

**MIMIC HUMAN SPEECH IN BAHASA INDONESIA USING  
SPEECH RECOGNITION AND SPEECH SYNTHESIS**

By

Valens Prabagita Ivan Susilo

A Thesis  
Submitted to the Faculty of Computing  
President University  
in Partial Fulfilment of the Requirements  
for the Degree of Bachelor of Science  
in Information Technology

Cikarang, Bekasi, Indonesia

October 2018

Copyright by  
Valens Prabagita Ivan Susilo  
2018

**MIMIC HUMAN SPEECH IN BAHASA INDONESIA USING  
SPEECH RECOGNITION AND SPEECH SYNTHESIS**

By

Valens Prabagita Ivan Susilo

Approved:

---

Dr. Tjong Wan Sen, S.T., M.T.  
Thesis Advisor

---

Drs. Nur Hadisukmana, M.Sc.  
Program Head of Information Technology

---

Ir. Rila Mandala, M.Eng., Ph.D.  
Dean of Faculty of Computing

## **ABSTRACT**

Everyday people use Speech recognition and speech synthesis unconsciously. The technologies help them with their activities. With each technology can produce any kinds software related to speech. Combine both of technologies can produce many more. One of the combinations is mimic human speech. This research will discuss about Convolutional Neural Network, Speech Recognition, and Speech Synthesis. The purpose of this research is to develop application to collect, train, and mimic speech in Bahasa Indonesia. User can participate record their speech. The collected speech will be train to be used in the application to recognize the speech. After the collected speech is trained, User can mimic their speech by identify or recognize the speech and generate or synthesis the speech. The application to collect and mimic speech develops in website application.

**Keywords:** Convolutional Neural Network, Speech Recognition, Speech Synthesis

## **ACKNOWLEDGMENTS**

The author gratefully acknowledges the amazing supports, guidance and advices from:

1. Bebe, Annisa Feryannie Widya.
2. My beloved family, Ayah, Ibu, Mas Risky and Dek Lia.
3. B35TFR13NDZONE, Nisa, Marsel, Bora, Fira, Rai, Pindy, Oscar, and Michel.
4. Sangar Team, Mas Andreas, Mba Diana, Nisa, Alfian, and Aci
5. DV 14, Vovo, Gojal, Ariyel, Rino, Willy, Azmi, Damara, Atisah, Risda, Sonya, Dini, and Arizha.
6. My cool thesis lecturer, Mr. Wan Sen.
7. My second family, Van Lith XXII.
8. And you, yes you.

## TABLE OF CONTENTS

	Page
ACKNOWLEDGMENTS .....	i
TABLE OF CONTENTS .....	ii
LIST OF TABLES .....	iv
LIST OF FIGURES .....	v
CHAPTER	
I INTRODUCTION .....	1
1.1 Background .....	1
1.2 Problem Statement .....	2
1.3 Research Objective .....	2
1.4 Scope and Limitation .....	2
1.5 Thesis Methodology .....	3
1.6 Thesis Outline .....	4
II LITERATURE STUDY .....	6
2.1 Machine Learning .....	6
2.1.1 Supervised Learning .....	7
2.1.2 Unsupervised Learning .....	7
2.2 Neural Network .....	9
2.3 Phonemic .....	9
2.4 Speech Recognition .....	10
2.5 Pre-processing Speech .....	10
2.5.1 Sound Sampling .....	10
2.5.2 MFCC .....	11
2.6 Speech Synthesis .....	12
2.7 Related Work .....	12
2.7.1 Lyrebird .....	12
2.7.2 Google Translate .....	13
III SYSTEM ANALYSIS .....	15
3.1 System Overview .....	15
3.2 Functional Analysis .....	15
3.3 Software and System Requirements .....	16

CHAPTER	Page
3.4 System Architecture .....	17
3.4.1 Use-Case Diagram .....	17
3.4.2 Use-Case Narrative .....	18
3.4.3 Activity Diagram .....	21
IV SYSTEM DESIGN .....	23
4.1 User Interface Design .....	23
4.1.1 Home Screen .....	23
4.1.2 Identify Speech Screen .....	24
4.1.3 Generate Speech Screen .....	25
4.2 Class Diagram .....	26
4.2.1 Front-end Library .....	27
4.2.2 Back-end Library .....	29
IV SYSTEM DEVELOPMENT .....	31
5.1 User Interface Development .....	31
5.1.1 Collect Application .....	31
5.1.1.1 Home Page .....	31
5.1.1.2 Phonemes Page .....	32
5.1.2 Mimic Application .....	34
5.1.2.1 Home Page .....	34
5.1.2.2 Identify Page .....	35
5.1.2.3 Generate Page .....	38
5.2 Application Details .....	41
5.2.1 Server Code .....	41
5.2.2 Router Code .....	43
5.2.3 Database Code .....	51
5.2.4 Train Application Code .....	53
5.2.5 Library Package Code .....	58
IV SYSTEM TESTING .....	59
6.1 Testing Environment .....	59
6.2 Testing Scenario .....	59
6.2.1 Collect Application .....	60
6.2.2 Train Application .....	61
6.2.3 Mimic Application .....	61
IV CONCLUSIONS AND FUTURE WORK .....	59
7.1 Conclusion .....	59
7.2 Future Work .....	59
REFERENCES .....	vi

## LIST OF TABLES

TABLE	Pag
3.1 Functionality Table .....	20
3.2 Use-Case Narrative – Make Speech ID .....	22
3.3 Use-Case Narrative – Identify Speech .....	23
3.4 Use-Case Narrative – Select Speech ID .....	24
3.5 Use-Case Narrative – Generate Speech .....	24
4.1 Home Screen Description .....	28
4.2 Identify Speech Screen Description .....	29
4.3 Generate Speech Screen Description .....	30
6.1 Collect application scenarios .....	28
6.2 Train application scenarios .....	29
6.3 Mimic application scenarios .....	30
6.4 Tester 1 speech recognition on identify process results.....	24
6.5 Tester 2 speech recognition on identify process results.....	24
6.4 Tester 3 speech recognition on identify process results .....	24



## LIST OF FIGURES

FIGURE	Page
2.1 Neural network to find estimated price from specific input .....	8
2.2 Phonemes, graphemes, and letters in the word “spoon” .....	9
2.3 Sound sampling process .....	11
2.4 100 samples in numbers from “Hello” sound with sample rate of 16kHz (16 samples per seconds) .....	11
2.5 Screenshot of Lyrebird in the website .....	13
2.6 Screenshot of Google Translate in website .....	14
3.1 Use-Case Diagram .....	18
3.2 Activity Diagram .....	22
4.1 Home Screen .....	24
4.2 Identify Speech Screen .....	25
4.3 Generate Speech Screen .....	26
4.4 Front-end Class Diagram .....	28
4.5 Back-end Class Diagram .....	30
5.1 Collect Home Page .....	22
5.2 Collect Phonemes page with ‘a’ phoneme .....	24
5.3 Collect Phonemes page after record process is finish .....	25
5.4 Mimic Home page .....	26
5.5 Mimic Identify page .....	28
5.6 Mimic Identify page alert user indicating identify process in the server is finish .....	30

FIGURE	Page
5.7 Mimic Identify page alert user indicating there is no value in Input element .....	18
5.8 Mimic Generate page .....	22
5.9 Mimic Generate page play the generate speech after generate process success .....	24
5.10 Mimic Generate page alert user indicating there is no value in Textarea element .....	25
5.11 Mimic Generate page alert user indicating there is error in the generate process .....	26
5.12 serverHandler on Server code .....	28
5.13 Create server command on Server code .....	13
5.14 uploadCollect on Collect Router code .....	14
5.15 router on Collect Router code .....	18
5.16 extract on Identify Router code .....	22
5.17 uploadSpeech on Identify Router code .....	24
5.18 identifySpeech on Identify Router code .....	25
5.19 router on Identify Router code .....	26
5.20 loadSpeech on Generate Router code .....	28
5.21 generateSpeech on Generate Router code .....	30
5.22 router on Generate Router code .....	13
5.23 connection on Database code .....	14
5.24 model on Database code .....	18
5.25 extractWav on Train application code .....	22

FIGURE	Page
5.26 extractFiles on loadData Train application code .....	24
5.27 load and split command on loadData Train application code .....	25
5.28 getData on nextBatch Train application code .....	26
5.29 getLabel on nextBatch Train application code .....	28
5.30 saveModelDB on Train application code .....	30
5.31 modelML on model Train application code .....	24
5.32 preparation command on train Train application code .....	25
5.33 train on train Train application code .....	26
6.1 Some screenshot from the Collect application scenarios .....	26
6.2 Some screenshot from the Train application scenarios .....	26
6.3 Some screenshot from the Mimic application scenarios .....	26
6.4 Some screenshot from the Tester 1 results .....	26
6.5 Some screenshot from the Tester 2 results .....	26
6.6 Some screenshot from the Tester 3 results .....	26

## **CHAPTER I**

### **INTRODUCTION**

#### **1.1 Background**

“Ok Google, play some music”. “Siri, what should I eat for lunch?”. Everyday people use their virtual assistance to boost their activities. People very like to use it because they just asked to their device and then in seconds, the wish is granted. It seems like, people are talking to the computer. The truth is, speech recognition takes big role with the help of machine learning. Google Assistance, Apple Siri, Microsoft Cortana, Amazon Alexa, and others have thousands of speech data to be analysed with the machine learning and they easily add data by collecting people speech from the assistance with permission.

If speech recognition is the process to get data by analysed speech, the opposite of speech recognition is speech synthesis, the process to produce artificial speech. Therefore, speech recognition is known as speech-to-text and speech synthesis is known as text-to-speech. “Hey Cortana, read my email” command make virtual assistance generate speech from the email text. With each technology can produce any kinds software related to speech. Combine both of technologies can produce many more. One of the combinations is mimic human speech.

## **1.2 Problem Statement**

This research aims to develop application which can be used to collect speech data, train the collected data and mimic speech in Bahasa Indonesia. The application to collect and mimic speech develops in website application. The application can recognize the speech and generate speech from text.

## **1.3 Research Objective**

This research sees an opportunity to implement speech recognition and speech synthesis to create a mimic speech.

## **1.4 Scope and Limitation**

This research focuses on developing an application which will be able to:

1. Perform collecting data.
2. Perform train the collected data.
3. Perform speech recognition.
4. Perform speech synthesis.

The limitations of this application are as following:

1. There are 10 selected phonemes to be used in the application, a, i, t, na, ma, mu, di, ri, and ku.
2. Recorded speech in 1 second, with sample rate 16000 and mono sound.
3. Speech recognition data is taken from recorded speech and in human speech in Bahasa Indonesia.
4. Speech synthesis data is taken from saved speech, result from speech recognition.
5. Application is developed as website application.

6. Collecting, testing, and running the application is used in the same hardware.

## **1.5 Thesis Methodology**

Rapid Application Development (RAD) methodology will be used in the development of this application. The RAD method, which was first developed by James Martin, is a Software Development Life Cycle method that gains its popularity in recent years due to its suitability to manage web application projects. The features of RAD were designed to overcome most of the shortcomings found in traditional waterfall model. Some of these features are: fast prototyping and capability to deal with change in requirements.

The RAD model implemented in this thesis will consists of four major phases:

### **1. Requirement Planning Phase**

This is the where system planning and analyses are done. System requirements are established, including the view range of the camera, the numbering schema for parking lots, and the general category of cars to be detected. Algorithms are proposed to solve the problem along with the overall outline of the program.

### **2. User Design Phase**

During this phase, the model of system's processes is the main focus. Models for input, output, process and user interface is built, and represented in different parts that include diagrams visualization. The system design will refer

to the plan created in previous stage. The design will be continuously discussed, reviewed, and updated until the best version is found.

### 3. Development Phase

This phase is where the ideas and plans are executed. The application is developed according to the predefined features standard. All the components of image processing and object detection are put together into one program to perform the work from beginning to the final output. Unit testing will also be done here. It focuses on application development, including: coding, unit integration, and testing.

### 4. Cut Over Phase

In this final phase there will be some test to evaluate the program's ability to determine which area of parking lot is empty. There will be certain test cases of parking areas that will produce different images to be processed. An evaluation will be done to see how far the application is capable to detect the area. Bugs fixing will also be done in this step. Another part of cutover phase is to create installation and operating manuals to allow people operate the program in their environment.

## **1.6 Thesis Outline**

The thesis consists of seven chapters, which are as follow:

### 1. Chapter I: Introduction

This chapter introduce the research background, problem, and objective. It also explains the research scope and limitation, method to achieve the objective.

## 2. Chapter II: Literature Study

This chapter contain the literature study that related to the research background.

## 3. Chapter III: System Analysis

This chapter explains the analysis of the application – both in its function and behaviour, in order to fulfil the prescribed requirements.

## 4. Chapter IV: System Design

This chapter explains the system design of interfaces and class diagram based on the previous chapter that will be used in the next chapter.

## 5. Chapter V: System Development

This chapter explains the system development of interfaces and code details on the application.

## 6. Chapter VI: System Testing

This chapter ensures the application system runs well by evaluating all the features, and making sure the system fulfils its function requirements.

## 7. Chapter VII: Conclusion and Future Work

This chapter sums up this research and also suggestion for future research work.



## **CHAPTER II**

### **LITERATURE STUDY**

#### **2.1 Machine Learning**

Machine learning is a form of AI that enables a system to learn from data rather than through explicit programming [1]. Others state that machine learning usually refers to the changes in systems that perform tasks associated with artificial intelligence (AI) [2].

Basically, machine learning consists of 2 words, Machine and Learning. As a noun, there are many definitions related to machine [3] and one of the definitions, machine is a computer. Meanwhile, learning can be interpreted as the activity of obtaining knowledge or knowledge obtained by study [4]. Combine both words, machine learning can be interpreted as a computer that does an activity of obtaining knowledge.

In general, there are 2 types of machine learning Supervised Learning and Unsupervised Learning.

##### *2.1.1 Supervised Learning*

In supervised machine learning, a system is trained with data that has been labelled. The labels categorise each data point into one or more groups, such as ‘apples’ or ‘oranges’. The system learns how this data – known as training data – is structured, and uses this to predict the categories of new – or ‘test’ – data [5]. It can be considered

the learning is guided by a teacher. One has a dataset which acts as a teacher and its role is to train the model or the machine. Once the model gets trained it can start making prediction or decision when new data is given to it [6].

It can be concluded that supervised machine learning is the learning way where the machine is guided by labelled data and it needs to learn how to classify the model data to fits the labelled data.

### *2.1.2 Unsupervised Learning*

Unsupervised learning is learning without labels. It aims to detect the characteristics that make data points more or less similar to each other, for example by creating clusters and assigning data to these clusters [5]. Suppose images of apples, bananas and mangoes are presented to the model, so what it does, based on some patterns and relationships it creates clusters and divides the dataset into those clusters. Now if a new data is fed to the model, it adds it to one of the created clusters [6].

It can be concluded that unlike supervised learning, unsupervised learning is the learning way where the machine is the figure out data by cluster the data.

## **2.2 Neural Network**

A neural network is an approach to machine learning in which small computational units are connected in a way that is inspired by connections in the brain [5].

