**TOPIC MODELING FOR**

**INDONESIAN TEXT**

**UNDERGRADUATE THESIS**

Submitted to the Faculty of Computer Science Universitas Amikom Yogyakarta as a
partial fulfillment of the requirement for bachelor degree

**By**

**Rasyiid Indra Parmadi**

**15.61.0042**

**BACHELOR DEGREE**
**STUDY OF INFORMATICS**
**FACULTY OF COMPUTER SCIENCE**
**UNIVERSITAS AMIKOM YOGYAKARTA**
**YOGYAKARTA**
**2019**

# APPROVAL

## THESIS

### TOPIC MODELING FOR
### INDONESIAN TEXT

prepared and compiled by

**Rasyiid Indra Parmadi**

**15.61.0042**

has been approved by the Thesis Supervisor

on February 9, 2019

Supervisor,

Arief Setyanto, Dr., S. Si, M. T.

NIK. 190302036

# ATTESTATION

## THESIS

### TOPIC MODELING FOR
### INDONESIAN TEXT

prepared and compiled by

**Rasyiid Indra Parmadi**

**15.61.0042**

has been maintained in front of the Board of Examiners

on February 28, 2019
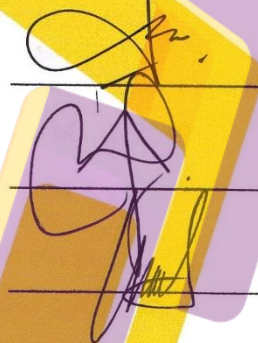
**Board of Examiners**

| Examiners | Signature |
|---|---|
| **Sudarmawan, S. T., M. T.** | |
| NIK. 190302035 | |
| **Andi Sunyoto, M. Kom.** | |
| NIK. 190302052 | |
| **Ike Verawati, M. Kom.** | |
| NIK. 190302237 | |

This thesis has been accepted as one of the requirements to obtain a Bachelor of

Computer degree on March 19, 2019

**DEAN OF THE FACULTY OF COMPUTER SCIENCE**

**Krisnawati, S. Si, M. T.**

**NIK. 190302038**

# STATEMENT

I, the undersigned, hereby declare that, this thesis is my own (ORIGINAL) work, and the contents of this thesis do not have works submitted by others to obtain an academic degree at any higher education institution, and as long as my knowledge is not works or opinions that have been written and / or published by others, except those written in this text and mentioned in the bibliography.

Everything related to the manuscript and the work that has been made is my personal responsibility.

Yogyakarta, 3 March 2019

Rasyiid Indra Parmadi

NIM. 15.61.0042

# MOTTO

"Il est bon à savoir, it is good to know."

*(Den Dhimas)*

*"Do the best and pray. God will take care of the rest."*

*(Den Dhimas)*

*"Use your youth as good as possible."*

*(Den Dhimas)*

*"Everything will be okay in the end, if its not okay, its not the end."*

*(Den Dhimas)*

# DEDICATION

I thank Allah God Almighty for giving His blessings, mercy and guidance so that I can finish this Thesis well. I also feel grateful to the people around me who have directly or indirectly helped me in working on this Thesis. . I present this thesis to :

1. My father, Suparmadi, My mother, Murni Indrawati, and one and only sister Zharifah Indra Parmadi who always prayed for, encouraged, and supported me.

2. Mr Arief Setyanto, Dr., S. Si, M. T. as a supervisor who always provides input and guidance in completing the Thesis.

3. To Raka Wichaksana and Ripto Sudiyarno as a mentor who has given his knowledge to help conduct research.

4. My WTF friends, who always give encouragement and motivation to immediately complete the Thesis.

5. My friends in playing game who always help entertain and help in distress.

6. 15BCI classmates, AMIKOM assistants and organizations that have become my friends during college.

As well as all those who have helped and supported me that I cannot mention one by one.

# ACKNOWLEDGEMENTS

Thank you we pray to Allah SWT for His blessings and gifts so that the writer can complete the thesis report in time with the title "TOPIC MODELING FOR INDONESIAN TEXT" This thesis was prepared to complete the final assignment of college and fulfill the graduation requirements of the Bachelor Informatics Education program at Amikom University Yogyakarta. During the education of Bachelor Informatics up to the Thesis completion process, various parties have provided facilities, assisted, fostered, and guided the writer for that especially to:

1. Mr. Prof. Dr. M. Suyanto, MM as Chancellor of the Amikom University in Yogyakarta who has provided many facilities in completing education..

2. Mr. Arief Setyanto, Dr., S. Si, M. T. as a supervisor who has spent a lot of time and energy guiding the author during the preparation of this thesis..

3. Mr / Mrs. Lecturer at Amikom University Yogyakarta who has provided writers with several useful disciplines..

4. Friends in arms of Bachelor Informatics Students in 2015, who have discussed and collaborate with writers during their education..

The author realizes, this thesis still has many weaknesses and weaknesses. Therefore, constructive criticism and suggestions will be welcomed, hopefully, the existence of this thesis can be useful and increase our insight, especially about Web Security.

Yogyakarta, 3 March 2019

Author

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# ABSTRACT

*Reading is the activity of perceptual, analyzing, and interpreting what is done by the reader to get the message to be conveyed by the author in the writing media. Along with current technological developments, technology has changed the way a person reads a text, especially text that is in electronic media. Generally, a reader must read the text that is in a particular file thoroughly to find out what topics can be obtained from the text. This becomes a problem in terms of the time and magnitude of the effort used by the reader.*

*This research will make a document summarization program, especially in Indonesian text using the Maximum Marginal Relevance or MMR algorithm. The MMR algorithm is a simple and efficient text summarization algorithm. With text summarizing programs, change someone to find and find out what topics are available in the text without having to read the text thoroughly.*

*This program will produce topics or conclusions from a text and help someone to find out what topics are in the text without reading it. This program will help someone who doesn't like the process of reading a lot of text. After getting the topic, he will look for other explanations about the topic from the internet to learn manually.*

**Keywords : Summarization, Indonesian Text, Maximum Marginal Relevance (MMR).**