

**FEATURE EXTRACTION AND
CLASSIFICATION OF MALAY SPEECH
VOWELS**

SHAHRUL AZMI BIN MOHD YUSOF

UNIVERSITI MALAYSIA PERLIS

2010

© This item is protected by original copyright



**Feature Extraction and Classification of
Malay Speech Vowels**

by

Shahrul Azmi bin Mohd Yusof

(0640610049)

A Thesis submitted

In fulfillment of the requirements for the degree of

Doctor of Philosophy (Mechatronic Engineering)

School of Mechatronic Engineering

UNIVERSITI MALAYSIA PERLIS

2010

ACKNOWLEDGEMENT

I would like to thank my main supervisor Prof. Dr. Sazali Yaacob for his guidance and comments during my work on the thesis. He has always been so patient with me, continuously giving me encouragement and always confident that I will complete this journey successfully. I have always respected you as my supervisor, my boss, my mentor and my friend. Next, I would like to thank my next supervisor, Assoc. Prof. Dr. Paulraj Murugesu Pandiyan for all the knowledge he has shared with me. He has always treated me more like a colleague than his student. For that, I sincerely thank you from the bottom of my heart.

I wish to thank my employer, Universiti Utara Malaysia, for funding me to further my studies. I also wish to express gratitude to my colleagues from Universiti Utara Malaysia, namely Assoc. Prof. Fadzilah Siraj, Assoc. Prof. Dr. Lim Kong Teong, Dr. Nur Idayu Mahat, Dr. Azizan Saaban for sharing their valuable knowledge in their expert fields of Artificial Intelligence, Statistics, Mathematics and Pattern Recognition. I will always cherish the wonderful discussions we had and that I am really grateful you had find the time to discuss with me on my work. I would like to also thank distinguished individuals like Prof. Dr. Nagarajan Ramachandran and Prof. Dr. Farid Ghani from UniMAP for giving me useful tips on how to improve my work.

Work is only one part of life. The other part is personal. It is impossible to work with a light heart if the personal life is disturbed. Therefore, I would like to thank all of my family members and friends for being so supportive during this long journey. I would also like to thank my parents and my family in giving me the support and confidence in

taking this journey. I have also been blessed with a number of friends with whom I have enjoyed their company in making my PhD trip less stressful and more bearable. I thank them for making my life as wonderful as it is.

In the end, I would like to thank my beautiful wife, Normiyah Jusoh, for walking beside me in this journey, and for making every step in the path of life worth taking. Her spirit, devotion and sacrifices were the fuel which kept me going. I would also like to thank my six beautiful children, Muhammad Faris, Muhammad Hakimi, Nur Sarah Aisyah, Aziz Imran, Sufiya Najihah and Shaza Fatini for filling my life with happiness and joy. I could have never achieved this without all your love, support and encouragement.

© This item is protected by original copyright

TABLE OF CONTENTS

DECLARATION	I
ACKNOWLEDGEMENT	II
TABLE OF CONTENTS	IV
LIST OF TABLES.....	VII
LIST OF FIGURES.....	IX
LIST OF ABBREVIATIONS	XI
ABSTRAK.....	XII
ABSTRACT	XIII
1. INTRODUCTION	1
1.1 OVERVIEW.....	1
1.2 PROBLEM STATEMENT AND ITS SIGNIFICANCE	2
1.3 RESEARCH OBJECTIVES	5
1.4 MOTIVATION FOR THE PRESENT WORK	5
1.5 CONTRIBUTIONS OF THE THESIS	7
1.6 ORGANIZATION OF THE THESIS.....	8
2 LITERATURE REVIEW.....	10
2.1 INTRODUCTION	10
2.2 AUTOMATIC SPEECH RECOGNITION (ASR)	10
2.3 PREVIOUS RESEARCH WORKS IN ASR AND VOWEL RECOGNITION	12
2.4 SPEECH PRODUCTION AND ACOUSTICS.....	18
2.5 LINEAR PREDICTION OF SPEECH	19
2.6 MALAY PHONEMES.....	22
2.7 SINGLE FRAME VS. MULTI FRAME ANALYSIS.....	23
2.8 SPECTRAL ENVELOPE	25
2.9 FREQUENCY SCALE	27
2.9.1 <i>Bark Scale</i>	27
2.9.2 <i>Mel scale</i>	28
2.10 CLASSIFICATION TECHNIQUES.....	30
2.10.1 <i>Neural Network</i>	30
2.10.2 <i>Logistic Regression</i>	31
2.10.3 <i>Linear Discriminant Analysis</i>	32
2.10.4 <i>K-Nearest Neighbors</i>	33
2.11 CONCLUSIONS.....	33
3 METHODOLOGY	36
3.1 INTRODUCTION	36
3.2 VOWEL RECOGNITION PROCESS	36
3.3 DATA ACQUISITION	37
3.4 PREPROCESSING.....	38
3.4.1 <i>End Point Detection Technique</i>	39
3.4.2 <i>Pre-Emphasizing</i>	42
3.4.3 <i>Windowing</i>	42