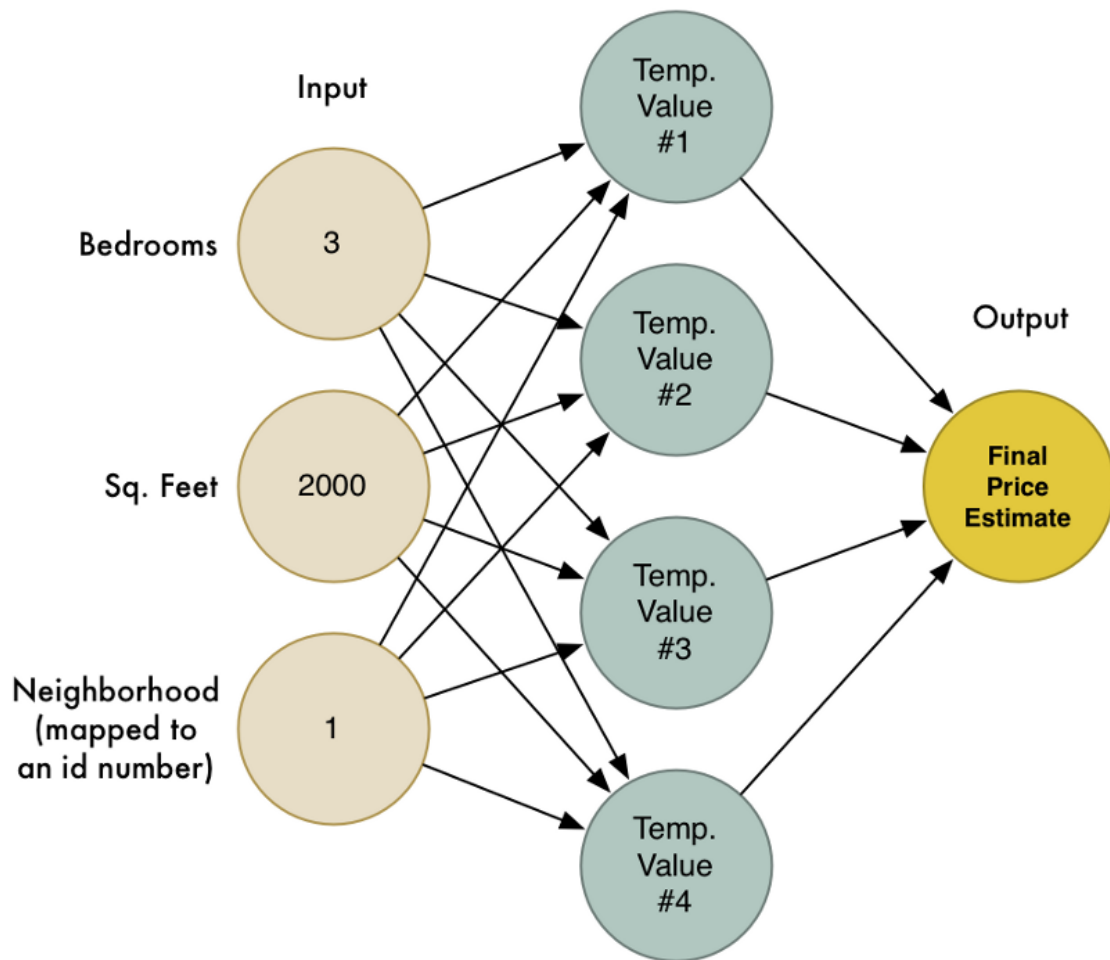


Figure 2.1 Neural network to find estimated price from specific input [7].

Every neural network model is basically a three-layered system, which are Input layer, Hidden Layer and Output Layer [8]. Input layer, where the inputs of the problem are received, hidden layers, where the relationship between the inputs & outputs are determined & represented by synaptic weights, & an output layer which emits the outputs of the problem [9].

*...Example and explanation of Neural Network...*

## 2.3 Phonemic

*...Explanation of phonemic...*

A phoneme is the basic unit of phonology. It is the smallest unit of sound that may cause a change of meaning within a language, but that doesn't have meaning by itself. For example, in the words "bake" and "brake," only one phoneme has been altered, but a change in meaning has been triggered. The phoneme /r/ has no meaning on its own, but by appearing in the word it has completely changed the word's meaning [10].

In the other hand, grapheme is individual letters and groups of letters that represent single phonemes, like the "s" and the "oo" in "spoon". Understanding how letters are used to encode speech sounds in written language is crucial in learning to decode unfamiliar words [11].

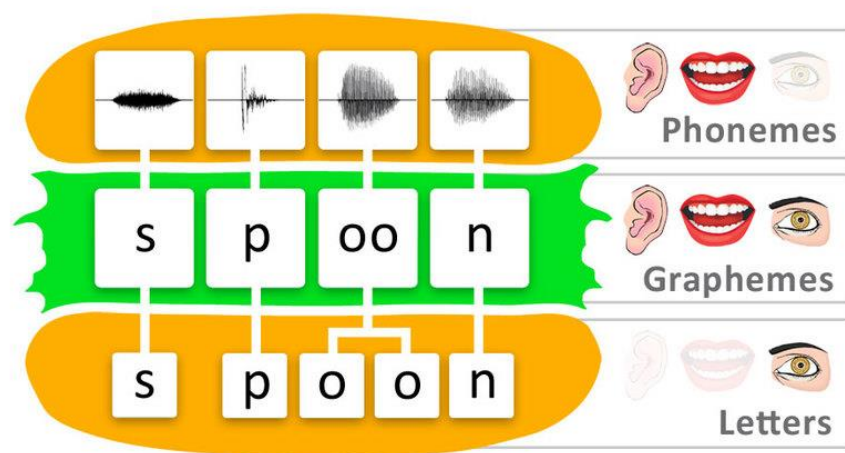


Figure 2.2 Phonemes, graphemes, and letters in the word "spoon" [11].

In Bahasa Indonesia phoneme is distributed into *vokal*, *diftong*, *konsonan*, *gugus konsonan* [12].

*...Explanation of silabel...*

## **2.4 Speech Recognition**

The process of automatically recognizing spoken words of speaker based on information in speech signal is called Speech Recognition [21]. Other definition of speech recognition, also known as Automatic Speech Recognition (ASR), or computer speech recognition, is the process of converting a speech signal to a sequence of words, by means of an algorithm implemented as a computer program [22].

*...Google Convolutional Neural Network for Small-footprint Keyword Spotting...*

## **2.5 Pre-processing Speech**

As sound is transmitted as waves, and computer understand numbers, it's necessary to pre-processing speech so that computer understand and can be feed as input into neural network.

*...Additional explanation to convert to CNN...*

### **2.5.1 Sound Sampling**

Sound sampling is taking a reading thousands of times a second and recording a number representing the height of the sound wave at that point in time. Basically, all an uncompressed .wav audio file [23].

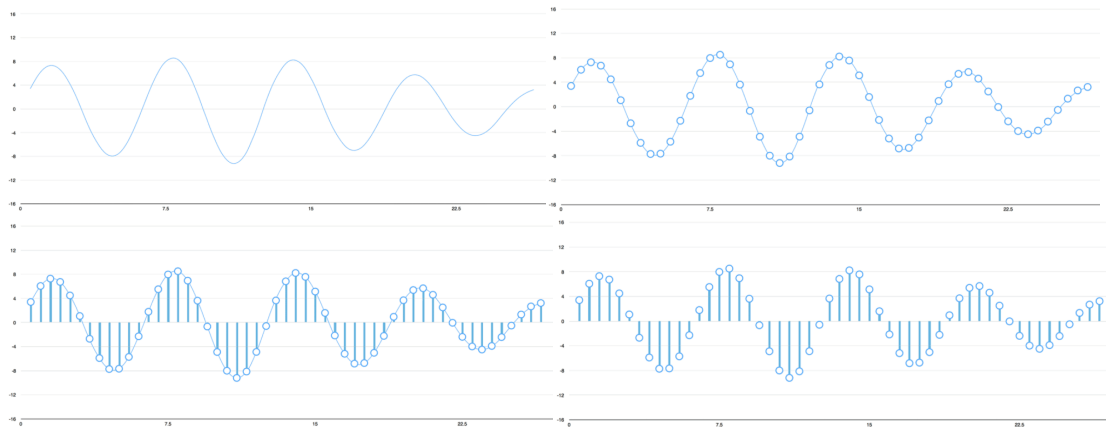


Figure 2.3 Sound sampling process [23].

```
[-1274, -1252, -1160, -986, -792, -692, -614, -429, -286, -134, -57, -41, -169, -456, -450, -541, -761, -1067, -1231, -1047, -952, -645, -489, -448,
-397, -212, 193, 114, -17, -110, 128, 261, 198, 390, 461, 772, 948, 1451, 1974, 2624, 3793, 4968, 5939, 6057, 6581, 7302, 7640, 7223, 6119, 5461,
4820, 4353, 3611, 2740, 2004, 1349, 1178, 1085, 901, 301, -262, -499, -488, -707, -1406, -1997, -2377, -2494, -2605, -2675, -2627, -2500, -2148, -
1648, -970, -364, 13, 260, 494, 788, 1011, 938, 717, 507, 323, 324, 325, 350, 103, -113, 64, 176, 93, -249, -461, -606, -909, -1159, -1307, -1544]
```

Figure 2.4 100 samples in numbers from “Hello” sound with sample rate of 16kHz  
(16 samples per seconds) [23].

Nyquist sampling theorem [24] provides a prescription for the nominal sampling interval required to avoid aliasing. The sampling frequency should be at least twice the highest frequency contained in the signal. In the case, where one has  $f_c = 3$  Hz, and so the Nyquist theorem tells that the sampling frequency,  $f_s$ , must be at least 6 Hz [25].

## 2.5.2 MFCC

...Explanation of MFCC...

## 2.6 Speech Synthesis

Speech synthesis is the artificial production of human speech [26]. The Synthesized speech can be created by concatenating pieces of recorded speech that are stored in a database [27]. From phoneme, text can be analysed by its phonemes and then with the phonemes concatenating speech can be done.

## 2.7 Related Work

The following are most related work to the research. Have relation to mimic speech, both speech recognition and speech synthesis.

### 2.7.1 *Lyrebird*

Lyrebird is website application, <https://lyrebird.ai>, that has 3 products: Custom Voice, Vocal Avatar, and Vocal Avatar API [28]. Custom voice is a product to create speech based on real people's speech, it can control the intonation, expression, and the emotion of the speech. Vocal avatar is a product to create own digital speech by read some English sentences, and then generate any sentences with own digital speech. Vocal avatar API is a product to provide API to use user's own vocal avatar.

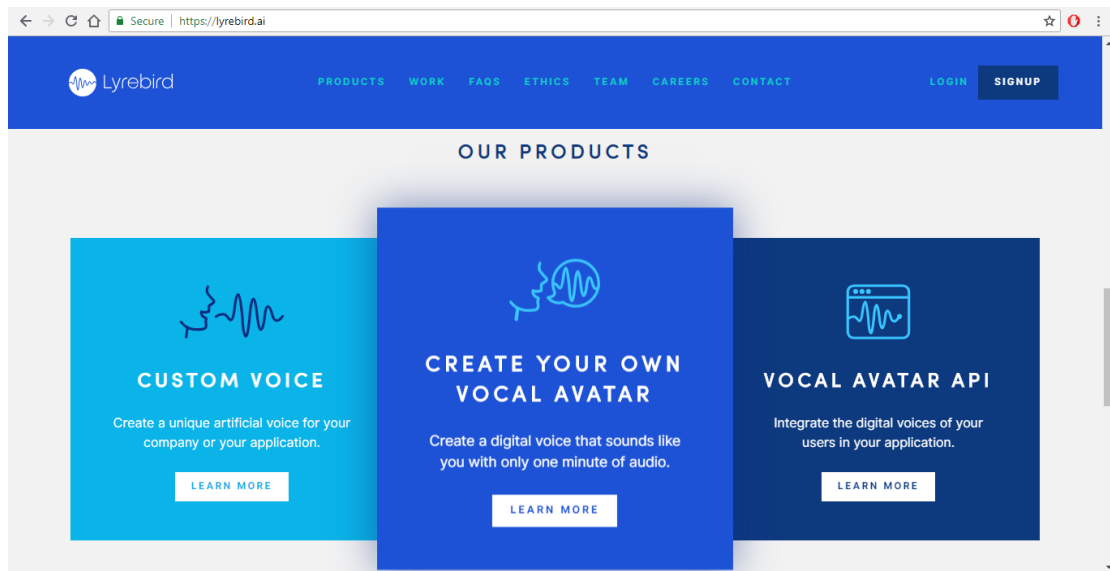


Figure 2.5 Screenshot of Lyrebird in the website [28].

### 2.7.2 Google Translate

Google Translate is one of Google products that is an application to translate languages. Google Translate can be access freely through <https://translate.google.com/>. It has many features [29] and one of them is Talk feature. Talk has function to input text from speech and generate speech from text in any languages.



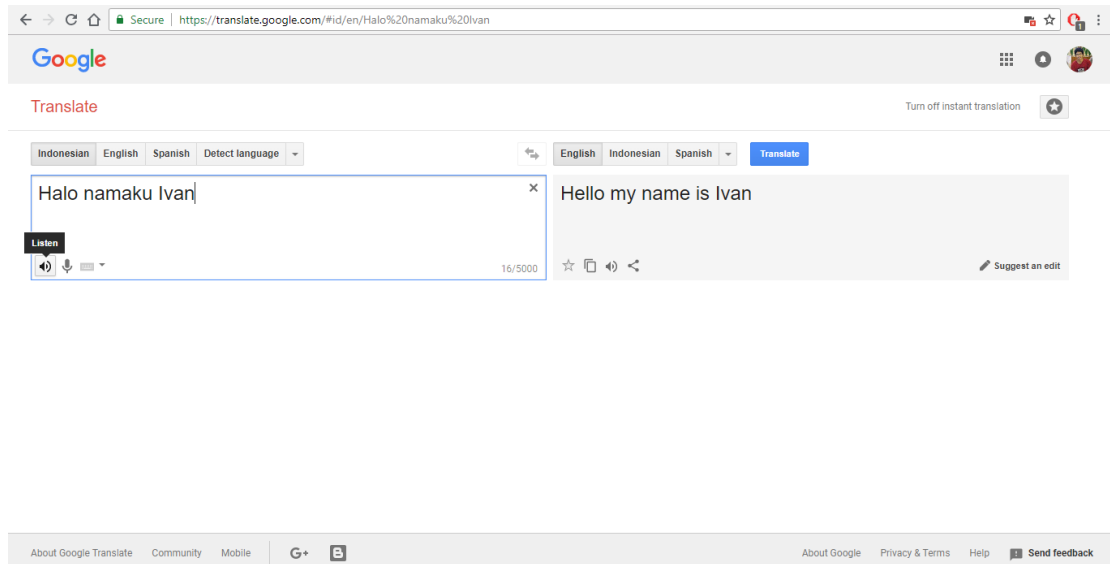


Figure 2.6 Screenshot of Google Translate in website [29].

From related work, it can be concluded that Lyrebird can mimic speech with its vocal avatar, but the speech is in English. In the other hand, Google Translate could mimic speech into any languages, but the speech vocal is from the Google Translate itself.

## CHAPTER III

### SYSTEM ANALYSIS

#### 3.1 System Overview

*...Update according to final app...*

This research is intended to implement speech recognition and speech synthesis into this research. This application will be trained to recognize the speech before used by user. After enough training, this application will identify speech from the user based on sentences that will be displayed. Then, with speech synthesis user can generate speech from identified speech that will become mimic speech. The objective of this research is to create a web-based application for mimic speech by identify user speech and then generated them.

#### 3.2 Functional Analysis

*... Update according to final app...*

There are several functions from this application listed in the Table 3.1.

Table 3.1 Functionality Table.

No	Function Description
1	Allow user to identify user's speech.
2	Allow user to select which speech data that will be used.
3	Allow user to generated speech.

### 3.3 Software and System Requirements

This research and application development should be supported by the following list requirement in order to write the research, build and run the application well.

