

**BAHASA MELAYU DIALECT TRANSLATOR AND DETECTOR
(MALAYFY)**



UNIVERSITI TEKNIKAL MALAYSIA MELAKA

BORANG PENGESAHAN STATUS LAPORAN

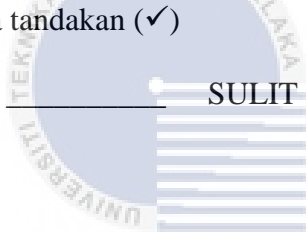
JUDUL: BAHASA MELAYU DIALECT TRANSLATOR AND DETECTOR (MALAYFY)

SESI PENGAJIAN: 2022 / 2023

Saya: ADELLA JAVA DIRGANTARI

mengaku membenarkan tesis Projek Sarjana Muda ini disimpan di Perpustakaan Universiti Teknikal Malaysia Melaka dengan syarat-syarat kegunaan seperti berikut:

1. Tesis dan projek adalah hakmilik Universiti Teknikal Malaysia Melaka.
2. Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dibenarkan membuat salinan untuk tujuan pengajian sahaja.
3. Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dibenarkan membuat salinan tesis ini sebagai bahan pertukaran antara institusi pengajian tinggi.
4. * Sila tandakan (✓)



_____ SULIT

(Mengandungi maklumat yang berdarjah keselamatan atau kepentingan Malaysia seperti yang termaktub di dalam AKTA RAHSIA RASMI 1972)

_____ TERHAD

(Mengandungi maklumat TERHAD yang telah ditentukan oleh organisasi / badan di mana penyelidikan dijalankan)

_____ ✓ _____ TIDAK TERHAD

(TANDATANGAN PELAJAR)

Alamat tetap: JL. CANDI AGUNG 1/20
A, MALANG, JAWA TIMUR,
INDONESIA, 65142

(TANDATANGAN PENYELIA)

Ts. DR. HALIZAH BINTI BASIRON

Nama Penyelia

Tarikh: 13/09/2023

Tarikh: 13/09/2023

CATATAN: * Jika tesis ini SULIT atau TERHAD, sila lampirkan surat daripada pihak

BAHASA MELAYU DIALECT TRANSLATOR AND DETECTOR
(MALAYFY)

ADELLA JAVA DIRGANTARI



This report is submitted in partial fulfillment of the requirements for the
Bachelor of Computer Science (Artificial Intelligence) with Honours.

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

FACULTY OF INFORMATION AND COMMUNICATION TECHNOLOGY
UNIVERSITI TEKNIKAL MALAYSIA MELAKA

2023

DECLARATION

I hereby declare that this project report entitled
BAHASA MELAYU DIALECT TRANSLATOR AND DETECTOR
(MALAYFY)

is written by me and is my own effort and that no part has been plagiarized
without citations.

STUDENT :  Date : 17/09/2023
ADELLA JAVA DIRGANTARI


I hereby declare that I have read this project report and found

this project report is sufficient in term of the scope and quality for the award of
Bachelor of Computer Science (Artificial Intelligence) with Honours.

SUPERVISOR :  Date : 17/09/2023
Ts. DR. HALIZAH BINTI BASIRON

DEDICATION

This final-year study is dedicated to my supervisor, Ts. Dr. Halizah binti Basiron, for her support and guidance. I also want to thank my parents and siblings for their financial and emotional support in helping me complete this project. Last but not least, I want to thank my friends for their constant encouragement and emotional support when I was a university student.



ACKNOWLEDGEMENTS

I would like to express my gratitude to my supervisor, Ts. Dr. Halizah binti Basiron, for her invaluable support and guidance throughout the project. She provided me with valuable feedback on my research methods, identified areas of improvement, and provided me with the necessary resources to successfully complete the project. I would also like to thank my beloved parents who have also been a constant source of support and encouragement, and I am extremely grateful for their unwavering belief in me. Additionally, I would like to extend our sincere appreciation to my seniors and friends who provided me with valuable insights and feedback that helped me to improve my project.