3.4.4	<i>Spectral Envelope</i>	43
3.5	DETERMINING THE FRAME SIZE AND DURATION	46
3.5.1	<i>Frame Shifted</i>	46
3.5.2	<i>Frame Expanding Analysis</i>	47
3.5.3	<i>Determining Vowel Frequency Range</i>	49
3.6	FEATURE EXTRACTION METHODS	49
3.6.1	<i>First Formant Bandwidth (F1BW)</i>	50
3.6.2	<i>Spectral Delta (SpD)</i>	54
3.6.3	<i>Bark Intensity (BrKI)</i>	55
3.6.4	<i>Fixed Frequency Band (FFB)</i>	56
3.6.5	<i>Formant Frequency Difference (FFD)</i>	57
3.7	CLASSIFICATION METHODOLOGY	58
3.7.1	<i>Levenberg-Marquardt (LM) Network</i>	58
3.7.2	<i>Multinomial Logistic Regression (MLR)</i>	58
3.7.3	<i>Linear Discriminant Analysis (LDA)</i>	59
3.7.4	<i>K-Nearest Neighbors (KNN)</i>	59
3.8	ROBUSTNESS ANALYSIS	60
3.8.1	<i>Training and Testing Procedure</i>	62
3.9	STATISTICAL ANALYSIS.....	65
3.10	CONCLUSIONS.....	67
4	FEATURE EXTRACTION FOR VOWEL RECOGNITION.....	69
4.1	INTRODUCTION	69
4.2	PROPOSED FEATURE EXTRACTION METHODS.....	69
4.2.1	<i>First Formant Bandwidth (F1BW)</i>	70
4.2.2	<i>Spectral Delta (SpD)</i>	72
4.2.3	<i>Fixed Frequency Band (FFB)</i>	73
4.2.4	<i>Formant Frequency Difference (FFD)</i>	77
4.2.5	<i>Bark Intensity (BrKI)</i>	78
4.3	MEL-FREQUENCY CEPSTRUM COEFFICIENTS (MFCC).....	79
4.4	CONCLUSIONS.....	79
5	VOWEL CLASSIFICATION	81
5.1	INTRODUCTION	81
5.2	CLASSIFICATION RESULTS	81
5.2.1	<i>First Formant Bandwidth (F1BW)</i>	82
5.2.2	<i>Spectral Delta (SpD)</i>	84
5.2.3	<i>Fixed Frequency Band (FFB)</i>	86
5.2.4	<i>Formant Frequency Difference (FFD)</i>	88
5.2.5	<i>Bark Intensity (BrKI)</i>	89
5.2.6	<i>Mel-frequency cepstrum coefficients (MFCC)</i>	91
5.3	CONCLUSIONS.....	94
6	ROBUSTNESS ANALYSIS.....	98
6.1	INTRODUCTION	98
6.2	ROBUSTNESS ANALYSIS	98
6.2.1	<i>First Formant Bandwidth (F1BW)</i>	99
6.2.2	<i>Spectral Delta (SpD)</i>	104
6.2.3	<i>Fixed Frequency Band (FFB)</i>	108
6.2.4	<i>Bark Intensity (BrKI)</i>	112

6.2.5	Single Frame Mel-Frequency Cepstrum Coefficients (MFCC)	117
6.2.6	Comparison Analysis	122
6.3	CONCLUSIONS.....	126
7	CONCLUSION	128
7.1	INTRODUCTION	128
7.2	CONCLUSIONS.....	128
7.3	SPEECH APPLICATION.....	133
7.4	FUTURE DIRECTIONS	134
	REFERENCES	135
	APPENDIX.....	141
A.	PAPERS RELATED TO VOWEL RECOGNITION.....	141
B.	LIST OF PUBLICATIONS	148