1. Laptop / Personal Computer

Laptop or Personal Computer is used as the tool where operating system is run. In this research, ASUS A455LN is used with Windows 10 as the OS.

2. Browser

3. Microsoft Office Word

Microsoft Office Word is used to write the research documentation. In this research, Microsoft Office Word 2016 is used.

4. Node.js, JavaScript Run-Time Environment.

Node.js is an open source server environment – Node.js is free – Node.js runs on various platform (Windows, Linux, Unix, Mac OS X, etc) – Node.js uses JavaScript on the server [30]. In this research, Node.js v10.14.1 is used.

5. NoSQL document-oriented database

Document databases pair each key with a complex data structure known as a document. Documents can contain many different key-value pairs, or key-array pairs, or even nested documents [31]. In this research, MongoDB Community Server 4.0.2 as database server and MongoDB Compass 1.15.4 as MongoDB UI.

## 6. Integrated Development Environment (IDE)

IDE is used as application development environment. In this research, Visual Studio Code 1.30.0 is used.

## 7. Git

Git is used as version control. The repository is placed on local and cloud as preventive action. In this research, Git 2.18.0.windows.1 and GitLab as cloud repository is used.

### 3.4 System Architecture

*... Update according to final app...*

This sub-chapter discusses about the use-case diagram and narrative for this application in both point of view, the actors and the system.

#### 3.4.1 Use-Case Diagram

Use-Case Diagram defines the functionality of a system and explains it in user point of view. The actors in this research is the application user. The diagram will be shown in Figure 3.1.

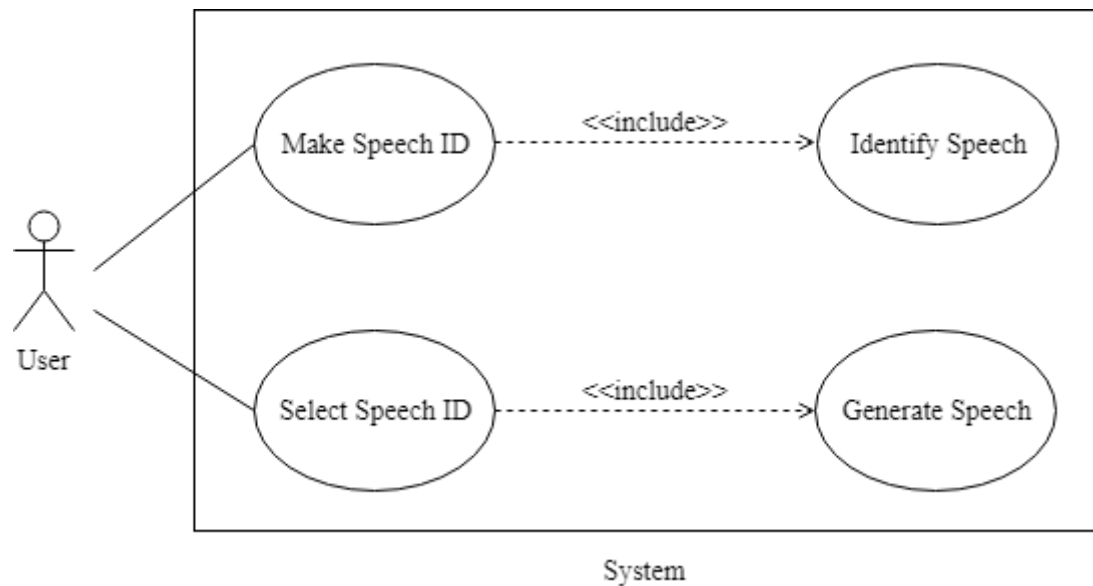


Figure 3.1 Use-Case Diagram.

### 3.4.2 Use-Case Narrative

Use-Case Narrative explains the interaction between the actors and the system. It describes the detail of use-cases such as name, description, pre-condition, post-condition, business rules, and the course of events that happened in the system. The Use-Case Narrative is shown in Table 3.2 and Table 3.5.

Table 3.2 Use-Case Narrative – Make Speech ID.

User Case Name	Make Speech ID
Use Case ID	UC01
Priority	High
Primary Business Actor	User
Primary System Actor	System
Another Participating Actor	None
Description	This use-case describes the event when user opens this application or in the home screen.

Precondition	None	
Trigger	User opens this application or user in the home screen.	
Typical Course of Event	<b>Actor Action</b>	<b>System Response</b>
	Choose Make Speech ID.	Start Make Speech ID activity.
Alternate Course	None	
Post Condition	Identify Speech screen is shown.	
Business Rule	None	
Implementation Constraint and Specifications	None	

Table 3.3 Use-Case Narrative – Identify Speech.

<b>User Case Name</b>	<b>Make Speech ID</b>	
Use Case ID	UC02	
Priority	High	
Primary Business Actor	User	
Primary System Actor	System	
Another Participating Actor	None	
Description	This use-case describes the event when Make Speech ID activity start.	
Precondition	User is from home screen.	
Trigger	User click Make Speech ID button in the home screen.	
Typical Course of Event	<b>Actor Action</b>	<b>System Response</b>
	Do Identify Speech.	Process Speech to Speech Data.
Alternate Course	<b>Actor Action</b>	<b>System Response</b>
	Finish Identify Speech.	Back to home screen.
Post Condition	User do identify speech again or Home screen is shown.	
Business Rule	None	
Implementation Constraint and Specifications	One speech ID for 1 speech data.	

Table 3.4 Use-Case Narrative – Select Speech ID.

<b>User Case Name</b>	<b>Select Speech ID</b>	
Use Case ID	UC03	
Priority	High	
Primary Business Actor	User	
Primary System Actor	System	
Another Participating Actor	None	
Description	This use-case describes the event when user opens this application or in the home screen.	
Precondition	None	
Trigger	User opens this application or user in the home screen.	
Typical Course of Event	<b>Actor Action</b>	<b>System Response</b>
	Select Speech ID.	Provide Speech ID.
Alternate Course	None	
Post Condition	Generate Speech screen is shown.	
Business Rule	None	
Implementation Constraint and Specifications	None	

Table 3.5 Use-Case Narrative – Generate Speech.

<b>User Case Name</b>	<b>Select Speech ID</b>	
Use Case ID	UC04	
Priority	High	
Primary Business Actor	User	
Primary System Actor	System	
Another Participating Actor	None	
Description	This use-case describes the event when Select Speech ID activity start.	
Precondition	User is from home screen.	
Trigger	User click Select Speech ID button in the home screen.	
Typical Course of Event	<b>Actor Action</b>	<b>System Response</b>
	Selected Speech ID.	Start Generate Speech Activity.
	Generate Speech.	Speech generated based on speech ID data.

Alternate Course	<b>Actor Action</b>	<b>System Response</b>
	Finish Generate Speech.	Back to home screen.
Post Condition	User do generate speech again or Home screen is shown.	
Business Rule	None	
Implementation Constraint and Specifications	None	

### 3.4.3 Activity Diagram

Activity Diagram presents a flowchart to represent the flow from one activity to another activity. As user open the application the home screen is shown. User could decide by choose to Make Speech ID or Select Speech ID. As long as user don't decide it stays in the home screen. If Make Speech ID is selected, Make Speech ID activity is started which is Identify Speech. Finish the activity will bring user back to home screen. If Select Speech ID is selected, Select Speech ID activity is started which is Generated Speech. Finish the activity will bring user back to home screen. The figure of the activity diagram can be seen in Figure 3.2.



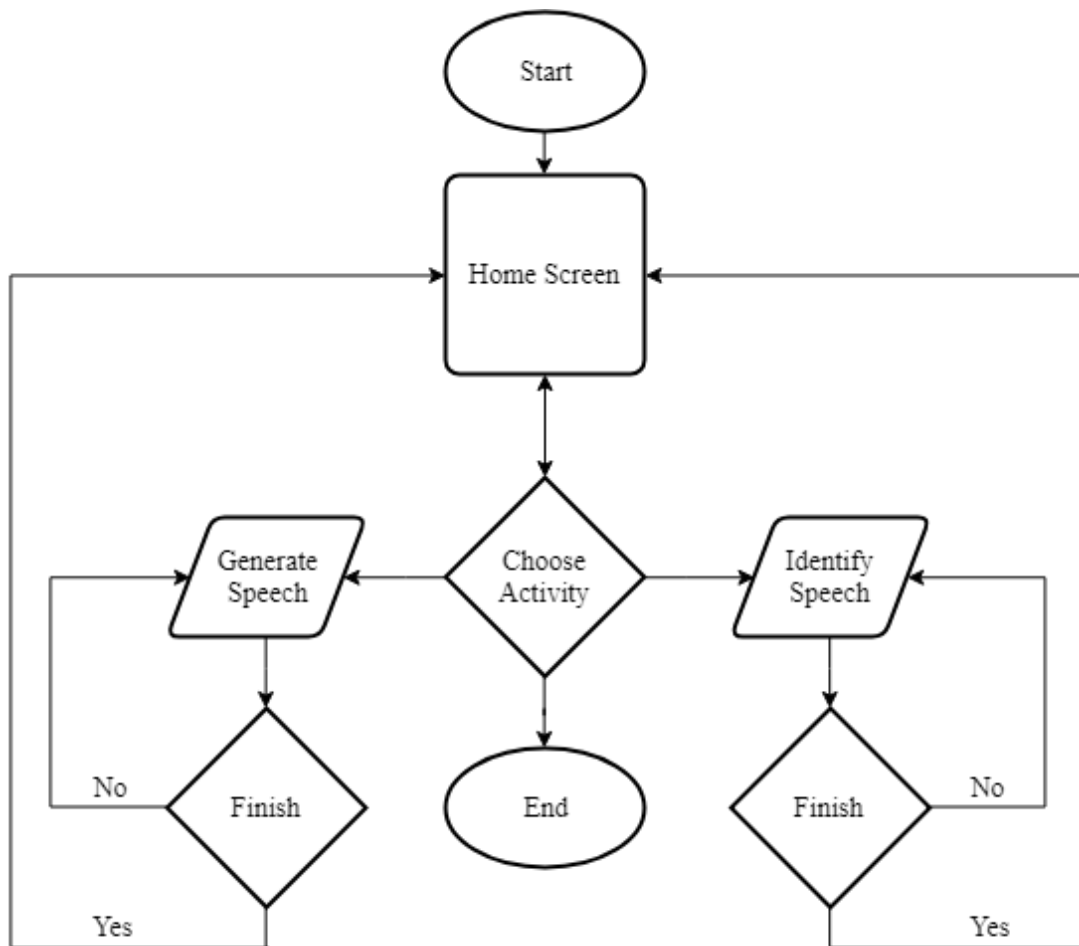


Figure 3.2 Activity Diagram.

## **CHAPTER IV**

### **SYSTEM DESIGN**

#### **4.1 User Interface Design**

*...Update according to final app...*

The User Interface (UI) design of this mobile application is divided into several features which are home screen, identify speech screen, and generate speech screen. The detail of every feature will be explained further below.

##### *4.1.1 Home Screen*

Figure 4.1 shows the design layout for home screen of the application. When user opens the application or finish Make Speech ID or Select Speech ID, it will show the home screen that consist of 2 buttons such as Make Speech ID and Select Speech ID. The description of the design layout is shown in Table 4.1.



Figure 4.1 Home Screen.

Table 4.1 Home Screen Description.

No	Description
1	Make Speech ID
2	Select Speech ID

#### 4.1.2 Identify Speech Screen

Figure 4.2 shows the design layout for identify speech screen of the application. When user choose or click Make Speech ID, it will show the identify speech screen that consist of 4 element such as Random Sentences, Record Button, Identify Button, and Finish Button. The description of the design layout is shown in Table 4.2.

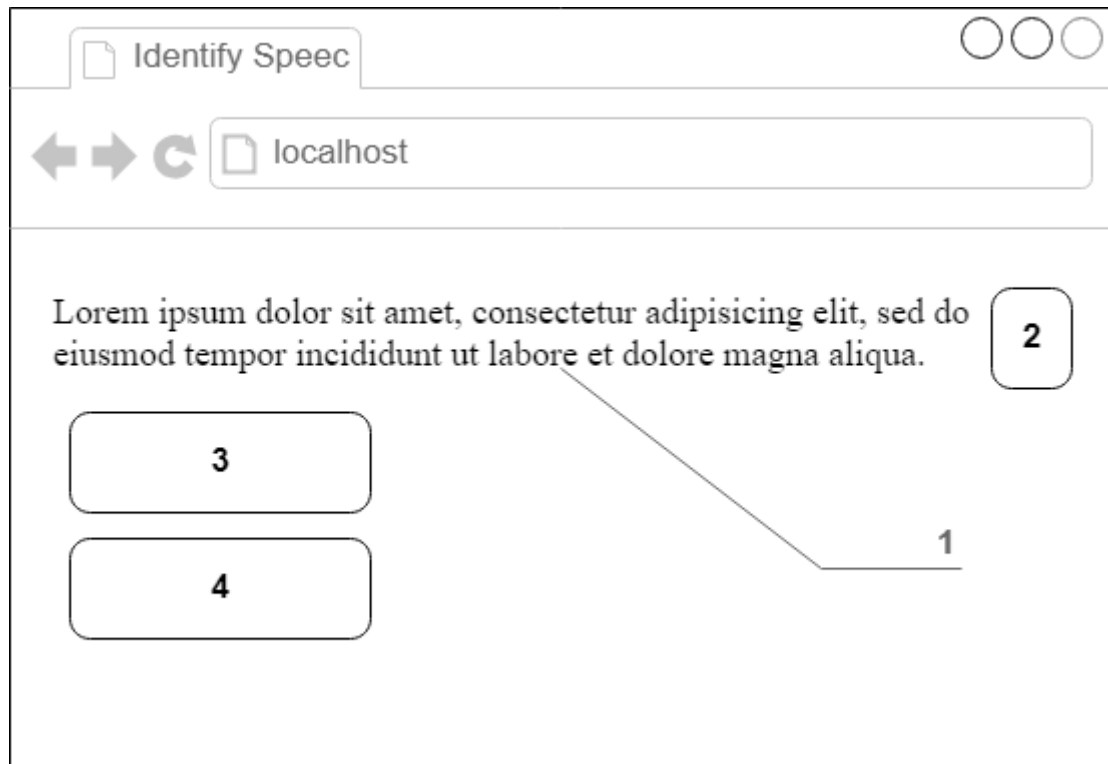


Figure 4.2 Identify Speech Screen.

Table 4.2 Identify Speech Screen Description.

No	Description
1	Random Sentences
2	Record Button
3	Identify Button
4	Finish Button

#### 4.1.3 Generate Speech Screen

Figure 4.3 shows the design layout for generate speech screen of the application. When user choose or click Select Speech ID, it will show the generate speech screen

that consist of 4 element such as ID Selector, Textbox Form, Generate and Play Button, and Finish Button. The description of the design layout is shown in Table 4.3.

