

Universiti Teknologi MARA

**Topic Based Search For Malay
Translated Al-Quran
Documents**

Yuhaini Binti Yusop

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

JANUARY 2014

ACKNOWLEDGEMENT

Bismillahirrahmanirrahim

In the name of Allah. The most Gracious and Most Merciful

First and foremost, I would like to praise Allah, because of His Almighty and His utmost blessings, I finally able to finish this research within the time duration given. This project is fulfilment of the requirements for Bachelor of Computer Science (Hons). I am thankful for the opportunity to conduct this research and for the strength that Allah S.A.W gave me throughout the process.

Secondly, special thanks to my supervisor Puan Norasiah Mohammad for her outstanding guidance and valuable advice, who took keen interest for all my work all over the year, both in the implementation part and writing this thesis. I appreciated her time for helping, support and encouraging me throughout the project is conducted.

Thirdly, thanks to my beloved family for always supporting and trusting me whatever I do especially my parent. Thanks for the love and kindness and always be there for me, they are always in my heart.

Last but not least, I would like to give my gratitude to my dearest friend for the collaboration and for everything that they have done for me during this research is carried out.

ABSTRACT

Document classification retrieval is a method of selecting the right document that contains useful information. One of the data corpuses that can be stored in the document is Al-Quran Malay translation collection. Al-Quran is a holy book and can be extremely important document to Muslims. It is used as a guidance to lead a righteous life. Most people read the Al-Quran along with its translation to understand the content or the meaning of the verses in each of the Surah (chapters). Each chapter describes information that can be referred to by Muslims with the help of the translation. However, it is not an easy task to search for information in the Malay translated Al-Quran since it contains many pages. Most of people find it hard to search the phrases manually in the Al-Quran, and it require a lot time consuming. Moreover, most of the searches are not related to what the user want. Therefore, the objective of this project is to develop Malay search query prototype for the Malay translated Al-Quran document where the user can search the information based on topic. Besides that, this project will also focus on the Malay text language, and the data to be used are the Malay translation of chapter Al-‘Imran. Thus, to accomplish these, the technique of Latent Dirichet Allocation (LDA) is proposed to extract the topic in the document. The significance of this project is to ease the effort of people in finding specific topic in the Al-Quran and the related information from those topics. Moreover, it can also reduce the time to search for the word and to provide fast retrieval of data. The expected outcome of this project is Malay search query prototype for the Malay Translation Al-Quran Document to extract documents which has verses containing the topic.

Keywords— Document Classification, Malay Translated Al-Quran Document, Latent Dirichlet Allocation, Topic.

TABLE OF CONTENTS

CONTENT	PAGE
SUPERVISOR’S APPROVAL	ii
DECLARATION	iii
ACKNOWLEDGEMENT	iv
ABSTRACT	v
TABLE OF CONTENT	vi
LIST OF FIGURE	xi
LIST OF TABLES	xiii
CHAPTER ONE: INTRODUCTION	1
1.0 Introduction	1
1.1 Research background	5
1.2 Problem statement	6
1.3 Significant of study	7
1.4 Research Objective	8
1.5 Research Scope	8
1.5.1 Stakeholder	8
1.6 Research Element	9
1.7 Expected Outcome	9
1.8 Possible Solution	9
1.9 Summary	10
CHAPTER TWO: LITERATURE REVIEW	11
2.0 Introduction	11

2.1	Document classification	11
2.2	Document Representation	13
2.3	Malay Stemmer	14
	2.3.1 Ahmad's Stemmer	15
	2.3.2 Othman's Stemmer	16
	2.3.3 Idris' Stemmer	16
2.4	Topic Model	18
2.5	Common Technique In Document Classification	19
	2.3.1 Latent Semantic Indexing	19
	2.3.2 Latent Dirichlet Allocation	20
	2.3.3 Pachinko Allocation Model	21
2.6	Advantages and Disadvantages of the Technique	22
2.7	Overview of Latent Dirichlet Allocation	24
	2.7.1 Uses of Latent Dirichlet Allocation	25
2.8	Related Work	25
	2.8.1 Geometric Latent Dirichlet Allocation on a Matching Graph For Large Scale Image Dataset	25
	2.8.2 Latent Dirichlet Allocation for Text, Image and Music	26
	2.8.3 LDA-Based Document Models For Ad-Hoc Retrieval	26
	2.8.4 DiscLDA: Discriminative Learning for Dimensionality Reduction and Classification	27