

**Universiti Teknologi MARA**

**Automated Agriculture Terms Dictionary Using  
String Parsing**

**NUR SYAHIDA BINTI AHMAD ZUBAIDILLAH**

**Thesis submitted in fulfillment of the requirement for  
Bachelor of Computer Science (Hons)  
Faculty of Computer and Mathematical Sciences**

**JULY 2014**

## **ACKNOWLEDGEMENT**

Alhamdulillah, praise and thank to Allah because of His Almighty and His utmost blessing, I was able to finish this research within the time duration given. Firstly, my special thanks go to my supervisor, Dr. NurSuriati Binti Jamil because guide me from starting until the end. Special thanks to my FYP lecturer Dr. Siti Salwa Binti Salleh for teaching this subject.

A special appreciation also goes to my beloved parents Ahmad Zubaidillah Bin Othman and . . . Thanks for giving support.

Last but not least, I would like to give my gratitude to my dearest friend in CS230. Thank you for your helps.

## **ABSTRACT**

This paper proposed the task of generating terms dictionary in the Agriculture domain. Its purpose is to generate only important terms using in Agriculture fields from Hypertext Markup Language (.html) file and Word Document (.doc) file as the input. The goal is to create a prototype which can automatically read the string from this type of file and then extracts their terms. The limitation of the previous manual work for generating terms are known to be labor-intensive, limitation of applicability and time consuming. Terms extracted are the smallest fragments of texts in documents, rather than the entire documents that contain the texts. After analyzing the characteristics of terms in Agriculture Dictionary, we propose a string parsing method to deal with the task. We also design and develop a parser on identified parsing method. To test the functionality of the prototype, compared between the actual terms from document and the generated terms is determined. The final prototype achieved 100% accuracy between actual terms and generated terms from html document and achieves 86.27% accuracy from word document file. However, to extract terms is difficult because of the type of files that cannot easily be read by the default library package. For future work, Automated Agriculture Terms Dictionary will be fixed to get better result.

## Table of Contents

CONTENTS	PAGE
SUPERVISOR'S APPROVAL .....	ii
DECLARATION .....	iii
ACKNOWLEDGEMENT.....	iv
ABSTRACT .....	v
Table of Contents.....	vi
List of Figures.....	ix
List of Tables.....	x
CHAPTER 1.....	1
1.1 Background of Study.....	1
1.2 Problem Statement .....	2
1.3 Objectives .....	4
1.4 Scope .....	4
1.5 Significance of Study .....	5
1.6 Summary.....	5
CHAPTER 2.....	6
2.1 Introduction.....	6
2.2 Information Extraction .....	6
2.3 Domain Application .....	7
2.4 Parser .....	10
2.4.1 Syntactic Parsing .....	10
2.4.2 Semantic Parsing .....	11
2.4.3 Full Parsing .....	12
<b>2.4.3.1 Fully Parsed Tree using Context Free Grammar (CFG) .....</b>	<b>13</b>
<b>2.4.3.2 Parsing with Dependency Grammar .....</b>	<b>13</b>
2.4.4 String Parsing .....	14
2.5 Terms.....	16
2.5.1 Dictionary .....	16
2.5.2 Hypertext Markup Language (html) .....	16

2.5.3 Word Document (doc) .....	17
2.6 Current Development on Information Extraction .....	18
2.6.1 A Novel Use of Statistical Parsing to Extract Information from Text.....	18
2.6.2 Automated reference resolution in legal texts .....	19
2.6.3 Extracting an Arabic Lexicon from Arabic Newspaper Text.....	20
2.6.4 Learning Dictionaries for Information Extraction by Multi-Level Bootstrapping.....	21
2.7 Summary.....	21
CHAPTER 3.....	22
3.1 Introduction.....	22
3.2 Research Framework.....	22
3.3 Preliminary Study.....	25
3.4 Requirement Analysis .....	26
3.5 Prototype Design.....	27
3.5.1 Pre-Processing of Data.....	27
3.5.1.1 Read Document .....	28
3.5.1.2 Extract Document .....	29
3.5.1.2.1 String Parsing .....	29
<b>3.5.1.2.1.1 Pseudo code for Extract Html Document</b> .....	30
<b>3.5.1.2.1.2 Pseudo code for Extract Word Document</b> .....	31
3.5.2 Development of Prototype .....	34
3.5.2.1 Interface Design Development .....	36
3.5.2.2 Application Development.....	36
3.6 Coding and Debugging.....	37
3.7 Testing of Prototype.....	37
3.8 Summary.....	38
CHAPTER 4.....	39
4.1 Introduction.....	39
4.2 Sample Data Design .....	39
4.2.1 Hypertext Markup Language (html) data.....	39
4.2.2 Word Document (doc) data.....	40
4.3 System Design .....	41