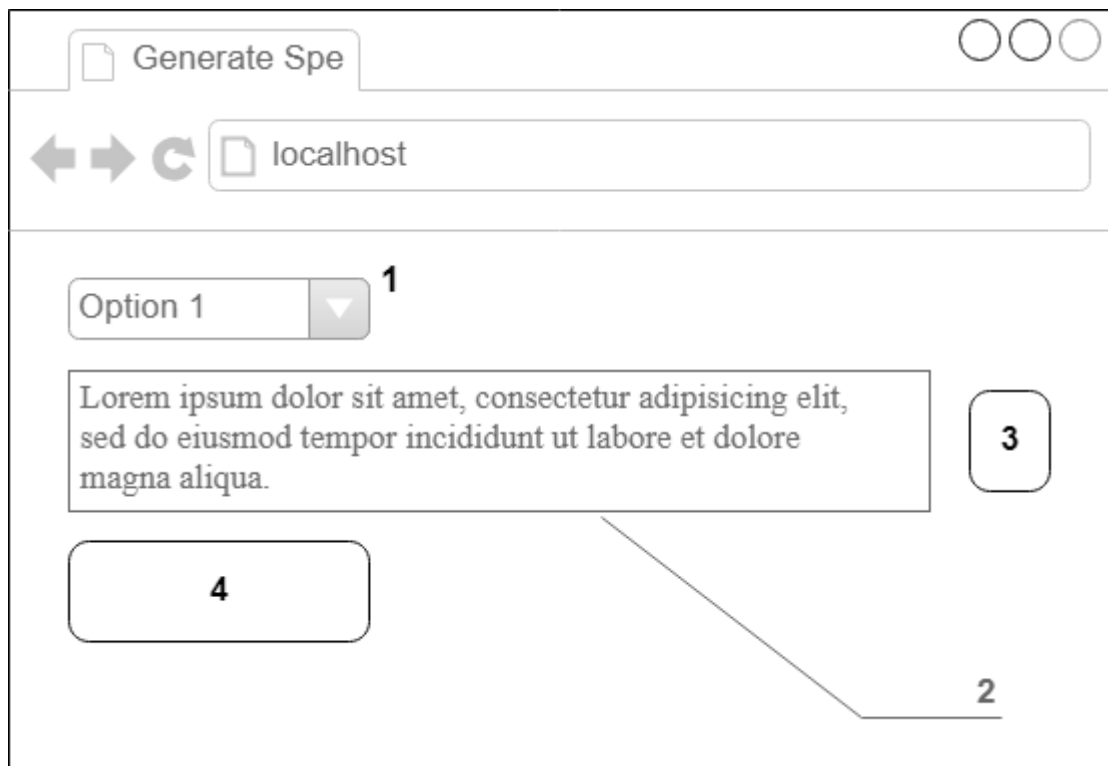


Figure 4.3 Generate Speech Screen.

Table 4.3 Generate Speech Screen Description.

No	Description
1	ID Selector
2	Textbox Form
3	Generate and Play Button
4	Finish Button

## 4.2 Class Diagram

*...Update according to final app...*

The class diagram is the structure of the system used toward this research. In this application there are 2 main category libraries, front-end and back-end. The detail of every libraries will be explained further below.

#### 4.2.1 *Front-end Library*

Front-end library is used as the application UI code. The front-end class diagram details are shown in Figure 4.4 Approximately there are 3 front-end classes as listed below.

##### 1. Home

Home class is class used to render Home screen. It consisting of 2 main methods, `identifyButton` and `generateButton`.

The `identifyButton` is method to start Make Speech ID activity which open Identify Speech screen. The `generateButton` is method to start Select Speech ID activity which open Generate Speech screen.

##### 2. IdentifySpeech

`IdentifySpeech` class is Make Speech ID activity. It is class used to render Identify Speech screen. It consisting of 4 main methods, `randomSentences`, `recordButton`, `identifyButton`, and `finishButton`.

The `randomSentences` is method to randomize sentence that used to sentence user need to be said. The `recordButton` is method to input user

recorded speech. The `identifyButton` is method to start identify user recorded speech and start `randomSentences` method again. The `finishButton` is method to finish Make Speech ID activity which open Home screen.

### 3. GenerateSpeech

`GenerateSpeech` class is Select Speech ID activity. It is class used to render Generate Speech screen. It consisting of 3 main methods, `selectID`, `generateButton`, and `finishButton`.

The `selectID` is method to select and define the speech ID from listed speech ID. The `generateButton` is method to generate and play speech based on inputted speech ID and text. The `finishButton` is method to finish Select Speech ID activity which open Home screen.

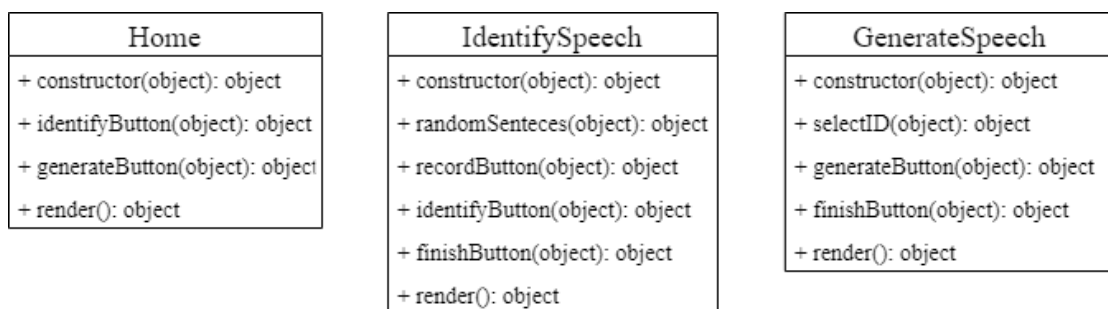


Figure 4.4 Front-end Class Diagram.

#### 4.2.2 *Back-end Library*

Back-end library is used as speech recognition process, speech synthesis process, create-read-update-delete (CRUD) database process and database collections. The back-end class diagram details are shown in Figure 4.5 Approximately there are 4 back-end classes as listed below.

1. Speech Recognition

`SpeechRecognition` class is class contains speech recognition algorithm. It containing `recognize` method to recognize speech from user and train speech.

2. Speech Synthesis

`SpeechSynthesis` class is class contains speech synthesis algorithm. It containing `generate` method to classify inputted text then generate the speech from classified text.

3. MongoDB

`MongoDb` class is class contains MongoDB process to established connection and data processing. It containing CRUD methods which is create, read, update, and delete data to MongoDB database.

4. MongoDB Collections

MongoDb collections is class collection that is used as data type. It containing 4 main collections, `SpeechCollection`, `TextCollection`, `SpeechDataCollection`, and `SpeechTrainedCollection`.



The **SpeechCollection** is collection used to store all speech data classify by grapheme and gender, and store speech file location and trained speech file location. It consisting 4 main variables, **grapheme**, **gender**, **speechFile**, and **trainedFile**. The **TextCollection** is collection used to store text data classify by word and grapheme. It consisting 2 main variables, **word** and **grapheme**.

The **SpeechDataCollection** is collection used to store speech data classify by speechID and gender, and store identified speech. It consisting 3 main variables, **speechID**, **gender**, and **speech**. The **SpeechTrainedCollection** is collection used to store trained speech data classify by grapheme and gender, and store identified speech. It consisting 3 main variable, **grapheme**, **gender**, and **data**.

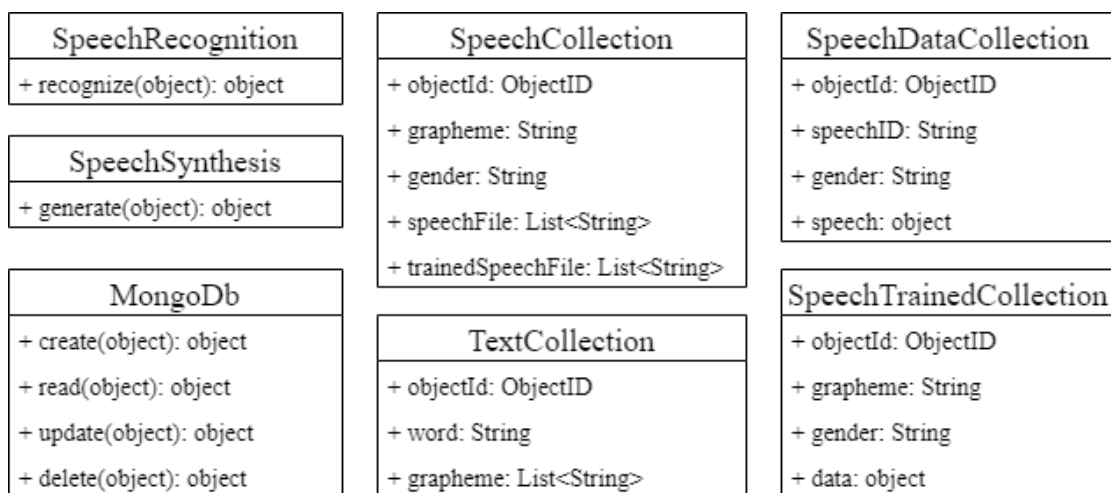


Figure 4.5 Bank-end Class Diagram.

## **CHAPTER V**

### **SYSTEM DEVELOPMENT**

#### **5.1 User Interface Development**

The UI is applied on Collect and Mimic application. There are two core pages in Collect application and three core pages in Mimic Application.

##### *5.1.1 Collect Application*

Two core pages in Collect application are Home page and Phonemes page.

##### *5.1.1.1 Home Page*

Home page is the home page of the Collect application. It shows information about how the collecting process works and information regarding to the collecting process. There are also 2 buttons, Start and Mimic Home button.

Start button will direct user to Phonemes page to start collecting. Mimic Home button will direct user to Mimic application Home page.



Figure 5.1 Collect Home page.

#### 5.1.1.2 Phonemes Page

Phonemes page is the main page of the Collect application. It shows information about which phoneme should be recorded and replayed the record. There are also 2 buttons, Record and Next button and an Audio element. Any response from the server is also printed in the console.

Record button will start the record process. When the record process is starting, Record button will be renamed to 'Say Now' which tell user that the record process is starting and indicate user to say the relevant phoneme. When the record process is finish, the recorded blob is uploaded to the server and Record button will be renamed to 'Record Again' which tell user that the record process is finish and user can do the next record. Next button will direct user to the next Phonemes page. In this application

there are 10 Phonemes. On the last phoneme, Next button will direct user back to Home page.

Audio element will be shown when the record process is finish with an audio from the record. User can hear the replay of the record from the Audio element.



Figure 5.2 Collect Phonemes page with 'a' phoneme.



Figure 5.3 Collect Phonemes page after record process is finish.

### 5.1.2 Mimic Application

Three core pages in Mimic application are Home page, Identify page, and Generate page.

#### 5.1.2.1 Home Page

Home page is the home page of the Mimic application. It shows information about how the mimic process works and information regarding to the identify and generate process. There are also 3 buttons, Identify, Generate, and Collect Home button.

Identify button will direct user to Identify page to start identifying speech. Generate button will direct user to Generate page to start generating speech. Collect Home button will direct user to Collect application Home page.



Figure 5.4 Mimic Home page.

#### 5.1.2.2 Identify Page

Identify page is the one of the main pages of the Mimic application. It shows information about which phoneme is recognizable and form of input and audio that will be used to identify the speech and update it to database. There are also 3 buttons, Record, Identify Speech and Finish button, and an Input text and an Audio element. Any response from the server is also printed in the console.

Record button will start the record process. When the record process is starting, Record button will be renamed to 'Say Now' which tell user that the record process is starting and indicate user to say the phoneme. When the record process is finish, the recorded blob is uploaded to the server Record button will be renamed to 'Record Again' which tell user that the record process is finish and user can do the next record.

Identify Speech button will take the value of Input element and file source of Audio element, then, send them to the server to be identify. The indication of identify process in the server is success or error is there is alert on the browser that show any information regarding to identify process in the server. If there is no value of Input element or file source of Audio element, Identify Speech button will alert user that there is not any value. Finish button will direct user back to the Home page.

Audio element will be shown when the record process is finish with an audio from the record. User can hear the replay of the record from the Audio element. Input element is used as input of user speech ID value.

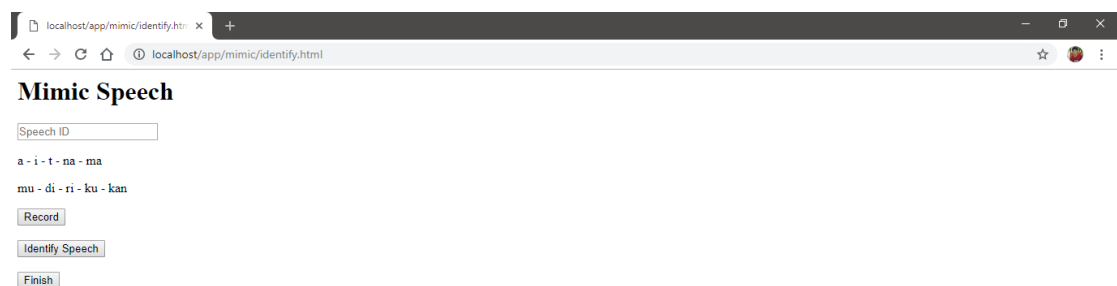


Figure 5.5 Mimic Identify page.

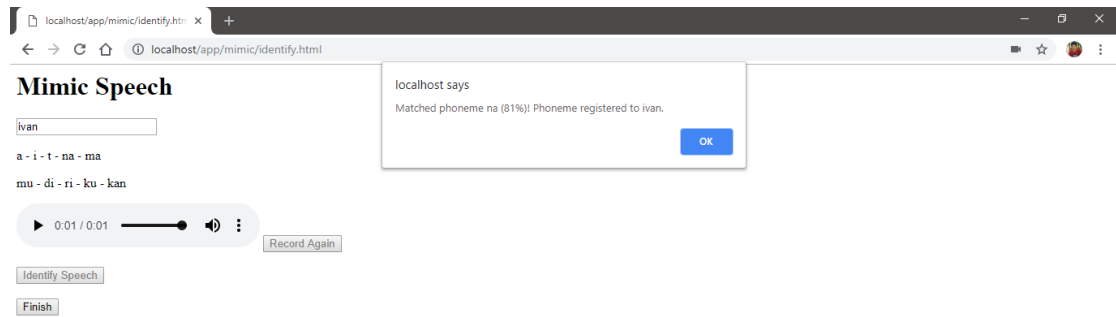


Figure 5.6 Mimic Identify page alert user indicating identify process in the server is finish.

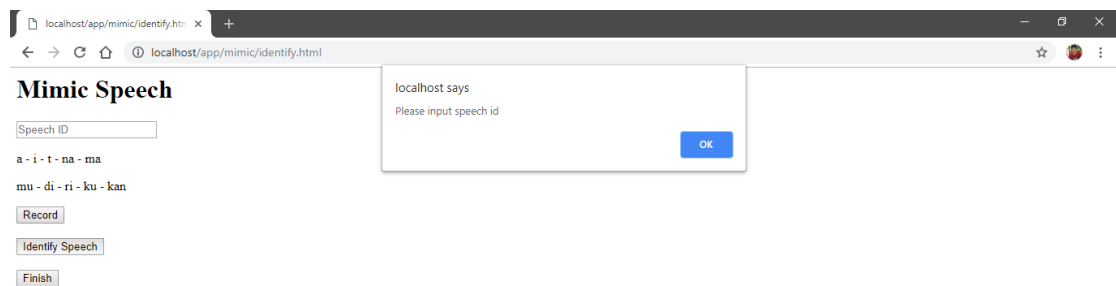


Figure 5.7 Mimic Identify page alert user indicating there is no value in Input element.



### 5.1.2.3 *Generate Page*

Generate page is the one of the main pages of the Mimic application. It shows form of data list and text area that will be used to generate the speech. When the page is loading, it also loads speech data from the database and then store it on data list. If there is no speech data found, it will alert user and direct user back to Home page. There are also 2 buttons, Generate Speech and Finish button, and a Datalist element, a Textarea element and an Audio element. Any response from the server is also printed in the console.

Generate Speech button will take the value of Datalist element and Textarea element, then, send them to the server to be generate. The indication of generate process in the server is error is there is alert on the browser that show any error information regarding to generate process in the server. The indication of generate process in the server is success is the Audio element is shown with an audio from the generated speech. If there's no value of Datalist element or Textarea element, Identify Speech button will alert user that there is not any value. Finish button will direct user back to the Home page.