© This item is protected by original copyright

LIST OF TABLES

Table 2.1 Summary of Vowel Recognition Papers	14
Table 2.2 Recent Related Literature on Vowel Recognition.....	17
Table 2.3 Vowel system of standard Malay	23
Table 2.4 Bark Frequency Scale.....	28
Table 3.1 Data Collection Details.....	37
Table 3.2 Values of first and second formant of the recorded vowels	49
Table 3.3 BrKI Frequency Ranges	55
Table 3.4 Number of features by frequency subband of FFB	56
Table 3.5 Different Sets of Formants	57
Table 4.1 Spectrum Details and Frequency Ranges for Each Vowel.....	70
Table 4.2 ANOVA Analysis of F1BW Features	71
Table 4.3 Individual vowel classification performance from different SpD number of features.....	72
Table 4.4 ANOVA Analysis of SpD Features.....	73
Table 4.5 Number of features by frequency subband of FFB	74
Table 4.6 MLR Vowel Classification Rates for FFB	74
Table 4.7 KNN Vowel Classification Results for FFB	75
Table 4.8 ANOVA Analysis of FFB Features.....	76
Table 4.9 Spectrum Mean and Standard Deviation of F1 and F2 for different vowels.....	77
Table 4.10 ANOVA Analysis of FFD Features	78
Table 4.11 ANOVA Analysis of BrKI Features.....	78
Table 4.12 ANOVA Analysis of MFCC Features.....	79
Table 5.1 Classification rate of F1BW features using different classifiers.....	83
Table 5.2 Best and worst vowel classification result for F1BW features.....	84
Table 5.3 Classification rate of SpD features using different classifiers.....	85
Table 5.4 Best and worst vowel classification result for SpD features	86
Table 5.5 FFB classification rate using multiple classifiers (Tabulated Results)	86
Table 5.6 Best and Worst Vowel Classification Result for FFB features	88
Table 5.7 Formant Feature Combinations	88
Table 5.8 Formant classification by individual vowels using formant combinations	89
Table 5.9 BrKI Classification Rate using Multiple Classifiers (Tabulated Results).....	89
Table 5.10 Best and Worst Vowel Classification Result for BrKI features.....	91
Table 5.11 MFCCs Classification Rate using Multiple Classifiers (Tabulated Results).....	92
Table 5.12 Best and Worst Vowel Classification Result for MFCCs features.....	93
Table 5.13 MFCCf Classification Rate using Multiple Classifiers	94
Table 5.14 Best and Worst Vowel Classification Result for MFCCf features.....	94
Table 5.15 Overall Vowel Classification Performance of Features by Classifiers	97
Table 6.1 Comparison of Overall F1BW Classification Rate by Different SNR level (Tabulated Result)	102
Table 6.2 Overall vowel classification Rate of Vowels on F1BW features using Clean Training Data (Tabulated Results).....	104
Table 6.3 Comparison of Overall SpD Classification Rate by Different SNR level (Tabulated Result)	107
Table 6.4 Overall vowel classification Rate of Vowels on SpD features using Clean Training Data (Tabulated Results).....	108
Table 6.5 Comparison of Overall FFB Classification Rate by Different SNR level (Tabulated Result)	110

Table 6.6 Overall vowel classification Rate of Vowels on FFB features using Clean Training Data (Tabulated Results).....	112
Table 6.7 Comparison of Overall BrKI Classification Rate by Different SNR level (Tabulated Result)	115
Table 6.8 Overall vowel classification Rate of Vowels on BrKI features using Clean Training Data (Tabulated Results).....	117
Table 6.9 Comparison of Overall MFCC Classification Rate by Different SNR level (Tabulated Result)	120
Table 6.10 Overall vowel classification Rate of Vowels on MFCC features using Clean Training Data (Tabulated Results).....	122
Table A.1 Papers Related to Vowel Recognition	141