ABSTRACT

The Bahasa Melayu Dialect Translator and Detector (Malayfy) system is an innovative solution designed to overcome language barriers and facilitate effective communication between users of different dialects in the context of Bahasa Melayu. This system consists of three modules: a language translator from English to Bahasa Melayu, a Bahasa Melayu dialect translator, and a Bahasa Melayu dialect detector. The project's primary goal is to create a reliable and precise system capable of translating English text into Bahasa Melayu while also translating standard Bahasa Melayu into dialects. Additionally, the dialect detection module employs the Random Forest algorithm, renowned for its ability to handle complex classification tasks, to predict the dialect of user input sentences. The project encompasses various stages, including a literature review, data collection, system design, and implementation. The system's effectiveness is assessed based on its accuracy, efficiency, and user satisfaction. The current issues associated with the absence of dialect detection and translation in existing machine translation tools for Bahasa Melayu are anticipated to be addressed with the development of this system. Effective communication with Bahasa Melayu speakers is crucial in various fields, including tourism, education, healthcare, and business. The successful implementation of this system, driven by the power of machine learning algorithms like Random Forest, could have significant implications in these fields and others.

ABSTRAK

Sistem Penterjemah dan Pengesan Dialek Bahasa Melayu (Malayfy) ialah penyelesaian inovatif yang direka untuk mengatasi halangan bahasa dan memudahkan komunikasi berkesan antara pengguna dialek berbeza dalam konteks Bahasa Melayu. Sistem ini terdiri daripada tiga modul: penterjemah bahasa daripada Bahasa Inggeris ke Bahasa Melayu, penterjemah dialek Bahasa Melayu, dan pengesan dialek Bahasa Melayu. Matlamat utama projek ini adalah untuk mencipta sistem yang boleh dipercayai dan tepat yang mampu menterjemah teks Inggeris ke dalam Bahasa Melayu sambil juga menterjemah Bahasa Melayu standard ke dalam dialek. Selain itu, modul pengesanan dialek menggunakan algoritma Random Forest, yang terkenal dengan keupayaannya untuk mengendalikan tugas pengelasan yang kompleks, untuk meramalkan dialek ayat input pengguna. Projek ini merangkumi pelbagai peringkat, termasuk kajian literatur, pengumpulan data, reka bentuk sistem dan pelaksanaan. Keberkesanan sistem dinilai berdasarkan ketepatan, kecekapan dan kepuasan pengguna. Isu semasa yang berkaitan dengan ketiadaan pengesanan dialek dan terjemahan dalam alat terjemahan mesin sedia ada untuk Bahasa Melayu dijangka dapat ditangani dengan pembangunan sistem ini. Komunikasi yang berkesan dengan penutur Bahasa Melayu adalah penting dalam pelbagai bidang, termasuk pelancongan, pendidikan, penjagaan kesihatan dan perniagaan. Kejayaan pelaksanaan sistem ini, didorong oleh kuasa algoritma pembelajaran mesin seperti Random Forest, boleh mempunyai implikasi yang ketara dalam bidang ini dan lain-lain.

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

TABLE OF CONTENTS

| | PAGE |
|-------------------------------------|-------------|
| DECLARATION..... | II |
| DEDICATION..... | III |
| ACKNOWLEDGEMENTS..... | IV |
| ABSTRACT | V |
| ABSTRAK | VI |
| TABLE OF CONTENTS..... | VII |
| LIST OF TABLES | XII |
| LIST OF FIGURES | XIII |
| LIST OF ABBREVIATIONS | XVI |
| LIST OF ATTACHMENTS..... | XVII |
| CHAPTER 1: INTRODUCTION..... | 1 |
| 1.1 Introduction..... | 1 |
| 1.2 Problem Statement(s)..... | 2 |
| 1.3 Objective(s)..... | 2 |
| 1.4 Project Scope | 2 |
| 1.4.1 Module(s)..... | 2 |
| 1.4.2 Target User(s) | 3 |
| 1.5 Project Significant..... | 4 |
| 1.6 Expected Outcome | 4 |