Audio element will be shown when the generate process is success with an audio from the generated speech. User can hear the generated speech from the Audio element. Datalist element is shows all speech data from database and it is used as input of user speech ID value. Textarea element is used as input of words that user want to generate.

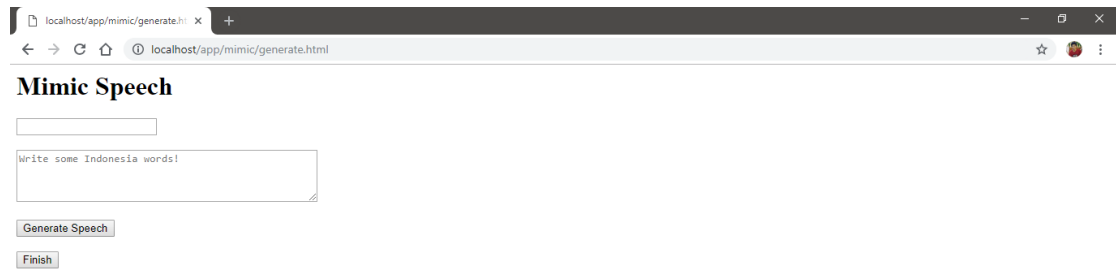


Figure 5.8 Mimic Generate page.

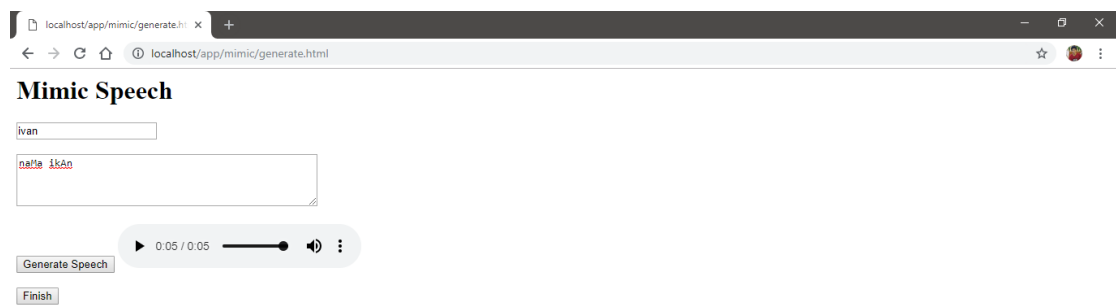


Figure 5.9 Mimic Generate page play the generate speech after generate process success.

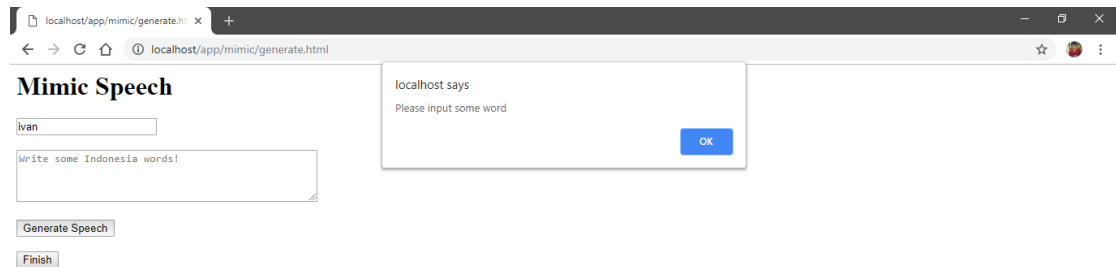


Figure 5.10 Mimic Generate page alert user indicating there is no value in Textarea element.

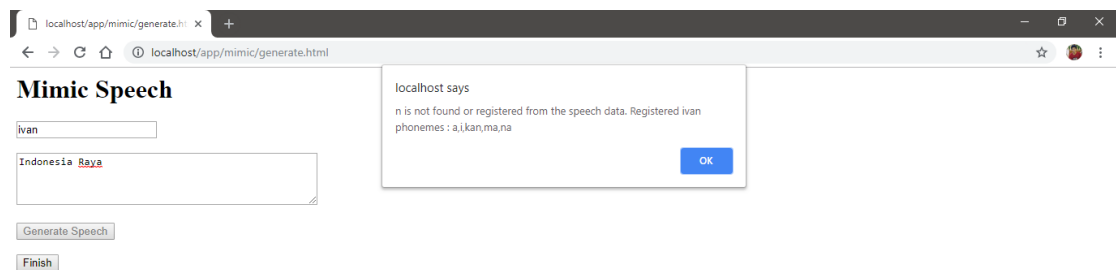


Figure 5.11 Mimic Generate page alert user indicating there is error in the generate process.

## 5.2 Application Details

The application details explain on server code, router code, database code, Train application code and library package that is used in the application. The figures that show the code also included comment that help explanation for corresponding line.

### 5.2.1 *Server Code*

The server code is used to create http server in the computer. It contains `serverHandler` and create server command itself.

`serverHandler` handle all request to the server each response corresponding to each request. There are 3 type main router `collect`, `identify`, `generate` to handle any request corresponding to them. When the request match it returns 200 alongside with the response corresponding to the router. When the request doesn't match any of those, it checks if the file URL as `filename` is existed in working directory or not. When the `filename` not found it return 404 Not Found. When found it read the `file` and return it with 200 so that it can be displayed by the browser. It returns 500 if error occurred when reading the `file`. used to create http server in the computer core pages in Mimic application are Home page, Identify page, and Generate page.

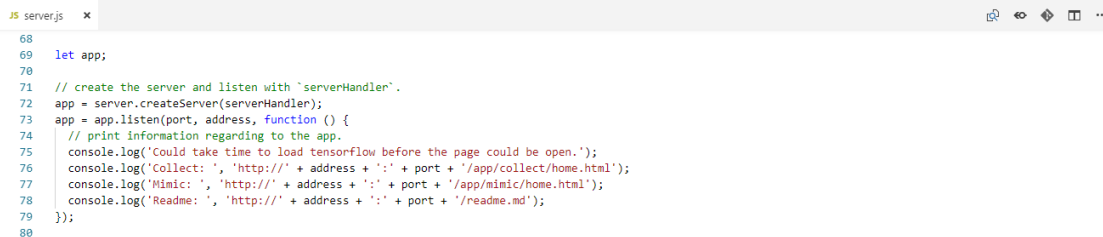
```

9
10
11 /**
12  * `serverHandler` handle all request to the server each response corresponding to each request.
13  * There are 3 type main router `collect`, `identify`, `generate` to handle any request corresponding to them.
14  * When the request match it return 200 alongside with the response corresponding to the router.
15  * When the request doesn't match any of those, it check if the file url as `filename` is existed in working
16  * directory or not. When the `filename` not found it return 404 Not Found. When found it read the `file` and return it with 200
17  * so that it can be displayed by the browser. It return 500 if error occurred when reading the `file`.
18  */
19 async function serverHandler(request, response) {
20   const uri = url.parse(request.url).pathname;
21   const filename = path.join(process.cwd(), uri);
22
23   const collect = await require('./router/collect')(address, port, filename, request);
24   if (typeof collect !== 'undefined') {
25     response.writeHead(200);
26     response.write(collect);
27     response.end();
28   }
29
30   const identify = await require('./router/identify')(address, port, filename, request);
31   if (typeof identify !== 'undefined') {
32     response.writeHead(200);
33     response.write(identify);
34     response.end();
35   }
36
37   const generate = await require('./router/generate')(address, port, filename, request);
38   if (typeof generate !== 'undefined') {
39     response.writeHead(200);
40     response.write(generate);
41     response.end();
42   }
43
44   fs.exists(filename, function (exists) {
45     if (!exists) {
46       response.writeHead(404);
47       response.write('404 Not Found: ' + filename + '\n');
48       response.end();
49     }
50
51     fs.readFile(filename, 'binary', function (err, file) {
52       if (err) {
53         response.writeHead(500);
54         response.write(err + '\n');
55         response.end();
56       }
57
58       response.writeHead(200);
59       response.write(file, 'binary');
60       response.end();
61     });
62   });
63 }
64
65
66
67
68

```

Figure 5.12 serverHandler on Server code.

Create server command accept `serverHandler` as the parameter. It will listen on defined variable `localhost` and `port`. When the server established, it prints information in console.



```

68
69 let app;
70
71 // create the server and listen with `serverHandler`.
72 app = server.createServer(serverHandler);
73 app = app.listen(port, address, function () {
74   // print information regarding to the app.
75   console.log('Could take time to load tensorflow before the page could be open.');
```

Figure 5.13 Create server command on Server code.

### 5.2.2 Router Code

The router is used to handle categorized request from the server. It split into 3 categorized files, collect, identify, and generate.

Collect router handle all collect app request, each response corresponding to each request. There are a lot of request categorized by collect app but all of them is basically uploadCollect. When the request doesn't match anything, it returns nothing. uploadCollect upload file in the request based on defined directory and then convert its path location to URL. It returns the URL as fileURL.

```

3
4  /**
5   * `uploadCollect` upload file in the `request` based on defined directory and
6   * then convert its path location to url.
7   *
8   * Return the url as `fileURL`.
9   */
10 function uploadCollect(address, port, label, request) {
11   const dirData = '\\data\\';
12   const dirCollect = '\\collect\\';
13   const dirLabel = '\\ ' + label + '\\';
14
15   const form = new formidable.IncomingForm();
16   form.uploadDir = process.cwd() + dirData + dirCollect + dirLabel;
17   form.keepExtensions = true;
18   form.maxFieldsSize = 10 * 1024 * 1024;
19   form.maxFields = 1000;
20   form.multiples = false;
21
22   return new Promise(function (resolve) {
23     form.parse(request, function (err, fields, files) {
24       const file = util.inspect(files);
25
26       const fileName = file.split('path:')[1].split('\\')[0].split(dirData)[1].split(dirCollect)[1].split(dirLabel)[1].toString().replace(/\\/g, '').replace(
27         (/\\/g, ''));
28       const fileURL = 'http://' + address + ':' + port + '/data/collect/' + label + '/' + fileName;
29
30       console.log('fileURL: ', fileURL);
31       resolve(JSON.stringify({ fileURL: fileURL }));
32     });
33   });
34 }

```

Figure 5.14 uploadCollect on Collect Router code.

```

34
35  /**
36   * `router` handle all collect app request, each response corresponding to each request.
37   * There are a lot of request categorized by collect app but all of them is basically `uploadCollect`.
38   * When the request doesn't match anything it return nothing.
39   */
40 async function router(address, port, filename, request) {
41   let response;
42
43   if (request.method === 'POST') {
44     if (filename.toString().indexOf('\\uploadCollectA') !== -1) response = await uploadCollect(address, port, 'a', request);
45     else if (filename.toString().indexOf('\\uploadCollectI') !== -1) response = await uploadCollect(address, port, 'i', request);
46     else if (filename.toString().indexOf('\\uploadCollectT') !== -1) response = await uploadCollect(address, port, 't', request);
47     else if (filename.toString().indexOf('\\uploadCollectNa') !== -1) response = await uploadCollect(address, port, 'na', request);
48     else if (filename.toString().indexOf('\\uploadCollectMa') !== -1) response = await uploadCollect(address, port, 'ma', request);
49     else if (filename.toString().indexOf('\\uploadCollectMu') !== -1) response = await uploadCollect(address, port, 'mu', request);
50     else if (filename.toString().indexOf('\\uploadCollectDi') !== -1) response = await uploadCollect(address, port, 'di', request);
51     else if (filename.toString().indexOf('\\uploadCollectRi') !== -1) response = await uploadCollect(address, port, 'ri', request);
52     else if (filename.toString().indexOf('\\uploadCollectKu') !== -1) response = await uploadCollect(address, port, 'ku', request);
53     else if (filename.toString().indexOf('\\uploadCollectKan') !== -1) response = await uploadCollect(address, port, 'kan', request);
54   }
55
56   return response;
57 }
58

```

Figure 5.15 router on Collect Router code.

Identify router handle all mimic app identify section request, each response corresponding to each request. There are 2 type requests categorized by mimic app identify section, uploadSpeech and identifySpeech. When the request doesn't match anything, it returns nothing. uploadSpeech upload file in the request based

on defined directory and then convert its path location to URL. It returns the URL and path location as `fileURL` and `filePath`. `identifySpeech` load latest saved trained model from MongoDB `models` collection. Use it to identify the speech file from `filePath` by `extract` then feed it to `model`. Matched speech will register the phoneme to corresponding name and update it on MongoDB. Error connection will return message. Not found model will return message. Unsatisfied identification, results below 0.75 (75%) from `model` will return message. None of those is occurred will update the MongoDB `speechDatas` collection on its name and identified phoneme `filePath` and return message. The returned message whether just information or even error is return as `status`. `extract` is support function decode wav based on `filePath`, frame it, extract by FFT and MFCC and convert it to TensorFlow input shape. It returns 4d tensor.

```

12
13 /**
14  * 'extract' decode wav based on 'filePath', frame it, extract by FFT and MFCC and
15  * convert it to tensorflow input shape.
16  *
17  * Return 4d tensor.
18  */
19 function extract(filePath) {
20   const sampleRate = 16000;
21   const frameSize = 25 / 1000 * sampleRate; // 400
22   const frameShift = 10 / 1000 * sampleRate; // 160
23   const fftSize = 256; // based on the next power of 2 from 400 -> 512 as the input.
24   const melCount = 40;
25   const lowHz = 300;
26   const highHz = 8000;
27
28   const mfcc = require('node-mfcc/src/mfcc').construct(fftSize, melCount, lowHz, highHz, sampleRate);
29
30   // decode wav.
31   const sample = wav.decode(fs.readFileSync(filePath)).channelData[0];
32
33   // framing the sample.
34   const frameSample = [];
35   for (i = 0; i <= sample.length - frameSize; i = i + frameShift) {
36     const frame = new Float32Array(fftSize * 2);
37     for (j = 0; j < frame.length; j++) {
38       if (j < frameSize) frame[j] = sample[i + j];
39       else frame[j] = 0;
40     }
41     frameSample.push(frame);
42   }
43
44   // extract by FFT and MFCC.
45   const extractedSample = [];

```



```

43
44 // extract by FFT and MFCC.
45 const extractedSample = [];
46 for (i = 0; i < frameSample.length; i++) {
47   const fftSample = fft(frameSample[i]);
48   extractedSample.push(new Float32Array(mfcc(fftSample)));
49 }
50
51 // define tensorflow input shape.
52 const time = extractedSample.length;
53 const freq = extractedSample[0].length;
54 const data = new Float32Array(time * freq);
55 const shape = [1, time, freq, 1];
56
57 // convert extracted wav to tensorflow input data.
58 for (i = 0; i < time; i++) {
59   const melFreq = extractedSample[i];
60   const offset = i * freq;
61
62   data.set(melFreq, offset);
63 }
64
65 // return 4d tensor.
66 return tf.tensor4d(data, shape);
67 }
68