© This item is protected by original copyright

LIST OF FIGURES

Figure 2.1 ASR System	11
Figure 2.2 Tongue Position	22
Figure 2.3 An example of spectrum and spectral envelope of Vowel /ə/ taken using a sampling frequency of 8 KHz.	25
Figure 2.4 Linear-Scaled Mean Spectrum Envelope of Vowels (0-4000Hz) after Low Pass Filtering	26
Figure 2.5 Log-Scaled Mean Spectrum Envelope of Vowels (0-4000Hz) after Low Pass Filtering	26
Figure 2.6 Critical band filters used in MFCC computation and their outputs	29
Figure 3.1 Vowel Recognition Process	36
Figure 3.2 Data Acquisition Process	38
Figure 3.3 Vowel Extraction Process	40
Figure 3.4 Example of clipped waveform	41
Figure 3.5 Spectrum Envelope of Vowels (Linear y-axis)	45
Figure 3.6 Mean Spectrum Envelope of Vowels (Original).....	45
Figure 3.7 Frame Shifted Waveform from 0-100%	46
Figure 3.8 Frame Shifted Spectrum from 0-100%	47
Figure 3.9 Frame Expanding Waveform from 0-100%.....	48
Figure 3.10 Frame Expanding Spectrum from 0-100%	48
Figure 3.11 Bandwidth Example.....	51
Figure 3.12 F1BW Methodology.....	52
Figure 3.13 Determining Extraction parameters for Vowel /a/	53
Figure 3.14 Spectrum Envelope of Vowels at Different SNR.....	61
Figure 3.15 Robustness Analysis Methodology (Training together with noisy Data) ...	63
Figure 3.16 Robustness Analysis Methodology (Training with only raw data).....	64
Figure 4.1 Mean Spectrum Envelope of Standard Malay Vowels (0-2.5KHz).....	70
Figure 4.2 Overall CR% using MLR Classifier	75
Figure 4.3 Overall CR% using KNN Classifier	76
Figure 5.1 Result of F1BW Classification Rate using Multiple Classifiers	83
Figure 5.2 Spectral Delta classification based on number of features.....	84
Figure 5.3 Result of SpD Classification Rate using Multiple Classifiers	85
Figure 5.4 Result of FFB Classification Rate using Multiple Classifiers	87
Figure 5.5 Result of BrKI Classification Rate using Multiple Classifiers	90
Figure 5.6 Result of MFCC Classification Rate using Multiple Classifiers	92
Figure 5.7 Result of MFCCf Classification Rate using Multiple Classifiers	93
Figure 5.8 Overall Vowel Classification Performance of Features by Classifiers Performance.....	96
Figure 6.1 Overall F1BW Classification Rate by Different SNR level (MLR)	99
Figure 6.2 Overall F1BW Classification Rate by Different SNR level (KNN)	100
Figure 6.3 Overall F1BW Classification Rate by Different SNR level (LDA).....	101
Figure 6.4 Comparison of Overall F1BW Classification Rate by Different SNR level	102
Figure 6.5 Overall F1BW Classification Rate of Vowels based on Classifiers and Training Conditions using Clean Training Data.....	103
Figure 6.6 Overall SpD Classification Rate by Different SNR level (MLR).....	105
Figure 6.7 Overall SpD Classification rate by different SNR level (KNN)	105

Figure 6.8 Overall SpD Classification Rate by Different SNR level (LDA)	105
Figure 6.9 Comparison of Overall SpD Classification Rate by Different SNR level ..	106
Figure 6.10 Overall SpD Classification Rate of Vowels based on Classifiers and Training Conditions using Clean Training Data.....	107
Figure 6.11 Overall FFB Classification Rate by Different SNR level (MLR).....	109
Figure 6.12 Overall FFB Classification Rate by Different SNR level (KNN).....	109
Figure 6.13 Overall FFB Classification Rate by Different SNR level (LDA)	109
Figure 6.14 Comparison of Overall FFB Classification Rate by Different SNR level	110
Figure 6.15 Overall FFB Classification Rate of Vowels based on Classifiers and Training Conditions.....	111
Figure 6.16 Overall BrKI Classification Rate by Different SNR level (MLR).....	113
Figure 6.17 Overall BrKI Classification Rate by Different SNR level (KNN).....	113
Figure 6.18 Overall BrKI Classification Rate by Different SNR level (LDA)	114
Figure 6.19 Comparison of Overall BrKI Classification Rate of Vowels by Different SNR level.....	115
Figure 6.20 Overall BrKI Classification Rate of Vowels based on Classifiers and Training Conditions using Clean Training Data.....	116
Figure 6.21 Overall MFCC Classification Rate by Different SNR level (MLR).....	117
Figure 6.22 Overall MFCC Classification Rate by Different SNR level (KNN).....	118
Figure 6.23 Overall MFCC Classification Rate by Different SNR level (LDA)	119
Figure 6.24 Comparison of Overall MFCC Classification Rate by Different SNR level	120
Figure 6.25 Overall vowel classification Rate of Vowels using MFCC features.....	121
Figure 6.26 KNN Model Trained with Noisy Data	122
Figure 6.27 LDA Model Trained with Noisy Data	123
Figure 6.28 MLR Model Trained with Noisy Data	123
Figure 6.29 KNN model trained with raw data	124
Figure 6.30 LDA model trained with raw data.....	125
Figure 6.31 MLR model trained with raw data	125
Figure 7.1 GUI for Malays Pronunciation Test Application.....	133