| | | |
|---|--|----|
| 1.7 | Conclusion | 5 |
| CHAPTER 2: LITERATURE REVIEW AND PROJECT METHODOLOGY . 6 | | |
| 2.1 | Introduction..... | 6 |
| 2.2 | Facts and Findings | 6 |
| 2.2.1 | Domain | 6 |
| 2.2.2 | Existing System | 7 |
| 2.2.3 | Technique | 11 |
| 2.3 | Project Methodology..... | 13 |
| 2.3.1 | Requirements Phase..... | 13 |
| 2.3.2 | Planning Phase..... | 14 |
| 2.3.3 | Design Phase..... | 14 |
| 2.3.4 | Implementation Phase..... | 14 |
| 2.3.5 | Testing and Integration Phase..... | 15 |
| 2.4 | Project Requirements | 15 |
| 2.4.1 | Software Requirements..... | 15 |
| 2.4.2 | Hardware Requirements..... | 17 |
| 2.5 | Project Schedule and Milestones | 18 |
| 2.6 | Conclusion | 20 |
| CHAPTER 3: ANALYSIS..... 21 | | |
| 3.1 | Introduction..... | 21 |
| 3.2 | Problem Analysis | 21 |
| 3.3 | Requirement Analysis | 22 |
| 3.3.1 | Bahasa Melayu Translator using Existing Model..... | 22 |
| 3.3.2 | Malaysian Dialect Classification Technique | 23 |

| | | |
|---------------------------------------|--|-----------|
| 3.3.3 | Dialect Translator using Rule-Based System | 25 |
| 3.3.4 | Data Requirement | 26 |
| 3.3.5 | Functional Requirement..... | 29 |
| 3.3.6 | Non-functional Requirement | 30 |
| 3.3.7 | Others Requirement | 30 |
| 3.4 | Conclusion | 32 |
| CHAPTER 4: DESIGN | | 33 |
| 4.1 | Introduction..... | 33 |
| 4.2 | High-Level Design..... | 33 |
| 4.2.1 | System Architecture for Expert System/DSS/Simulation | 33 |
| 4.2.2 | Bahasa Melayu Translator Module..... | 35 |
| 4.2.3 | Dialect Detection Module..... | 37 |
| 4.2.4 | Dialect Translator Module..... | 39 |
| 4.2.5 | User Interface Design for Expert System/DSS/Simulation | 41 |
| 4.2.5.1 | Navigation Design..... | 42 |
| 4.2.5.2 | Input Design..... | 43 |
| 4.2.5.3 | Technical Design | 44 |
| 4.2.5.4 | Output Design..... | 46 |
| 4.3 | Conclusion | 47 |
| CHAPTER 5: IMPLEMENTATION..... | | 48 |
| 5.1 | Introduction..... | 48 |
| 5.2 | Software Development Environment Setup..... | 48 |
| 5.3 | Software Configuration Management..... | 49 |
| 5.3.1 | XAMPP..... | 49 |

| | | |
|---------------------------------|--|-----------|
| 5.3.2 | Django..... | 51 |
| 5.3.3 | Cloud Setup using Linux Fedora Server..... | 55 |
| 5.4 | Implementation Status | 56 |
| 5.5 | Python Coding | 57 |
| 5.5.1 | Bahasa Melayu Translator | 57 |
| 5.5.2 | Dialect Translator | 58 |
| 5.5.3 | Dialect Detector | 62 |
| 5.6 | Conclusion | 67 |
| CHAPTER 6: TESTING | | 68 |
| 6.1 | Introduction..... | 68 |
| 6.2 | Test Plan..... | 68 |
| 6.2.1 | Test Organization..... | 69 |
| 6.2.2 | Test Environment..... | 69 |
| 6.2.3 | Test Schedule..... | 70 |
| 6.3 | Test Strategy | 70 |
| 6.3.1 | Classes of Tests..... | 70 |
| 6.4 | Test Implementation | 71 |
| 6.4.1 | Test Description..... | 71 |
| 6.4.2 | Test Data..... | 73 |
| 6.5 | Test Results and Analysis | 73 |
| 6.5.1 | Bahasa Melayu Translator Test and Analysis | 74 |
| 6.5.2 | Dialect Translator Test and Analysis..... | 74 |
| 6.5.2.1 | Dialect Translator Accuracy Result..... | 76 |
| 6.5.3 | Dialect Detector Test and Analysis | 76 |
| 6.5.3.1 | Dialect Detector Accuracy Test..... | 77 |