```

Figure 5.16 extract on Identify Router code.

```

68
69 /**
70  * 'uploadSpeech' upload file in the 'request' based on defined directory and
71  * then convert its path location to url.
72  *
73  * Return the url and path location as 'fileURL' and 'filePath'.
74  */
75 function uploadSpeech(address, port, request) {
76   const dirData = '\\data\\';
77   const dirSpeech = '\\speech\\';
78
79   const form = new formidable.IncomingForm();
80   form.uploadDir = process.cwd() + dirData + dirSpeech;
81   form.keepExtensions = true;
82   form.maxFieldsSize = 10 * 1024 * 1024;
83   form.maxFields = 1000;
84   form.multiples = false;
85
86   return new Promise(function (resolve) {
87     form.parse(request, function (err, fields, files) {
88       const file = util.inspect(files);
89
90       const fileName = file.split('path:')[1].split('\\')[0].split(dirData)[1].split(dirSpeech)[1].toString().replace(/\\/g, '').replace(/\\/g, '');
91       const fileURL = 'http://' + address + ':' + port + '/data/speech/' + fileName;
92       const filePath = files.file.path;
93
94       console.log('fileURL: ', fileURL);
95       console.log('filePath: ', filePath);
96       resolve(JSON.stringify({ fileURL: fileURL, filePath: filePath }));
97     });
98   });
99 }
100

```

Figure 5.17 uploadSpeech on Identify Router code.

```

100 JS identify.js
101 /**
102  * `identifySpeech` load latest saved trained model from MongoDB `models` collection.
103  * Use it to identify the speech file from `filePath` by `extract` then feed it to `model`.
104  * Matched speech will register the phoneme to corresponding `name` and update it on MongoDB.
105  *
106  * Error connection will return message. Not found model will return message. Unsatisfied identification,
107  * results below 0.75 (75%) from `model` will return message. None of those is occurred will update the MongoDB
108  * `speechDatas` collection on its `name` and identified phoneme `filePath` and return message. The returned
109  * message wheter just information or even error is return as `status`.
110  */
111 function identifySpeech(request) {
112   const form = new formidable.IncomingForm();
113
114   return new Promise(function (resolve) {
115     form.parse(request, async function (err, fields, files) {
116       // take `name` and `filePath` data from POST request.
117       const name = fields.name;
118       const filePath = fields.filePath;
119
120       let status;
121       let resultDB = {};
122
123       // load model location from MongoDB.
124       try {
125         model.models.db = await connection.connect();
126         resultDB = await model.models.findOne({}).sort({ createdAt: -1 }).select('location _id').exec();
127         await connection.disconnect();
128
129         // return if nothing is found.
130         if (typeof resultDB.location === 'undefined') {
131           status = 'No train model found.';
132           console.log('Status:', status);
133
134           if (typeof resultDB.location === 'undefined') {
135             status = 'No train model found.';
136             console.log('Status:', status);
137
138             resolve(JSON.stringify({ status: status }));
139           }
140         } catch (err) { // return if connection error is occurred.
141           status = err.errmsg;
142           console.log('Status:', status);
143
144           resolve(JSON.stringify({ status: status }));
145         }
146
147         // load the model as `tfModel` from model location and extract speech as `data` from `filePath`.
148         const locationURL = resultDB.location;
149         const tfModel = await tf.loadModel(locationURL + '/model.json');
150         const data = extract(filePath);
151
152         // identify the speech and process data from it.
153         const prediction = tfModel.predict(data).dataSync();
154         let predictionLabel = 0;
155         for (i = 0; i < prediction.length; i++) {
156           if (prediction[predictionLabel] <= prediction[i]) predictionLabel = i;
157         }
158         const label = labels[predictionLabel]
159         const labelData = prediction[predictionLabel]
160
161         // return if result unsatisfied, below 75%.
162         if (labelData < 0.75) {
163           status = 'No phonemes matched.';
164           console.log('Status:', status);
165
166           resolve(JSON.stringify({ status: status }));
167         }
168       }
169     });
170   });
171 }

```

```

159     status = 'No phonemes matched.';
160     console.log('Status:', status);
161
162     resolve(JSON.stringify({ status: status }));
163   } else {
164     // find if 'name' is already registered in MongoDB or not.
165     try {
166       model.speechDatas.db = await connection.connect();
167       resultDB = await model.speechDatas.findOne({ name: name }).select('phonemes _id').exec();
168
169       // create the document if not found, update if found before update to MongoDB.
170       let updateData = {};
171       if (resultDB === null) {
172         updateData.name = name;
173         updateData.phonemes = {};
174         updateData.phonemes[label] = filePath
175       } else {
176         updateData = resultDB;
177         updateData.phonemes[label] = filePath
178       }
179       await model.speechDatas.updateOne({ name: name }, updateData, { upsert: true });
180
181       await connection.disconnect();
182     } catch (err) { // return if connection error is occurred.
183       status = err.errmsg;
184       console.log('Status:', status);
185
186       resolve(JSON.stringify({ status: status }));
187     }
188
189     status = 'Matched phoneme ' + label + ' (' + (labelData.toFixed(2) / 1 * 100) + '%)! Phoneme registered to ' + name + '.';
190     console.log('Status:', status);
191     status = 'Matched phoneme ' + label + ' (' + (labelData.toFixed(2) / 1 * 100) + '%)! Phoneme registered to ' + name + '.';
192     console.log('Status:', status);
193
194     // return success identify speech information.
195     resolve(JSON.stringify({ status: status }));
196   }
197 });
198 }
199

```

Figure 5.18 identifySpeech on Identify Router code.

```

JS identify.js x
199 /**
200  * "router" handle all mimic app identify section request, each response corresponding to each request.
201  * There are 2 type request categorized by mimic app identify section, 'uploadSpeech' and 'identifySpeech'.
202  * When the request doesn't match anything it return nothing.
203  */
204
205 async function router(address, port, filename, request) {
206   let response;
207
208   if (request.method === 'POST') {
209     if (filename.toString().indexOf('\\uploadSpeech') !== -1) response = await uploadSpeech(address, port, request);
210     if (filename.toString().indexOf('\\identifySpeech') !== -1) response = await identifySpeech(request);
211   }
212
213   return response;
214 }
215

```

Figure 5.19 router on Identify Router code.

generate router handle all mimic app generate section request, each response corresponding to each request. There are 2 type requests categorized by mimic app identify section, loadSpeech and generateSpeech. When the request doesn't match anything, it returns nothing. loadSpeech load all registered name from MongoDB

speechDatas collection. It returns the array of name as data and error message if nothing is found or error. generateSpeech load registered phonemes from MongoDB speechDatas collection based on name. Extract words as extractWord. Decode all corresponding extractWord based on the database result. them and encode to wav file and take its fileURL location. It returns the URL as fileURL and error message if connection error or unregistered phoneme is found.

```

JS generate.js x
7
8  /**
9   * `loadSpeech` load all registered name from MongoDB `speechDatas` collection.
10  *
11  * Return the array of name as `data` and `error` message if nothing is found or
12  * connection error.
13  */
14  async function loadSpeech() {
15    let resultDB = [];
16    let error;
17
18    try {
19      model.speechDatas.db = await connection.connect();
20      resultDB = await model.speechDatas.find({}).select('name -_id').exec();
21      await connection.disconnect();
22
23      if (resultDB.length === 0) {
24        error = 'No speech data found.';
25        console.log('Error:', error);
26      }
27    } catch (err) {
28      error = err.errmsg;
29      console.log('Error:', error);
30    }
31
32    const data = resultDB.map((data) => data.name);
33
34    console.log('speechData:', data);
35    return JSON.stringify({ speechData: data, error: error });
36  }
37

```

Figure 5.20 loadSpeech on Generate Router code.

```

37
38
39 /**
40  * `generateSpeech` load registered phonemes from MongoDB `speechDatas` collection based on `name`.
41  * Extract `words` as `extractWord`. Decode all corresponding `extractWord` based on the database result.
42  * Combine them and encode to wav file and take its `fileURL` location.
43  *
44  * Return the url as `fileURL` and `error` message if connection error or unregistered phoneme is found.
45  */
46 function generateSpeech(address, port, request) {
47   const form = new formidable.IncomingForm();
48
49   return new Promise(function (resolve) {
50     form.parse(request, async function (err, fields, files) {
51       // take `name` and `words` data from POST request.
52       const name = fields.name;
53       const words = fields.words.toLowerCase();
54
55       let fileURL;
56       let error;
57       let resultDB = {};
58       let buffer;
59
60       // load registered phonemes from MongoDB.
61       try {
62         model.speechDatas.db = await connection.connect();
63         resultDB = await model.speechDatas.findOne({ name: name }).select('phonemes _id').exec();
64         await connection.disconnect();
65       } catch (err) { // return if connection error is occurred.
66         error = err.errmsg;
67         console.log('Error:', error);
68
69         resolve(JSON.stringify({ fileURL: fileURL, error: error }));
70       }
71
72       resolve(JSON.stringify({ fileURL: fileURL, error: error }));
73     });
74
75     // take just the phonemes from `resultDB`.
76     const phonemes = [];
77     for (phoneme in resultDB.phonemes) phonemes.push(phoneme);
78
79     // extract the `words`.
80     const extractWord = [];
81     for (i = 0; i < words.length; i++) {
82       if (words[i] === ' ') { // check if it is space.
83         extractWord.push(' ');
84       } else if (i + 3 <= words.length && phonemes.includes(words.substring(i, i + 3))) { // check if it the next 3 char is in `phonemes`.
85         extractWord.push(words.substring(i, i + 3));
86         i = i + 2;
87       } else if (i + 2 <= words.length && phonemes.includes(words.substring(i, i + 2))) { // check if it the next 2 char is in `phonemes`.
88         extractWord.push(words.substring(i, i + 2));
89         i = i + 1;
90       } else if (phonemes.includes(words.substring(i, i + 1))) { // check if it the next char is in `phonemes`.
91         extractWord.push(words.substring(i, i + 1));
92       } else { // return if unregistered phoneme is found.
93         error = words[i] + ' is not found or registered from the speech data. ' +
94           'Registered ' + name + ' phonemes : ' + phonemes;
95         console.log('Error:', error);
96
97         resolve(JSON.stringify({ fileURL: fileURL, error: error }));
98       }
99     }
100
101     // combine each `extractWord` into 1 `channelData`.
102     const sampleRate = 16000;
103     const channelData = [new Float32Array(0)];
104     // iterate each `extractWord`.
105     for (i = 0; i < extractWord.length; i++) {

```

```

100 // iterate each `extractWord`.
101 for (i = 0; i < extractWord.length; i++) {
102   const prevAmpli = channelData[0];
103   const offset = channelData[0].length;
104
105   let nextAmpli;
106   if (extractWord[i] === ' ') { // if space it fill the `nextAmpli` with 0 `sampleRate` times.
107     nextAmpli = new Float32Array(sampleRate).map(() => 0);
108   } else { // else it fill from the decode result based on `resultDB` and corresponding `extractWord`.
109     buffer = fs.readFileSync(resultDB.phonemes[extractWord[i]]);
110     nextAmpli = wav.decode(buffer).channelData[0].slice(0, sampleRate);
111   }
112
113   // arrange the ampli and redefined the `channelData`.
114   const ampli = new Float32Array(sampleRate + offset);
115   ampli.set(prevAmpli);
116   ampli.set(nextAmpli, offset);
117   channelData[0] = ampli;
118 }
119
120 // define upload directory.
121 const dirData = '\\data\\';
122 const dirGenerate = '\\generate\\';
123 const fileName = 'upload_' + cryptoRandomString(32) + '.wav';
124 const file = process.cwd() + dirData + dirGenerate + fileName;
125
126 // encode wav based on `channelData` and `sampleRate`
127 buffer = wav.encode(channelData, { sampleRate: sampleRate });
128 fs.writeFileSync(file, Buffer.from(buffer));
129
130 // return the file url as `fileURL`.
131 fileURL = 'http://' + address + ':' + port + '/data/generate/' + fileName;
132 console.log('fileURL:', fileURL);
133 // return the file url as `fileURL`.
134 fileURL = 'http://' + address + ':' + port + '/data/generate/' + fileName;
135 console.log('fileURL:', fileURL);
136 resolve(JSON.stringify({ fileURL: fileURL }));
137 }
138 }
139 }

```

Figure 5.21 generateSpeech on Generate Router code.

```

JS generate.js
137
138 /**
139  * `router` handle all mimic app generate section request, each response corresponding to each request.
140  * There are 2 type request categorized by mimic app identify section, `loadSpeech` and `generateSpeech`.
141  * When the request doesn't match anything it return nothing.
142  */
143 async function router(address, port, filename, request) {
144   let response;
145
146   if (request.method === 'GET') {
147     if (filename.toString().indexOf('\\loadSpeech') !== -1) response = await loadSpeech();
148   } else if (request.method === 'POST') {
149     if (filename.toString().indexOf('\\generateSpeech') !== -1) response = await generateSpeech(address, port, request);
150   }
151
152   return response;
153 }
154


```

Figure 5.22 router on Generate Router code.

### 5.2.3 Database Code

The database code is used to create connection to MongoDB and create MongoDB collection model scheme. It contains connection and model.

`connection` define 2 functions `connect` and `disconnect`. `connect` establish connection to MongoDB `mimic_speech` database on localhost. Established connection is saved in `connection`. `disconnect` close established connection to MongoDB.



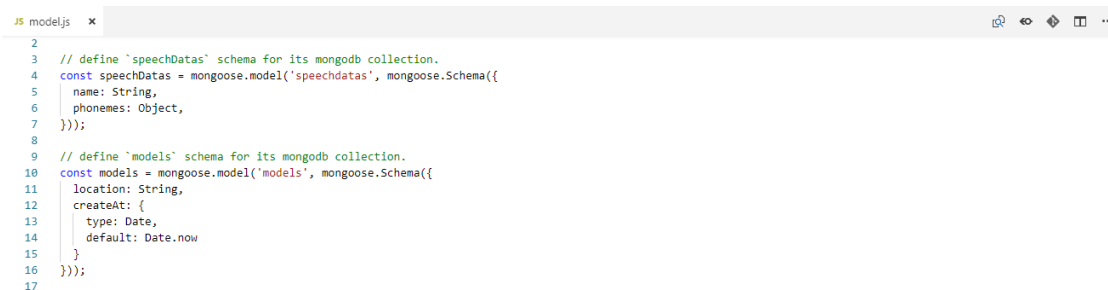
```

4
5
6  /**
7   * `connect` establish connection to mongodb mimic_speech database on localhost.
8   *
9   * Established connection is saved in `connection`.
10  */
11  async function connect() {
12    connection = await mongoose.connect('mongodb://localhost/mimic_speech', { useNewUrlParser: true });
13    return connection;
14  };
15
16  /**
17   * `disconnect` close established `connection` to mongodb.
18   */
19  async function disconnect() {
20    await connection.disconnect();
21  }

```

Figure 5.23 connection on Database code.

`model` define 2 schemes variable `speechDatas` and `models`. `speechDatas` define scheme with `name` as string and `phonemes` as object. `models` define scheme with `location` as string and `createAt` as date with current date being default value when create or update collection when `createAt` is not defined.



```

2
3  // define `speechDatas` schema for its mongodb collection.
4  const speechDatas = mongoose.model('speechdatas', mongoose.Schema({
5    name: String,
6    phonemes: Object,
7  }));
8
9  // define `models` schema for its mongodb collection.
10 const models = mongoose.model('models', mongoose.Schema({
11   location: String,
12   createdAt: {
13     type: Date,
14     default: Date.now
15   }
16 }));
17

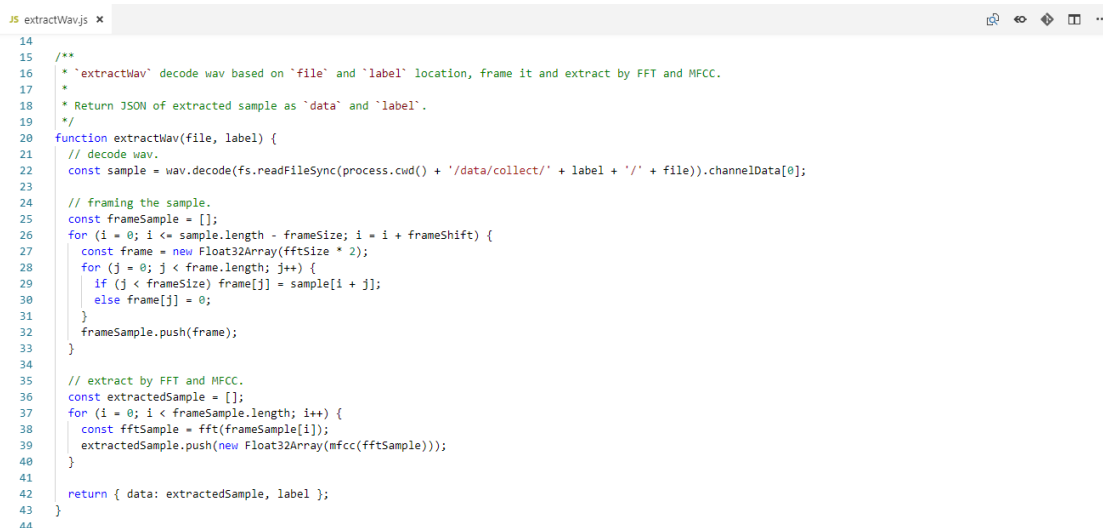
```

Figure 5.24 model on Database code.