LIST OF ABBREVIATIONS

ANOVA	Analysis of Variance
ASR	Automatic Speech Recognition
BRKI	Bark Intensity
df	Degree of Freedom
F1BW	First Formant Bandwidth
FFB	Fixed Frequency Band
FFD	Formant Frequency Difference
FFT	Fast Fourier Transform
KNN	K-Nearest Neighbours
LDA	Linear Discriminant Analysis
LM	Levenberg-Marquardt (LM) Network
LPC	Linear Predictive Coding
MFCCs	Single Framed Mel-frequency cepstrum coefficients
MFCCf	Multi-Framed Mel-frequency cepstrum coefficients
MLR	Multinomial Logistic Regression
SNR	Signal-to-Noise Ratio
SPSS	A computer program used for statistical analysis
SpD	Spectral Delta
ZCR	Zero-Crossing Rate

ABSTRAK

EKSTRAKSI CIRI DAN KLASIFIKASI VOWAL SUARA MELAYU

Dalam bahasa manusia, fonem ialah unit struktural terkecil yang membezakan makna. Biasanya, bahasa seperti bahasa Inggeris umumnya menggabungkan fonem untuk membentuk sesuatu perkataan. Dalam banyak bahasa, unit konsonan-vowal (CV) mempunyai frekuensi kejadian yang tertinggi di antara pelbagai bentuk unit sub-perkataan. Oleh kerana itu, pengecaman unit CV dengan ketepatan yang baik adalah sangat penting untuk pembangunan sistem pengenalan suara. Ada juga banyak aplikasi yang berdasarkan kepada fonem vokal. Diantaranya ialah sistem terapi bicara yang dapat meningkatkan sebutan perkataan terutama kepada anak-anak. Ada juga sistem yang mengajar pesakit cacat pendengaran untuk bercakap dengan sebutan betul dengan mengucapkan kata-kata pada tahap kefahaman yang tinggi. Semua sistem ini memerlukan kebolehan pengenalan vokal yang sememangnya menjadi fokus di dalam tesis ini. Tesis ini menyumbangkan empat kaedah ekstraksi ciri yang diperbaiki untuk pengecaman vokal berdasarkan intensiti jalur frekuensi saringan. Cadangan-cadangan yang baru itu adalah Jalur Lebar Forman Pertama (F1BW), Jalur Frekuensi Forman Tetap (FFB), Intensiti Bark (BrKI), Delta Spektral (SpD) dan Perbezaan Frekuensi Forman (FFD). Kemampuan keempat-empat kaedah ini akan dibandingkan dengan tiga kaedah ekstraksi ciri konvensional iaitu Koefisien Cepstral Frekuensi Mel satu frem (MFCCs), Koefisien Cepstral Frekuensi Mel banyak frem (MFCCf) dan 3-Forman Pertama. Klasifikasi-klasifikasi yang dianalisa adalah Regresi Logistik Multinomial (MLR), Rangkaian Levenberg-Marquardt (LM), Jiran Terhampir-k (KNN) dan Analisa Diskriminan Linear (LDA). Ada empat sumbangan utama dari tesis ini. Pertama adalah korpus vokal baru yang terdiri daripada lebih dari 1300 huruf vokal dirakam dari 100 individu Malaysia. Kedua adalah lima kaedah ekstraksi ciri yang telah menunjukkan prestasi lebih baik berbanding MFCC jika dianalisa dengan frame tunggal. Ketiga adalah prestasi dan analisa kerobusan menggunakan klasifikasi yang berbeza dan tahap kebisingan Gaussian yang berbeza. Sumbangan keempat adalah kriteria untuk melakukan analisis vokal terpencil.