| | | |
|--|--|-----------|
| 6.5.4 | User Interface Test and Analysis..... | 78 |
| 6.6 | Conclusion | 79 |
| CHAPTER 7: PROJECT CONCLUSION | | 80 |
| 7.1 | Observation on Weakness and Strengths..... | 80 |
| 7.1.1 | Strengths | 80 |
| 7.1.2 | Weaknesses..... | 81 |
| 7.2 | Propositions for Improvement | 82 |
| 7.3 | Project Contribution..... | 83 |
| 7.4 | Conclusion | 84 |
| REFERENCES..... | | 85 |
| APPENDICES | | 87 |



اونيورسيتي تيكنيكل مليسيا ملاك

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

LIST OF TABLES

| | PAGE |
|---|-----------|
| Table 1.1: Modules in the Malayfy Application | 3 |
| Table 2.1: Software Requirements | 15 |
| Table 2.2: Hardware Requirements | 17 |
| Table 2.3: Project Schedule for FYP 1 | 18 |
| Table 2.4: Project Schedule for FYP 2 | 19 |
| Table 3.1: dialect_list.csv | 27 |
| Table 3.2: dict_perak.csv | 28 |
| Table 3.3: Others Requirement | 30 |
| Table 3.4: Python Libraries | 31 |
| Table 5.1: Step of installing XAMPP for this project | 50 |
| Table 5.2: Implementation Status..... | 56 |
| Table 5.3: Bahasa Melayu Translator coding explanation | 58 |
| Table 5.4: Coding explanation about the dataset | 59 |
| Table 5.5: Dialect translation process | 61 |
| Table 5.6: Prediction process | 65 |
| Table 6.1: Test Organization..... | 69 |
| Table 6.2: Test Schedule | 70 |
| Table 6.3: Classes of Tests | 70 |
| Table 6.4: Test Description..... | 71 |
| Table 6.5: Bahasa Melayu translator test result | 74 |
| Table 6.6: Dialect translator test result..... | 75 |
| Table 6.7: Dialect detector test result..... | 76 |
| Table 6.8: User interface test result..... | 78 |
| Table 7.1: The Details of the Improvements..... | 82 |

LIST OF FIGURES

| | PAGE |
|--|------|
| Figure 2.1: Agile Methodology | 13 |
| Figure 3.1: Sample data of Kelantan dialect phrases | 24 |
| Figure 3.2: Accuracy, precision, and recall result of dialect detector | 25 |
| Figure 3.3: Dialect dictionary dataset | 29 |
| Figure 4.1: System Architecture | 34 |
| Figure 4.2: Bahasa Melayu Translator Flowchart..... | 35 |
| Figure 4.3: Bahasa Melayu Translator Interface..... | 36 |
| Figure 4.4: Dialect Detection Flowchart | 37 |
| Figure 4.5: One of the form's respondents..... | 38 |
| Figure 4.6: Dialect detector interface | 39 |
| Figure 4.7: Dialect Translator Flowchart | 39 |
| Figure 4.8: Dialect translator interface..... | 40 |
| Figure 4.9: Website Flowchart..... | 41 |
| Figure 4.10: Website Navigation Design | 42 |
| Figure 4.11: User input for dialect detection | 44 |
| Figure 4.12: User input for translation | 44 |
| Figure 4.13: Import translation model..... | 45 |
| Figure 4.14: Import Random Forest algorithm | 45 |
| Figure 4.15: Dialect prediction output | 47 |
| Figure 4.16: Translated text output..... | 47 |
| Figure 5.1: Overall software development setup..... | 49 |
| Figure 5.2: Setup localhost using XAMPP..... | 49 |
| Figure 5.3: Create a Django project (Malayfy) | 51 |
| Figure 5.4: Run Django using terminal | 51 |