### 5.2.4 Train Application Code

Train application code is used to train collected data and save the model on n iteration and last iteration so that later can be used to identify the speech. It contains `extractWav`, `loadData`, `nextBatch`, `saveModelDB` as support code and `model` and `train` as the main code.

`extractWav` decode wav based on `file` and `label` location, frame it and extract by FFT and MFCC. It returns JSON of extracted sample as `data` and `label`.



```

14
15 /**
16  * 'extractWav' decode wav based on 'file' and 'label' location, frame it and extract by FFT and MFCC.
17  *
18  * Return JSON of extracted sample as 'data' and 'label'.
19  */
20 function extractWav(file, label) {
21   // decode wav.
22   const sample = wav.decode(fs.readFileSync(process.cwd() + '/data/collect/' + label + '/' + file)).channelData[0];
23
24   // framing the sample.
25   const frameSample = [];
26   for (i = 0; i <= sample.length - frameSize; i = i + frameShift) {
27     const frame = new Float32Array(fftSize * 2);
28     for (j = 0; j < frame.length; j++) {
29       if (j < frameSize) frame[j] = sample[i + j];
30       else frame[j] = 0;
31     }
32     frameSample.push(frame);
33   }
34
35   // extract by FFT and MFCC.
36   const extractedSample = [];
37   for (i = 0; i < frameSample.length; i++) {
38     const fftSample = fft(frameSample[i]);
39     extractedSample.push(new Float32Array(mfcc(fftSample)));
40   }
41
42   return { data: extractedSample, label };
43 }
44

```

Figure 5.25 `extractWav` on Train application code.

`loadData` is used to load all collected data, then, extract them and split them into train data, validation data, and test data. http server in the computer. It contains `extractFiles` and `load` and `split` command itself. `extractFiles` extract wav from loaded files as `load`, assigned to corresponding `data`. It returns array of extracted



waves. Every iteration when extracting the loaded files is printed. If inconsistent data is met, the process stops and exit.

```

9
10
11 /**
12  * 'extractFiles' extract wav from loaded files as 'load', assigned to corresponding 'data'.
13  *
14  * Return array of extracted waves. Every iteration when extracting the loaded files is printed.
15  * If inconsistent data is met, the process stop and exit.
16  */
17 function extractFiles(load, data, type) {
18   const data = [];
19   let lastSpectrogram;
20
21   // shuffle 'load' to avoid overfitting model.
22   load = shuffle(load);
23
24   for (k = 0; k < load.length; k++) {
25     const spectrogram = extractWav(load[k].file, load[k].label);
26     if (k !== 0 && lastSpectrogram !== 'undefined') {
27       && lastSpectrogram.data.length !== spectrogram.data.length {
28         console.log(load[k - 1].file + ' (label : ' + load[k - 1].label + ') ' +
29           ' is different with data [time] ' +
30           load[k].file + ' (label : ' + load[k].label + ') . ' +
31           'Please check again. Default time is 1 second. Process exiting. ');
32         process.exit();
33       } else {
34         data.push(spectrogram);
35       }
36     }
37     lastSpectrogram = spectrogram;
38     console.log(new Date().toISOString() + ' | Step ' + (k + 1) + ' | Extracting ' + type + ' Wav: ' + load[k].file);
39   }
40   return data;
41 }
42

```

Figure 5.26 extractFiles on loadData Train application code.

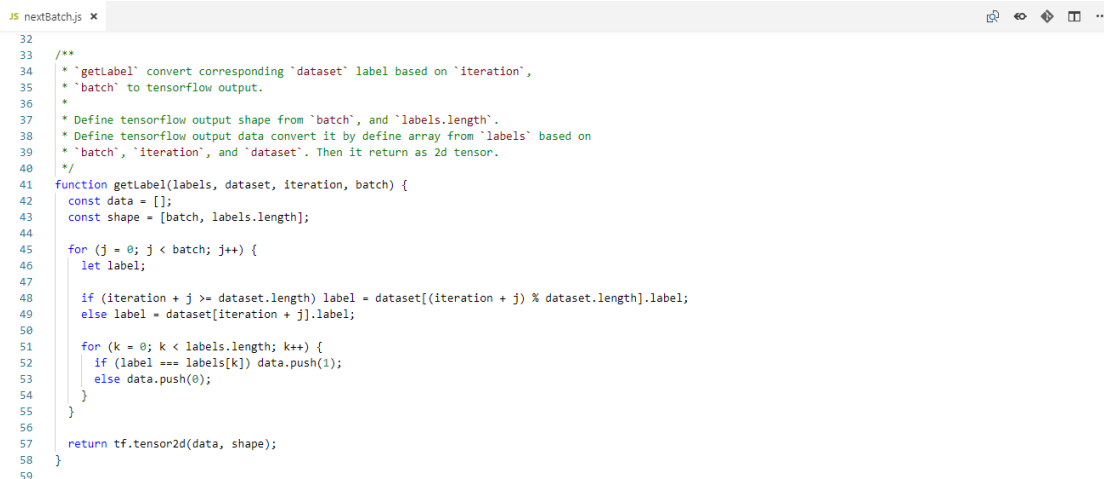
```

42
43 // iterate each 'labels' to load files. Every iteration when loading the files is printed.
44 for (i = 0; i < labels.length; i++) {
45   console.log(new Date().toISOString() + ' | Step ' + (i + 1) + ' | Load files from label: ' + labels[i]);
46
47   // shuffle 'files' to avoid overfitting model.
48   const files = shuffle(fs.readdirSync(process.cwd() + '/data/collect/' + labels[i]));
49
50   // 80% train data, 10% validation data, 10% test data.
51   const validation = Math.floor(10 / 100 * files.length);
52   const test = Math.floor(10 / 100 * files.length);
53   const train = files.length - validation - test;
54   for (j = 0; j < files.length; j++) {
55     const data = {
56       file: files[j],
57       label: labels[i]
58     };
59
60     if (j < train) loadTrain.push(data);
61     else if (j < train + validation) loadValidation.push(data);
62     else if (j < train + validation + test) loadTest.push(data);
63   }
64 }
65
66 const datatrain = extractFiles(loadTrain, datatrain, 'Train');
67 const datavalidation = extractFiles(loadValidation, datavalidation, 'Validation');
68 const datatest = extractFiles(loadTest, datatest, 'Test');
69
70 // put 'labels', 'datatrain', 'datavalidation', 'datatest' to JSON and export it.
71 const dataset = { labels, datatrain, datavalidation, datatest }

```

Figure 5.27 load and split command on loadData Train application code.

`nextBatch` is used to get the next batch data and then convert it to train model input or output. It contains `getData` and `getLabel`. `getData` convert corresponding `dataset` data based on `iteration`, `batch`, `time`, `freq` to TensorFlow input. It defines TensorFlow input shape from `batch`, `time`, and `freq`. It defines TensorFlow input data convert it by rearrange the data inside `dataset` based on `batch`, `iteration`, `time`, and `freq`. Then it returns as 4d tensor. `getLabel` convert corresponding `dataset` label based on `iteration`, `batch` to TensorFlow output. It defines TensorFlow output shape from `batch`, and `labels.length`. It defines TensorFlow output data convert it by define array from `labels` based on `batch`, `iteration`, and `dataset`. Then it returns as 2d tensor.



```

32
33
34  /**
35   * 'getLabel' convert corresponding 'dataset' label based on 'iteration',
36   * 'batch' to tensorflow output.
37   *
38   * Define tensorflow output shape from 'batch', and 'labels.length'.
39   * Define tensorflow output data convert it by define array from 'labels' based on
40   * 'batch', 'iteration', and 'dataset'. Then it return as 2d tensor.
41   */
42  function getLabel(labels, dataset, iteration, batch) {
43    const data = [];
44    const shape = [batch, labels.length];
45
46    for (j = 0; j < batch; j++) {
47      let label;
48
49      if (iteration + j >= dataset.length) label = dataset[(iteration + j) % dataset.length].label;
50      else label = dataset[iteration + j].label;
51
52      for (k = 0; k < labels.length; k++) {
53        if (label === labels[k]) data.push(1);
54        else data.push(0);
55      }
56    }
57    return tf.tensor2d(data, shape);
58  }
59

```

Figure 5.28 `getData` on `nextBatch` Train application code.

```

32
33
34  /**
35   * `getLabel` convert corresponding `dataset` label based on `iteration`,
36   * `batch` to tensorflow output.
37   *
38   * Define tensorflow output shape from `batch`, and `labels.length`.
39   * Define tensorflow output data convert it by define array from `labels` based on
40   * `batch`, `iteration`, and `dataset`. Then it return as 2d tensor.
41   */
42  function getLabel(labels, dataset, iteration, batch) {
43    const data = [];
44    const shape = [batch, labels.length];
45
46    for (j = 0; j < batch; j++) {
47      let label;
48      if (iteration + j >= dataset.length) label = dataset[(iteration + j) % dataset.length].label;
49      else label = dataset[iteration + j].label;
50
51      for (k = 0; k < labels.length; k++) {
52        if (label === labels[k]) data.push(1);
53        else data.push(0);
54      }
55    }
56
57    return tf.tensor2d(data, shape);
58  }
59

```

Figure 5.29 getLabel on nextBatch Train application code.

saveModelDB save location to models collection in mimic\_speech database in MongoDB.

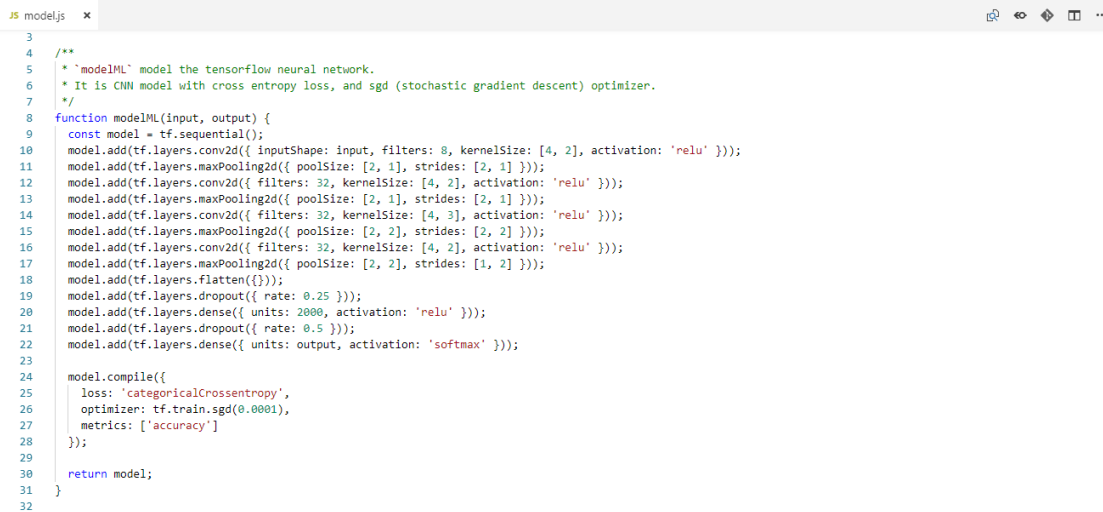
```

3
4
5  /**
6   * `saveModelDB` save `location` to `models` collection in mimic_speech database in MongoDB.
7   */
8  async function saveModelDB(location) {
9    try {
10      model.models.db = await connection.connect();
11      await model.models.create({ location });
12      await connection.disconnect();
13    } catch (err) {
14      error = err.errmsg;
15      console.log('Error:', error);
16    }
17  }

```

Figure 5.30 saveModelDB on Train application code.

model as modelML model the TensorFlow neural network. It is CNN model with cross entropy loss, and SGD (Stochastic Gradient Descent) optimizer.



```

3
4
5  /**
6   * `modelML` model the tensorflow neural network.
7   * It is CNN model with cross entropy loss, and sgd (stochastic gradient descent) optimizer.
8   */
9  function modelML(input, output) {
10   const model = tf.sequential();
11   model.add(tf.layers.conv2d({ inputShape: input, filters: 8, kernelSize: [4, 2], activation: 'relu' }));
12   model.add(tf.layers.maxPooling2d({ poolSize: [2, 1], strides: [2, 1] }));
13   model.add(tf.layers.conv2d({ filters: 32, kernelSize: [4, 2], activation: 'relu' }));
14   model.add(tf.layers.maxPooling2d({ poolSize: [2, 1], strides: [2, 1] }));
15   model.add(tf.layers.conv2d({ filters: 32, kernelSize: [4, 3], activation: 'relu' }));
16   model.add(tf.layers.maxPooling2d({ poolSize: [2, 2], strides: [2, 2] }));
17   model.add(tf.layers.conv2d({ filters: 32, kernelSize: [4, 2], activation: 'relu' }));
18   model.add(tf.layers.maxPooling2d({ poolSize: [2, 2], strides: [1, 2] }));
19   model.add(tf.layers.flatten({}));
20   model.add(tf.layers.dropout({ rate: 0.25 }));
21   model.add(tf.layers.dense({ units: 2000, activation: 'relu' }));
22   model.add(tf.layers.dropout({ rate: 0.5 }));
23   model.add(tf.layers.dense({ units: output, activation: 'softmax' }));
24
25   model.compile({
26     loss: 'categoricalCrossentropy',
27     optimizer: tf.train.sgd(0.0001),
28     metrics: ['accuracy']
29   });
30   return model;
31 }
32

```

Figure 5.31 modelML on model Train application code.

train contains preparation command for train and train itself. train train the model to fit the dataset.dataTrain and validate the model on dataset.dataValidation data by its loss and accuracy. The train iterate trainIteration times and test every testIteration in trainIteration. Every iteration loss and accuracy are printed. Every test is occurred and the last iteration, confusionMatrix is printed, model is saved, the saved model location is saved to MongoDB, loss and accuracy is saved as csv.

```

7
8 // load all data to `dataset`.
9 const dataset = require('./lib/loadData');
10
11 // define input and output tensorflow shape.
12 const time = dataset.datatrain[0].data.length;
13 const freq = dataset.datatrain[0].data[0].length;
14 const inputShape = [time, freq, 1];
15 const outputShape = dataset.labels.length;
16
17 const model = require('./model')(inputShape, outputShape); // load tensorflow model to `model`.
18 const batch = 32; // n inputs at the same time to be trained.
19 const epochs = 16; // n times the batches to be trained.
20 const trainIteration = 1500; // n times the train iteration.
21 const testIteration = 300; // Every n times in `trainIteration` the `model` is tested.
22
23 // define all test input and output from `dataset.datatest` to `testData` and `testLabel` for prediction.
24 const testData = nextBatch.getData(dataset.datatest, i, dataset.datatest.length, time, freq);
25 const testLabel = nextBatch.getLabel(dataset.labels, dataset.datatest, i, dataset.datatest.length);
26

```

Figure 5.32 preparation command on train Train application code.

```

29
30 /**
31  * `train` train the `model` to fit the `dataset.datatrain` and validate the `model` on
32  * `dataset.datavalidation` data by its loss and accuracy.
33  *
34  * Train iterate `trainIteration` times and test every `testIteration` in `trainIteration`.
35  * Every iteration `loss` and `accuracy` is printed. Every test is occurred and the last iteration,
36  * `confusionMatrix` is printed, `model` is saved, the saved model location is saved to MongoDB,
37  * `loss` and `accuracy` is saved as csv.
38  */
39 async function train() {
40   for (i = 0; i < trainIteration; i++) {
41     // define train input and output from `dataset.datatrain` and validation from `dataset.datavalidation`.
42     const batchData = nextBatch.getData(dataset.datatrain, i, batch, time, freq);
43     const batchLabel = nextBatch.getLabel(dataset.labels, dataset.datatrain, i, batch);
44     const batchValidation = nextBatch.getData(dataset.datavalidation, i, batch, time, freq);
45
46     // train the model and get `loss` and `accuracy`.
47     const history = await model.fit(
48       batchData,
49       batchLabel,
50       { batchSize: batch, validationData: batchValidation, epochs: epochs, verbose: 0 }
51     );
52     loss.push(history.history.loss[0].toFixed(10));
53     accuracy.push(history.history.acc[0].toFixed(5));
54
55     console.log(new Date().toISOString() + ' | Step ' + (i + 1) + ' | accuracy: ' + accuracy[i] + ' | loss: ' + loss[i]);
56
57     // check if it is n `testIteration`.
58     if ((i !== 0 && i % testIteration === 0) || (i + 1 === trainIteration)) {
59       // test the model and show `confusionMatrix` the test result.
60       const labels = tf.argMax(testLabel, 1);
61       const predictions = tf.argMax(model.predict(testData), 1);
62       const confusionMatrix = tf.math.confusionMatrix(labels, predictions, dataset.labels.length);
63

```

Figure 5.33 train on train Train application code.