ABSTRACT

FEATURE EXTRACTION AND CLASSIFICATION OF MALAY SPEECH VOWELS

In human language, a phoneme is the smallest structural unit that distinguishes meaning. Normally, language like English commonly combines phonemes to form a word. In many languages, the Consonant-Vowel (CV) units have the highest frequency of occurrence among different forms of subword units. Therefore, recognition of CV units with a good accuracy is crucial for development of a speech recognition system. There are also many applications that use vowels phonemes. Among them are speech therapy systems that improve utterances of word pronunciation especially to children. There are also systems that teach hearing impaired person to speak properly by pronouncing words with a good degree of intelligibility. All of these systems require high degree of vowel recognition capability in which this study focuses on. This thesis contributes five modified feature extraction methods for vowel recognition based on intensities of the Frequency Filter Bands. They are First Formant Bandwidth (F1BW), Fixed Formant Frequency Band (FFB), Spectral Delta (SpD), Bark Intensity (BrKI) and Formant Frequency Difference (FFD). The performance of these five proposed methods are compared with performance of three conventional feature extraction methods of single frame Mel-frequency cepstrum coefficients (MFCCs), multiple frame Mel-frequency cepstrum coefficients (MFCCf) and the first three formant features. The classifiers analysed in this study were Multinomial Logistic Regression (MLR), Levenberg-Marquardt (LM) network, k-Nearest Neighbors (KNN) and Linear Discriminant Analysis (LDA). There are four main contributions of this thesis. First is the new vowel corpus consisting of more than 1300 recorded vowels from 100 Malaysian speakers. Second are the five improved feature extraction methods which perform better than MFCC on single frame analysis. The third is the performance and robustness analysis using different classifiers and different Gaussian noise level. The fourth contribution is the frame analysis criteria for isolated vowel analysis.

1. INTRODUCTION

1.1 Overview

Automatic speech recognition (ASR) has made great strides with the development of digital signal processing hardware and software especially using English as the language of choice. Despite of all these advances, machines cannot match the performance of their human counterparts in terms of accuracy and speed, especially in case of speaker independent speech recognition. Victor Zue in 2004 stated that within 5 to 10 years, systems that can handle more complicated human-to-computer interactions, like processing a request for movie tickets at a particular theatre via speech recognition should be in use (Hoffman, 2009). Today, significant portion of speech recognition research focuses on speaker independent speech recognition problem. The reasons are its wide range of applications, and limitations of available techniques of speech recognition.

Speech processing analyses and processes speech signals for information retrieval, giving commands and speaker recognition. Due to the heavy dependency on digital signals, speech processing can be placed in the area of digital signal processing and natural language processing. Speech processing covers a broad area that relates to the following important research directions like speaker recognition (Campbell Jr, 1997), Speech enhancement (Ephraim, 1992), Speech coding (Furui, 2001), Voice analysis (Jilek, Marienhagen, & Hacki, 2004), Speech synthesis (Furui, 2001) and Speech recognition (Rabiner & Juang, 1993).

Automatic Speech Recognition (ASR) is the automatic conversion of speech sound waves to text and this automatic interaction between man and machine. This has

been an increasingly interesting problem for many years due to its vast potential like cell phone voice-dialing, bank transaction with a computer, or software dictating letters to word. The language of choice is usually English but researchers from different countries are focusing on doing speech recognition using their native language.

Speech recognition system performs two fundamental operations: signal modeling and pattern matching (Picone, Inc, & Dallas, 1993). Signal modeling represents process of converting speech signal into a set of parameters. Pattern matching is the task of finding parameter set from memory which closely matches the parameter set obtained from the input speech signal. The signal modeling involves four basic operations: spectral shaping, feature extraction, parametric transformation, and statistical modeling (Picone *et al.*, 1993).

Spectral shaping is the process of converting the speech signal from sound pressure wave to a digital signal; and emphasizing important frequency components in the signal. Feature extraction is the process of obtaining different features such as power, pitch, and vocal tract configuration from the speech signal. Parameter transformation is the process of converting these features into signal parameters through process of differentiation and concatenation. Statistical modelling involves conversion of parameters in signal observation vectors.

1.2 Problem Statement and Its Significance

In human language, a phoneme is the smallest structural unit that distinguishes meaning. Normally, language like English commonly combines phonemes to form a word. In many languages, the Consonant-Vowel (CV) units have the highest frequency

of occurrence among different forms of subword units. Therefore, recognition of CV units with a good accuracy is crucial for development of a speech recognition system. Recognition of these subword units is a large class set pattern classification problem because of the large number (typically, a few thousands) of units (Sekhar, Takeda, & Itakura, 2002). In this case, if ASR recognizes the vowel with a good accuracy, system can reduce region of search and improve accuracy and time.

English uses a combination of phonemes to form words which may not exactly follow the characters of the words. Because of this, a large database of vocabulary is needed in order to represent each individual word. Standard Malay (SM) on the other hand can be uttered properly based on the combination of CV phonemes. One advantage the Bahasa Malaysia has over English is the numbers of vowel phoneme that need to be considered. The proper Bahasa Malaysia has only 6 vowels phonemes which are /a/, /e/, /i/, /o/, /u/ and /ə/ (Maris, 1966) whereas typical American English has 20 vowel phonemes (Power, 2009).

