTreeGenerator in answering XML keyword query using MCLCA.



Abstrak tesis yang dikemukakan kepada Senat Universiti Putra Malaysia sebagai
memenuhi keperluan untuk ijazah Master Sains

**TINGKATKAN KECEKAPAN MENJAWAB PENCARIAN KATA KUNCI
XML MENGGUNAKAN STRUKTUR MCCTREE TIDAK PADAT**

Oleh

UMMU SULAIM BINTI SAZALY

November 2012

Pengerusi : Profesor Madya Mohd Hasan bin Selamat

Fakulti : Komputer Sains dan Teknologi Maklumat

Masyarakat kini hidup dalam dunia siber di mana semua perkara boleh dilakukan dengan hanya menaip pada papan kekunci dan sistem akan menyelesaikan proses tersebut. Oleh kerana interaksi diselesaikan atas talian, perkongsian data adalah servis yang penting dalam menghantar dan menerima maklumat. Lanjutan Bahasa Penanda (XML) telah dipilih sebagai medium perkongsian data yang utama kerana ianya mudah untuk ditafsirkan oleh manusia dan mesin. Oleh kerana kepentingan itu, banyak kajian telah dilakukan untuk meningkatkan keberkesanan dalam mendapatkan maklumat dari fail XML. Pelbagai kaedah dan teknik telah diperkenalkan terutamanya untuk memproses pertanyaan maklumat.

Aturan Am Padat Terendah (CLCA) dan Aturan Am Padat Maksimum Terendah (MCLCA) digunakan dalam algoritma CGTreeGenerator dan MCCTreeGenerator telah terbukti dalam mengembalikan keputusan yang tepat dalam menjawab pertanyaan kata kunci XML. CGTreeGenerator memadatkan pepohon XML dengan menghapuskan nod yang tidak relevan berdasarkan anggapan CLCA, yang

menghasilkan Pepohon Global Padat (CGTree). MCCTreeGenerator menggunakan CGTree untuk memilih pohon dipanggil Pohon Bersambung Padat Maksima (MCCTree) sebagai keputusan pertanyaan berdasarkan anggapan MCLCA. Walau bagaimanapun, MCCTree tidak boleh digunakan secara terus dalam kaedah penarafan kerana pengiraan dalam kaedah penarafan menggunakan struktur pohon sebelum ianya dipadatkan. Jika keputusan tidak boleh digunakan terus oleh kaedah penarafan, proses tersebut adalah tidak efektif. Tambahan pula, jika proses yang tidak efektif perlukan pemeriksaan semula pepohon asal, kecekapan proses akan berkurangan. Kajian ini adalah tindakbalas kepada kelemahan tersebut. Kajian ini mencadangkan algoritma baru, dinamakan XMCCTreeGenerator, untuk meningkatkan kecekapan algoritma CGTree-MCCTreeGenerator.

Kajian ini mengenalpasti proses-proses efektif yang diperlukan dalam menghasilkan keputusan pertanyaan XML menggunakan MCLCA serta tanpa memadatkannya. Kesemua proses itu membina algoritma XMCCTreeGenerator yang menghasilkan pohon yang sama seperti MCCTree dalam struktur yang berlainan. Pohon baharu yang dikembalikan ini dinamakan lanjutan MCCTree (XMCCTree) boleh digunakan terus oleh kaedah penarafan kerana ianya berada dalam struktur yang tidak padat. Eksperimen dijalankan menggunakan dataset XML di Pangkalan Data XML di laman sesawang Universiti Washington. Dua fail dipilih yang mengandungi struktur data berlainan dan dibahagikan kepada tiga saiz. Kata kunci dipilih rawak secara manual dari fail-fail tersebut dan diproses antara tiga sehingga lima kata kunci.

Dua prototaip dibangunkan menggunakan CGTree-MCCTreeGenerator dan XMCCTree. Oleh kerana kajian ini memberi tumpuan kepada kecekapan algoritma,

tempoh masa setiap pemprosesan diambil daripada eksperimen. Kesimpulannya, XMCCTreeGenerator yang dicadangkan adalah lebih berkesan berbanding CGTree-MCCTreeGenerator yang sebelumnya dalam menjawab pertanyaan kata kunci XML menggunakan MCLCA.



TABLE OF CONTENTS

	Page
ABSTRACT	ii
ABSTRAK	v
APPROVAL	viii
DECLARATION	x
LIST OF TABLES	xiii
LIST OF FIGURES	xiv
LIST OF ABBREVIATIONS	xvi
 CHAPTER	
1 INTRODUCTION	1
1.1 Research Background	1
1.2 Research Problem	4
1.3 Research Objective	5
1.4 Research Scope	5
1.5 Research Contribution	8
1.6 Thesis Organization	8
2 LITERATURE REVIEW	10
2.1 Introduction	10
2.2 Query in XML	11
2.3 Keyword Query in XML	15
2.3.1 Notion and Algorithm	16
2.3.2 Lowest Common Ancestor	18
2.3.3 CLCA and CGTree	22
2.3.4 Compact Connected Tree	26
2.3.5 MCLCA and MCCTree	27
2.4 Ranking Procedure	28
2.4.1 Scores on Structural Compactness	28
2.4.2 Scores on Text Similarity	32
2.4.3 Ranking the MCCTree	34
2.5 Summary	36
3 RESEARCH METHODOLOGY	38
3.1 Introduction	38
3.2 General Steps of Research Methodology	38
3.3 Research Proposal (Phase 1)	41
3.3.1 Literature Studies	41
3.3.2 Research Validation	42
3.4 Research Solution (Phase 2)	44
3.5 Research Experiment (Phase 3)	45
3.5.1 Objective of Experiment	45
3.5.2 XML Dataset	46

3.5.3	Instrumentation	47
3.5.4	Result	48
3.5.5	Data Analysis	49
3.6	Validity of the Result	50
3.6.1	Internal Validity	50
3.6.2	External Validity	51
3.6.3	Construct Validity	52
3.6.4	Conclusion Validity	52
3.7	Summary	53
4	XMCCTREEGENERATOR CONSTRUCTION	55
4.1	Introduction	55
4.2	Framework of XML Keyword Query	55
4.3	Analysis of Requirement	57
4.4	Identifying the Features	59
4.5	Algorithm Functionalities	61
4.6	XMCCTreeGenerator vs CGTreeGenerator and MCCTreeGenerator	69
4.6.1	Concept of CGTreeGenerator and MCCTreeGenerator	69
4.6.2	Concept of XMCCTreeGenerator	70
4.7	Process Visualization in XMCCTreeGenerator	71
4.8	Summary	80
5	RESULT AND DISCUSSION	82
5.1	Introduction	82
5.2	Experiment Remarks	82
5.3	Experiment Result	84
5.4	Performance Criteria	86
5.4.1	Performance of Both Prototypes in Executing a File - Range of 50 kb of File A	86
5.4.2	Performance of Both Prototypes in Executing a Different Size of File but With the Same Structure - File A	88
5.4.3	Performance of Both Prototypes in Executing a Different Structure of File – File B	90
5.5	Discussion	92
5.5.1	Number of Node Entering CMSet	94
5.5.2	Keyword Position in a File	94
5.5.3	Increase on the Size of the File	95
5.6	Summary	97
6	CONCLUSION AND FUTURE WORKS	99
6.1	Conclusion	99
6.2	Limitations and Weaknesses	101
6.3	Future Works	102
REFERENCES		103
BIODATA OF THE STUDENT		107
LIST OF PUBLICATIONS		108