## 5.2.5 Library Package Code

## **CHAPTER VI**

### **SYSTEM TESTING**

#### **6.1 Testing Environment**

The testing environment specify the environment during testing. The environment specification are as follows:

1. Windows 10.
2. Chrome browser.

There is additional environment only on the identify process. User that the records is trained by the application is called Tester 1. Male user that the record hasn't trained by the application called Tester 2. Female user that the record hasn't trained by the application called Tester 3. Every tester is test in each the following environment:

1. Noisy background (Loud music or people chit-chat).
2. Semi noisy background (AC noise or sound from the other rooms).
3. No noisy background.

#### **6.2 Testing Scenario**

The testing scenario conducted by evaluating all the features with a set of cases or scenario in the application based on its functionality requirement and defined testing environment.

The testing scenarios of the application is categorized based on Collect application, Train application, and Mimic application.

### 6.2.1 Collect Application

The Collect application has 3 subcategorized scenarios, Collect server, Home page and Phonemes page. The Collect server has scenarios to allow user to access Home page, Phonemes page in the browser and process the upload record request. The Home page has scenarios to allow user to direct to Phonemes page to start the Collect application and direct to Mimic application Home page to change to Mimic application. The Phonemes page has scenarios to allow user to upload the record to the server and direct to the next Phonemes page or direct back to Home page. Figures 6.1 shows some screenshot of expected result from the scenarios.

Table 6.1 Collect application scenarios.

No	Scenario	Expected Result	Result
1	Access Home page	Home page is shown	As expected
2	Access Phonemes page	Phonemes page is shown	As expected
3	Process upload record request	Printed process results information on the console	As expected
4	Direct to Phonemes page	Directed to Phonemes page	As expected
5	Direct to Mimic application Home page	Directed to Mimic application Home page	As expected
6	Upload the record	Audio element is shown with the user record	As expected
7	Direct to the next Phonemes page or direct back to Home page	Directed to the next Phonemes page or Directed back to Home Page	As expected

Figures 6.1 Some screenshot from the Collect application scenarios.

### 6.2.2 Train Application

The Train application has scenario to allows user to train collected data. When train collected data is run, it loads and splits collected data, extracts the data, gets the corresponding batch data, trains the TensorFlow model to fit the batch train data, tests the TensorFlow model from test data, and saves the TensorFlow model to database. Figures 6.2 shows some screenshot of expected result from the scenarios.

Table 6.2 Train application scenarios.

No	Scenario	Expected Result	Result
1	Load and split collected data	Printed information on the console	As expected
2	Extract the data	Printed information on the console	As expected
3	Get corresponding batch data	Printed information on the console	As expected
4	Train the TensorFlow model to fit batch train data	Printed information on the console	As expected
5	Test the TensorFlow model from test data	Printed information on the console	As expected
6	Save the TensorFlow model to database	Recorded document in mimic_speech database models collection	As expected

Figures 6.2 Some screenshot from the Train application scenarios.

### 6.2.3 Mimic Application

The Mimic application has 4 subcategorized scenarios, Mimic server, Home page, Identify page, and Generate page. The Mimic server has scenarios to allow user to access Home page, Identify page, and Generate page in the browser and process the upload record and identify speech on Identify request and load speech data and generate



speech on Generate request. The Home page has scenarios to allow user to direct to Identify page to start the identify speech application, direct to Generate page to start the generate speech application, and direct to Collect application Home page to change to Collect application. The Identify page has scenarios to allow user to upload the record to the server, send form request for identify speech to the server, and direct back to Home page. The Generate page has scenarios to allow user to load speech data from the server, send form request for generate speech to the server, and direct back to Home page. Figures 6.3 shows some screenshot of expected result from the scenarios.

Table 6.3 Mimic application scenarios.

No	Scenario	Expected Result	Result
1	Access Home page	Home page is shown	As expected
2	Access Identify page	Identify page is shown	As expected
3	Access Generate page	Generate page is shown	As expected
3	Process upload record on Identify request	Printed process results information on the console	As expected
4	Process identify speech on Identify request	Printed process results information on the console	As expected
5	Process load speech data on Generate request	Printed process results information on the console	As expected
6	Process generate speech on Generate request	Printed process results information on the console	As expected
7	Direct to Identify page	Directed to Identify page	As expected
8	Direct to Generate page	Directed to Generate page	As expected
9	Direct to Collect application Home page	Directed to Collect application Home page	As expected
10	Upload the record	Audio element is shown with the user record	As expected
11	Send form for identify speech	Alert information despite success or error in the process	As expected
12	Direct back to Home page from Identify page	Directed back to Home Page	As expected
13	Load speech data	Datalist element store the speech data or alert if no speech data is found	As expected

14	Send form for generate speech	Audio element is shown with the generated speech or alert if error is occurred	As expected
15	Direct back to Home page from Generate page	Directed back to Home Page	As expected

Figures 6.3 Some screenshot from the Mimic application scenarios.

The dataset during the speech recognition testing on identify process is 1000 male speech data on no noisy background. Corrected result shows from the more than 75% of model accuracy. The following tests table results below and figures that shows some screenshot of identify process results:

Table 6.4 Tester 1 speech recognition on identify process results.

No	Phoneme	Noisy Background	Semi Noisy Background	No Noisy Background	Result
1	a				/3
2	i				/3
3	t				/3
3	na				/3
4	ma				/3
5	mu				/3
6	di				/3
7	ri				/3
8	ku				/3
9	kan				/3
10	Unknown 1				/3
11	Unknown 2				/3
12	Unknown 3				/3

Figures 6.4 Some screenshot from the Tester 1 results.

Table 6.5 Tester 2 speech recognition on identify process results.

No	Phoneme	Noisy Background	Semi Noisy Background	No Noisy Background	Result
1	a				/3
2	i				/3
3	t				/3
3	na				/3
4	ma				/3
5	mu				/3
6	di				/3
7	ri				/3
8	ku				/3
9	kan				/3
10	Unknown 1				/3
11	Unknown 2				/3
12	Unknown 3				/3

Figures 6.5 Some screenshot from the Tester 2 results.

Table 6.6 Tester 3 speech recognition on identify process results.

No	Phoneme	Noisy Background	Semi Noisy Background	No Noisy Background	Result
1	a				/3
2	i				/3
3	t				/3
3	na				/3
4	ma				/3
5	mu				/3
6	di				/3
7	ri				/3
8	ku				/3
9	kan				/3
10	Unknown 1				/3
11	Unknown 2				/3
12	Unknown 3				/3

Figures 6.6 Some screenshot from the Tester 3 results.

## **CHAPTER VII**

### **CONCLUSIONS AND FUTURE WORK**

#### **7.1 Conclusion**

The following list sums up that the application is achieved based on this research objective:

1. This application enables to recognize speech in Bahasa Indonesia speech from record audio.
2. This application enables to generate speech in Bahasa Indonesia speech from text.
3. This application enables to mimic speech through the website.
4. This application enables to collect speech data through the website.
5. This application enables to train the collected speech data through the console.

#### **7.2 Future Work**

The following suggestion for further development and improvements of the research or application:

1. User Interface

Improvement on UI will always help the user experience. With some colourful theme, clear button or inputs, the application will comfortable to be used.

2. Speech Recognition

Improvement on machine learning model can be made. When there is no right or wrong in modelling the machine learning model, there is always optimal model to get the best accuracy. A research to find the optimal model or a research developing application with different machine learning model are big improvement in the application or even in Speech Recognition itself.

### 3. Speech Synthesis

Improvement on Speech Synthesis is when generating the speech. By removing some of silence or unused part of the speech and also reducing background will make generated speech more fluently and good to hear.

## REFERENCES

- [1] Hurwitz, J., & Kirsch, D. (2018). *Machine Learning For Dummies®*, IBM Limited Edition. New Jersey: John Wiley & Sons, Inc.
- [2] Nilsson, N. J., & Laboratory, R. (2005). *INTRODUCTION TO MACHINE LEARNING*. California: Nils J. Nilsson.
- [3] dictionary.cambridge.org. (2018). *MACHINE*. Retrieved from Cambridge English Dictionary:  
<https://dictionary.cambridge.org/dictionary/english/machine>
- [4] dictionary.cambridge.org. (2018). *LEARNING*. Retrieved from Cambridge English Dictionary:  
<https://dictionary.cambridge.org/dictionary/english/learning>
- [5] Society, T. R. (2017). *Machine learning: the power and promise*. London: The Royal Society.
- [6] edureka.com. (2018, October 18). *What is Machine Learning?* Retrieved from edureka: <https://www.edureka.co/blog/what-is-machine-learning/>
- [7] Geitgey, A. (2016, January 3). *Machine Learning is Fun! Part 2*. Retrieved from Medium: <https://medium.com/@ageitgey/machine-learning-is-fun-part-2-a26a10b68df3>
- [8] Magdi Zakaria, M. A.-S. (2014). Artificial Neural Network : A Brief Overview. *Int. Journal of Engineering Research and Applications*, 6.
- [9] A.D.Dongare, R. A. (2012). Introduction to Artificial Neural Network. *International Journal of Engineering and Innovative Technology*, 6.
- [10] courses.lumenlearning.com. (n.d.). *Introduction to Language*. Retrieved from lumen: <https://courses.lumenlearning.com/boundless-psychology/chapter/introduction-to-language/>
- [11] readingdoctor.com.au. (2016). *Phonemes, Graphemes and Letters: The Word Burger*. Retrieved from Reading Doctor:  
<http://www.readingdoctor.com.au/phonemes-graphemes-letters-word-burger/>
- [12] Yanti, N. T. (n.d.). FONEM BAHASA INDONESIA. *Academia.edu*, 10.
- [13] dictionary.cambridge.org. (2018). *VOWEL*. Retrieved from Cambridge English Dictionary: <https://dictionary.cambridge.org/dictionary/english/vowel>
- [14] puebi.readthedocs.io. (n.d.). *Huruf Vokal*. Retrieved from PUEBI Daring: <https://puebi.readthedocs.io/en/latest/huruf/huruf-vokal/>

- [15] dictionary.cambridge.org. (2018). *DIPHTHONG*. Retrieved from Cambridge English Dictionary:  
<https://dictionary.cambridge.org/dictionary/english/diphthong>
- [16] puebi.readthedocs.io. (n.d.). *Huruf Diftong*. Retrieved from PUEBI Daring:  
<https://puebi.readthedocs.io/en/latest/huruf/huruf-diftong/>
- [17] dictionary.cambridge.org. (2018). *CONSONANT*. Retrieved from Cambridge English Dictionary:  
<https://dictionary.cambridge.org/dictionary/english/consonant>
- [18] puebi.readthedocs.io. (n.d.). *Huruf Konsonan*. Retrieved from PUEBI Daring:  
<https://puebi.readthedocs.io/en/latest/huruf/huruf-konsonan/>
- [19] dictionary.cambridge.org. (2018). *CLUSTER*. Retrieved from Cambridge English Dictionary: <https://dictionary.cambridge.org/dictionary/english/cluster>
- [20] puebi.readthedocs.io. (n.d.). *Gabungan Huruf Konsonan*. Retrieved from PUEBI Daring: <https://puebi.readthedocs.io/en/latest/huruf/gabungan-huruf-konsonan/>
- [21] Dave, B., & Pipalia, P. D. (2014). SPEECH RECOGNITION: A REVIEW. *International Journal of Advance Engineering and Research*, 7.
- [22] M.A.Anusuya, & S.K.Katti. (2009). Speech Recognition by Machine: A Review. *International Journal of Computer Science and Information Security*, 25.
- [23] Geitgey, A. (2016, December 24). *Machine Learning is Fun Part 6*. Retrieved from Medium: <https://medium.com/@ageitgey/machine-learning-is-fun-part-6-how-to-do-speech-recognition-with-deep-learning-28293c162f7a>
- [24] en.wikipedia.org. (2018, October). *Nyquist–Shannon sampling theorem*. Retrieved from Wikipedia:  
[https://en.wikipedia.org/wiki/Nyquist%E2%80%93Shannon\\_sampling\\_theorem](https://en.wikipedia.org/wiki/Nyquist%E2%80%93Shannon_sampling_theorem)
- [25] Olshausen, B. A. (2000). Aliasing. *Redwood Center for Theoretical Neuroscience*, 6.
- [26] Rabiner, L. R., & Schafer, R. W. (1978). *Digital Processing of Speech Signals*. Prentice Hall: New Jersey.
- [27] Hande, S. S. (2014). A Review on Speech Synthesis an Artificial Voice Production. *International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering*, 8.
- [28] lyrebird.ai. (2018). *Our Voice Products*. Retrieved from Lyrebird:  
<https://lyrebird.ai/products>

- [29] [translate.google.com](https://translate.google.com/intl/en/about/languages/). (2018). *Languages*. Retrieved from Google Translate:  
<https://translate.google.com/intl/en/about/languages/>
- [30] [w3school.com](https://www.w3schools.com/nodejs/nodejs_intro.asp). (2018). *Node.js Introduction*. Retrieved from W3School:  
[https://www.w3schools.com/nodejs/nodejs\\_intro.asp](https://www.w3schools.com/nodejs/nodejs_intro.asp)
- [31] [mongodb.com](https://www.mongodb.com/nosql-explained). (2018). *NoSQL Databases Explained*. Retrieved from  
MongoDb: <https://www.mongodb.com/nosql-explained>