There are many speech recognition applications that can help physically handicapped or paralyzed individuals in their daily life. Activities such as opening doors or pressing a button may prove very difficult for these individuals. Voice command system may help to lessen their burden. There are limitless applications for voice command systems such as opening doors, ordering and purchasing items or even turning on electrical appliances. Current new technologies are enabling robots to assist human via voice command. Most of these applications use English speech recognition capability. In Malaysia, there are still many individuals who are unable to converse well in English. Their pronunciation in English may not be suitable for these speech

recognition systems which normally works well with American or British English spoken individuals. Systems that work well with Malaysians spoken English are still limited and there are even lesser systems that utilize Malay language.

In Malay language, children are taught to spell the words using a combination of consonants (C) and vowels (V) sounds. A computer system that can read CV combinations from a person who can properly pronounce any SM words have the capability to function without a proper database. There are also many applications that use vowels phonemes. All of these systems require high degree of SM vowel recognition capability.

Although there are studies concerning Malay phoneme recognition, it is still at its infancy (Rosdi & Ainon, 2008) and multiple frame analysis is mostly in use by Malaysian researcher. Accuracy and processing time is a concern when developing speech therapy systems. More efforts are needed to be taken in order to develop Malay speech recognition system and this study is an effort to improve Malay vowel recognition.

In Malaysia, researches in vowel recognition is still lacking especially in the usage of Malay vowels, independent speaker systems, recognition robustness and algorithm speed and accuracy. There is a need to develop a better algorithm of Malay vowel recognition in terms of accuracy and robustness. This thesis will address the issues of researching on Malay word and vowel databases, independent speaker systems, and robustness analysis and algorithm improvement.

1.3 Research Objectives

This study is an effort to increase Malay vowel recognition capability by using a new speech database that consist of words uttered by Malaysian speakers from the three major races, Malay, Chinese and Indians. Robust feature extraction methods need to be developed. The main objective of this study is to increase independent speaker Malay vowel recognition capability in terms of accuracy and robustness. In order to achieve this, three sub-objectives are listed below.

- i). To develop an improved feature extracting algorithms for Malay vowels using independent speaker database.
- ii). To study linear and non linear classifiers in classifying vowels.
- iii). To study robustness of feature extraction methods using different classifiers under different Gaussian noise level.

1.4 Motivation for the Present Work

English word pronunciation depends on a sequence of phonemes. Audio signals are broken up into acoustic components and translated into phonemes. These phoneme sequences are then compared with words from an English database that can make up of thousands of words. For Malay words, the approach is different. It is comprised of Consonant-Vowel (CV) and Consonant-Vowel-Consonant (CVC) combinations. It is possible that a Malay word can be spelled out by a computer similar to a human being. We believe that a computer can be taught to spell like a child and having a computer system able to translate and CV or CVC combinations into proper and understandable words.

Among the motivations for this work is to have voice command and speech recognition application which uses Malay words instead of the commonly used English words. This may allow Malaysian citizen who are not well versed in English language to be able to use this system to assist them in their daily life like purchasing items using an automated speech recognition telemarketing system. This system may also help children to improve their pronunciation of Malay words thus promoting our National Language of Bahasa Malaysia. A speech therapy system may improve utterances of word pronunciation especially to children. Hearing impaired person can learn to speak properly by pronouncing words with a good degree of intelligibility through the use of a visual therapy system and evaluate the pronunciation capability of the speaker and display the results. Even language learners from foreign countries may find it easier to learn Bahasa Malaysia with the assistance of an application that teaches Bahasa Malaysia by correcting the pronunciation of the speaker.

Another motivation for this work is to have a computer system to recognise Malay Language based on CV or CVC words. This capability may allow the speech recognition system to recognize any Malay word spoken through the sequence of consonant and vowels, independent of speaker's race, gender and age. All of this type of system requires a high degree of Standard Malay vowel recognition capability. A good application of this capability is in healthcare where foreign doctors are able to interact with a local Malaysian without the use of a translator. In tourism industry, this system may allow a foreigner to converse with a local through the use of a mobile electronic translator in order to find directions to place of interest. For example, a German tourist may be able to interact with a villager from Malaysia concerning some food delicacies from that community with an application that can recognise different

dialects. In other words, the language or even dialect barrier between different people from different culture and ethnic group can be brought down with the use of a good, reliable Malaysian speech recognition system. This will bring in more tourists and enrich the locals especially from villages and promote the tourism industry in Malaysia.

