



اَوْنُوْا سِيْقِيْ بِاَتِيْكَوْ لُوْ كِيْ فَا مَارَا
UNIVERSITI
TEKNOLOGI
MARA

Universiti Teknologi MARA

Search Engine for Books

Amierul Izzuddin Bin Azman

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

February 2022

DECLARATION

I certify that this report and the research to which it refers are the product of my own work and that any ideas or quotation from the work of other people, published or otherwise are fully acknowledged in accordance with the standard referring practices of the discipline.



.....

AMIERUL IZZUDDIN BIN AZMAN

2020968971

(29 January 2022)

ABSTRACT

Search engine is web-based software that searches for and identifies things in a collection that match the user's keywords or characters, and is mostly used to locate specific websites on the Internet. The basic idea of information retrieval (IR), as is well-known in the computer industry, is to search a given amount of data and retrieve those records that fulfil a set of criteria. For this project, it will be more on developing a search system that focus on the domain of books. Although many search engines available for book searching, most of them still have room for improvement. The reality to find books using search engine are quite challenging because, the users need to know the title of the book in order to retrieved relevant result from the search engine. Therefore, this project aims to propose a search engine that might produce better result by manipulating the indexing structure of the search engine. Software Development Life Cycle also known as SDLC was used as the methodology in this project development. This project applies the vector space model for the matching process, because the chosen software library for this project which is Apache Lucene is using the vector space model as its foundation. In addition, the Bag-of-Words approach was used as the basis for the indexing module. The indexing process indexed the information of the books such as the book title, the author name, publisher, year published, pages count and synopsis of the book. The search engine's assessment criteria include recall and precision. In IR, recall and precision have long been used as standard evaluation criteria. Keenly, the results of this project prove that index files that contain more domain information can increase the relevancy of a search engine.

TABLE OF CONTENTS

| CONTENTS | PAGE |
|---------------------------------------|-------------|
| SUPERVISOR'S APPROVAL | ii |
| DECLARATION | iii |
| ACKNOWLEDGEMENT | iv |
| ABSTRACT | v |
| TABLE OF CONTENTS | vi |
| LIST OF FIGURES | vii |
| LIST OF TABLES | viii |
| | |
| CHAPTER ONE: INTRODUCTION | |
| 1.1 Background of Study | 1 |
| 1.2 Problem Statement | 2 |
| 1.3 Project's Questions | 3 |
| 1.4 Project's Objective | 3 |
| 1.5 Project's Scope | 3 |
| 1.6 Significance | 4 |
| 1.7 Conclusion | 4 |
| | |
| CHAPTER TWO: LITERATURE REVIEW | |
| 2.1 Introduction | 5 |
| 2.2 Book | 6 |
| 2.3 Search Engine Components | 7 |
| 2.3.1 Crawling | 8 |
| 2.3.2 Indexing | 9 |
| 2.3.3 Query | 9 |

| | |
|---|----|
| 2.3.4 Matching | 10 |
| 2.3.5 Ranking | 11 |
| 2.4 Search Engine Model | 12 |
| 2.4.1 Boolean | 12 |
| 2.4.2 Probabilistic | 13 |
| 2.4.3 Vector-Space Model | 13 |
| 2.5 Search Engine Evaluation | 13 |
| 2.5.1 Recall | 14 |
| 2.5.2 Precision | 14 |
| 2.5.3 Mean Average Precision | 15 |
| 2.5.4 F-measure | 15 |
| 2.6 Existing Application Related to Search Engine | 16 |
| 2.6.1 Google | 16 |
| 2.6.2 Yahoo! | 17 |
| 2.6.3 Bing | 18 |
| 2.7 Related Research | 19 |
| 2.8 Justification of the Chosen Model, Component, Evaluation and Features | 22 |
| 2.8.1 Chosen Model of Search Engine | 22 |
| 2.8.2 Chosen Component of Search Engine | 22 |
| 2.8.3 Chosen Evaluation of Search Engine | 22 |
| 2.8.4 Features of Search Engine | 23 |
| 2.9 Summary | 23 |

CHAPTER THREE: METHODOLOGY

| | |
|----------------------------------|----|
| 3.1 Introduction | 24 |
| 3.2 Operational Framework | 24 |
| 3.2.1 Planning | 25 |
| 3.2.2 Information Gathering | 26 |
| 3.2.3 Implementation Development | 27 |