| | |
|---|-----------|
| Figure 5.5: Django successfully worked page in localhost | 52 |
| Figure 5.6: Create a new Django application | 52 |
| Figure 5.7: Second app's necessary components..... | 53 |
| Figure 5.8: Setting up the second_app URLs | 53 |
| Figure 5.9: Import Libraries at views.py | 54 |
| Figure 5.10: Dialect translator integration | 54 |
| Figure 5.11: Cloud Setup..... | 55 |
| Figure 5.12: Import the translation model. | 58 |
| Figure 5.13: User input..... | 58 |
| Figure 5.14: Translation output..... | 58 |
| Figure 5.15: import the Kelantan dictionary dataset | 59 |
| Figure 5.16: dict_kelantan.csv output | 59 |
| Figure 5.17: convert the csv file to dictionary | 60 |
| Figure 5.18: Rule-based method | 60 |
| Figure 5.19: user input and translation output | 60 |
| Figure 5.20: Import csv..... | 62 |
| Figure 5.21: User input..... | 62 |
| Figure 5.22: Append new sentence to the dataset | 63 |
| Figure 5.23: convert dialect class to numbers | 63 |
| Figure 5.24: Bag of words method..... | 64 |
| Figure 5.25: Bag of words output | 64 |
| Figure 5.26: Copying values of the DataFrame..... | 64 |
| Figure 5.27: Implementing Random Forest algorithm..... | 65 |
| Figure 5.28: Prediction result..... | 65 |
| Figure 5.29: Accuracy result | 66 |
| Figure 6.1: User accuracy feedback for dialect translator..... | 76 |
| Figure 6.2: User accuracy feedback of dialect detector..... | 77 |



LIST OF ABBREVIATIONS

| | | |
|-------------|---|---|
| FYP | - | Final Year Project |
| ISDM | - | Intelligent System Development Methodology |
| NLP | - | Natural Language Processing |
| NMT | - | Neural Machine Translation |
| SMT | - | Statistical Machine Translation |
| ANN | - | Artificial Neural Network |
| RNN | - | Recurrent Neural Network |
| LSTM | - | Long Short-Term Memory |
| GRU | - | Gated Recurrent Unit |

اونيورسيتي تيكنيكل مليسيا ملاك

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

LIST OF ATTACHMENTS

| | | PAGE |
|-------------------|----------------------------------|-------------|
| Appendix A | Google Form User Feedback | 85 |
| Appendix B | Python Coding | 91 |
| Appendix C | Website Coding | 99 |



CHAPTER 1: INTRODUCTION

1.1 Introduction

Effective communication across language borders is now more crucial than ever thanks to globalization and the world's growing interconnectedness. Among people, communities, and nations, language plays a crucial role in communication, connection, and understanding. Language diversity can, however, provide difficulties, obstructing clear communication and resulting in misunderstandings and misinterpretations. The communication barrier that develops when persons or organizations need to interact with or convey information to people who speak or comprehend only Bahasa Melayu is the issue statement for the requirement of a Bahasa Melayu translator.

To bridge the language barrier and facilitate effective communication between parties, whether in business, education, healthcare, or other professions, a Bahasa Melayu translator is required. Without a Bahasa Melayu translator, crucial information could be misunderstood, misconstrued, or misread, with potentially disastrous effects. By offering a reliable and effective method for dialect detection and translation in the Bahasa Melayu language, the Bahasa Melayu Dialect Translator and Detector system (Malayfy) intends to overcome these issues. Malay, often known as Bahasa Melayu, is a widely spoken language in Southeast Asia with numerous regional varieties.