1.5 Contributions of the Thesis

Among the contributions of this thesis are the five improved feature extraction methods for vowel recognition based on intensities of the frequency filter bands and single frame analysis. They are First Formant Bandwidth (F1BW), Fixed Formant Frequency Band (FFB), Spectral Delta (SpD), Bark Intensity (BrKI) and Formant Frequency Difference (FFD). These features are analysed using four different classifiers of Linear Discriminant Analysis (LDA), k-Nearest Neighbours (KNN), Multinomial Logistic Regression (MLR) and Levenberg-Marquardt (LM). Robustness analysis of these features and classifiers are also contributions of this thesis in providing robustness capability of different extraction methods compared to conventional methods of Mel-frequency cepstrum coefficients (MFCC) and the first three formant frequencies. This robustness analysis also includes the capability of the classifiers under different SNR condition.

To summarize the contributions of this thesis, four main contributions are listed below.

They are:

- i). New vowel corpus consists of more than 1300 recorded vowels from 100 Malaysian speakers.

- ii). Five improved feature extraction methods which perform better than MFCC on single frame analysis.
- iii). Performance and robustness analysis of Malay Vowel Classification under different Gaussian noise level.
- iv). Criteria for isolated Malay Vowel analysis.

1.6 Organization of the Thesis

This work has been divided into seven chapters:

Chapter 1 introduces the framework in which the thesis has been developed. First, overview of speech processing is presented, followed by an explanation of the Problem Statement and Significance, Research Objectives, Motivation for the present work and Contributions of the thesis.

Chapter 2 presents the literature review part of the thesis on the background studies and researches done by previous researchers in the field of speech recognition. First, an explanation of speech production and acoustics is given followed by Linear Prediction of Speech. Malay Phonemes is introduced together with a brief explanation about Standard Malay (SM). Next, an explanation of Spectral Envelope is presented followed by explanation of Frequency Scales of Bark and Mel. An explanation of current feature extraction methods based on Malay phonemes from studies done mostly by academicians from Malaysian universities is presented. A brief explanation of classification techniques is also presented.

Chapter 3 contains the methodology used in the research including the experimental setup, the feature extraction methods and experimental work. It presents the database

used in the experiments, along with an endpoint detector designed for cleaning the samples.

Chapter 4 explains the 4 new feature extraction methods of First Formant Bandwidth, Spectral Delta, Formant Frequency Difference, Bark Intensity and Fixed Frequency Band. Their performance will be compared with the first 3 formant features and Mel-frequency cepstrum coefficient features. Feature validation is also presented using ANOVA method to show the mean significance of each vowel against the features.

Chapter 5 discusses the classification results obtained in the experiments. Then, a detailed comparison between the performance of the proposed approach and that of conventional approach is offered. These results will be discussed and concluded at the end of this chapter.

Chapter 6 presents the study of robustness of the feature extraction methods and three classifiers. In terms of robustness analysis, the performance of different methods and different classifiers will be presented.

Chapter 7 makes a summary and discussion of the results from the study. An application utilising of the vowel recognizer will also be discussed and presented. Finally, the most important conclusions are extracted, and future work directions are suggested.

2 LITERATURE REVIEW

2.1 Introduction

Automatic Speech Recognition (ASR) belongs to the class of digital speech processing technologies that also includes speech synthesis (text-to-speech, language generation) and voice biometrics (speaker identification, speaker verification). In general, their aim is to allow a machine to replicate human ability to hear, identify, and utter natural human spoken language. For the past 35 years, robust ASR systems have made great progress over the years. Several factors have contributed to this rapid progress, such as the development of advanced signal processing techniques but also the continuously increasing computing power.

2.2 Automatic Speech Recognition (ASR)

The earliest attempts to devise ASR systems were made in 1950s and 1960s, when various researchers tried to exploit fundamental ideas of acoustic phonetics. Since signal processing and computer technologies were yet very primitive, most of the speech recognition systems investigated used spectral resonances during the vowel region of each utterance which were extracted from output signals of an analogue filter bank and logic circuits.

The goal of an Automatic Speech Recognition (ASR) system is to transcribe speech to text. As illustrated in Figure 2.1, the speaker's mind decides the source word sequence W that is delivered through his/her text generator. The source is passed through a noisy communication channel that consists of the speaker's vocal apparatus