1.2 Problem Statement(s)

- i. Language barriers can impede communication and lead to misunderstandings and misinterpretations.
- ii. Current machine translation tools for Malay are often inaccurate and unreliable.
- iii. There are challenges for businesses, organizations and individuals who need to communicate with Malay speakers in various contexts, such as tourism, education, health, and commerce.

1.3 Objective(s)

- i. To gather Bahasa Melayu words and their dialect.
- ii. To develop language translator from English to Bahasa Melayu and Bahasa Melayu dialect translator and detector.
- iii. To evaluate Bahasa Melayu translator and Bahasa Melayu dialect translator and detector.

1.4 Project Scope

The project scope defines the boundaries and objectives of the Bahasa Melayu Dialect Translator and Detector system. It describes the particular system modules that will be created in detail and specifies the target users who will profit from its features. The core modules and the target user base are highlighted in this chapter's overview of the project's scope.

1.4.1 Module(s)

The table below shows the three components that will be built for the Malayfy application. Natural Language Processing (NLP) techniques will be used in the modules that translate Bahasa Melayu using the Transformers model of deep learning,

detect dialects using random forest algorithms, and translate dialects using rule-based systems.

Table 1.1: Modules in the Malayfy Application

| Modules | Functions |
|---------------------------------|---|
| Bahasa Melayu Translator | Able to aid user in interacting and communicating with the native in Malaysia by help user to translate the words or sentences in English to Bahasa Melayu. |
| Dialect Detector | For the dialect detection feature, classification machine learning classifies the sentences of Malaysian dialects such as Johor, Kedah, Kelantan, Melaka, Negeri Sembilan, Pahang, and Penang. This able to assist the user in identifying sentences in a dialect and predicting which dialect the phrases belong to. |
| Dialect Translator | The general bahasa melayu will be translated to the selected dialect using this dialect translator module. With a rule-based approach, if the vocabulary is present in the selected dialect dictionary, the words from the sentences will transform to that dialect. |
| User interface module | This module will provide a user-friendly interface for users to interact with the system, such as a website. It will include input/output text fields and language selection options. |

1.4.2 Target User(s)

- i. Businesses that operate in or trade with Malaysia and need to communicate with Bahasa Melayu-speaking customers, clients, or partners.

- ii. Educational institutions that have Bahasa Melayu-speaking students, teachers, or researchers and need to communicate with them effectively.
- iii. Tourists and travelers who visit Malaysia and want to communicate with locals in Bahasa Melayu.

1.5 Project Significant

The project's successful completion will significantly impact cross-cultural communication, promoting understanding and enabling global collaborations. The importance of a Bahasa Melayu translation project rests in its capacity to break down language barriers and promote interlingual communication. This is crucial in a nation like Malaysia, where a diverse populace speaks various languages in addition to Bahasa Melayu, the official language. A Bahasa Melayu translator can support increased communication and cooperation between various populations inside and outside Malaysia. A competent Bahasa Melayu translation can also enable others who may not be fluent in the language to access information, education, and opportunities, encouraging inclusivity and equality. Moreover, the project will also contribute to developing natural language processing and machine learning technologies, which will have broad applications beyond Bahasa Melayu translation.

1.6 Expected Outcome

The goal of the Malayfy project is to create a reliable and effective dialect translator, dialect language detector, and English to Malay translation system. The initiative aims to improve comprehension and communication of Bahasa Melayu by successfully identifying and classifying various dialects and by providing accurate and trustworthy translations. The system is anticipated to have an intuitive user interface that makes it simple to access its functions and promotes seamless communication across language boundaries. If the project is successful, people, companies, and organizations will be able to get over language obstacles and improve their Bahasa Melayu communication skills.

1.7 Conclusion

In conclusion, this project will provide a solution for native speakers when they come to Malaysia. This website will help them communicate with locals and understand the dialects used by them, as there are many different dialects. A discussion of the project's methodology and a review of relevant studies will follow in the subsequent chapter.

