

SELF-ORGANIZING MAP AND MULTILAYER PERCEPTRON  
FOR MALAY SPEECH RECOGNITION

GOH KIA ENG

A thesis submitted in fulfilment of the  
requirements for the award of the degree of  
Master of Science (Computer Science)

Faculty of Computer Science and Information System  
Universiti Teknologi Malaysia

AUGUST 2006

To my beloved mother and father

## ACKNOWLEDGEMENT

First of all, I would like to thank my mother and father who have been supporting me and giving me lots of encouragements to complete this thesis. They've been so great and I know there would be no way I could have such a wonderful life without having the love and care from them. Thanks for always been there for me.

A special thank to my supervisor, Prof. Madya Abdul Manan bin Ahmad, for all his guidance and time. Thanks so much for his advices, comments and suggestions on how to improve this research and how to produce a good thesis. He is an understanding and helpful person in helping me to complete this research.

Not forgetting I also would like to take this opportunity to thank all my friends. All the motivations, helps, and supports are fully appreciated. Thanks for being there and listening to my complaints and lend me a helpful hand when I am in troubles.

Last but not least for those who were not mentioned above, I would like you to know that your countless effort and support will always remembered. All credits to everyone! Thank you very much.

## ABSTRACT

Various studies have been done in this field of speech recognition using various techniques such as Dynamic Time Warping (DTW), Hidden Markov Model (HMM) and Artificial Neural Network (ANN) in order to obtain the best and suitable model for speech recognition system. Every model has its drawbacks and weaknesses. Multilayer Perceptron (MLP) is a popular ANN for pattern recognition especially in speech recognition because of its non-linearity, ability to learn, robustness and ability to generalize. However, MLP has difficulties when dealing with temporal information as it needs input pattern of fixed length. With that in mind, this research focuses on finding a hybrid model/approach which combines Self-Organizing Map (SOM) and Multilayer Perceptron (MLP) to overcome as well as reduce the drawbacks. A hybrid-based neural network model has been developed to speech recognition in Malay language. In the proposed model, a 2D SOM is used as a sequential mapping function in order to transform the acoustic vector sequences of speech signal into binary matrix which performs dimensionality reduction. The idea of the approach is accumulating the winner nodes of an utterance into a binary matrix where the winner node is scaled as value “1” and others as value “0”. As a result, a binary matrix is formed which represents the content of an utterance. Then, MLP is used to classify the binary matrix to which each word corresponds to. The conventional model (MLP only) and the proposed model (SOM and MLP) were tested for digit recognition (“*satu*” to “*sembilan*”) and word recognition (30 selected Malay words) to find out the recognition accuracy using different values of parameters (cepstral order, dimension of SOM, hidden node number and learning rate). Both of the models were also tested using two types of classification: syllable classification and word classification. Finally, comparison and discussion was made between conventional and proposed model based on their recognition accuracy. The experimental results showed that the proposed model achieved higher accuracy.

## ABSTRAK

Banyak penyelidikan telah dijalankan dalam bidang pengecaman suara menggunakan pelbagai teknik seperti *Dynamic Time Warping* (DTW), *Hidden Markov Models* (HMM), *Artificial Neural Network* (ANN) dan sebagainya. Namun demikian, setiap teknik mempunyai kelemahannya masing-masing. Hal ini menyebabkan sistem menjadi kurang tepat. *Multilayer Perceptron* (MLP) merupakan satu rangkaian neural yang terkenal bagi pengecaman suara. Walau bagaimanapun, MLP mempunyai kelemahan di mana akan melemahkan prestasi sistem. Oleh itu, penyelidikan ini menumpu terhadap pembangunan satu model *hybrid* yang menggabungkan dua rangkaian neural iaitu *Self-Organizing Map* (SOM) dan *Multilayer Perceptron* (MLP). Satu model berasaskan rangkaian neural *hybrid* telah dibangunkan bagi sistem pengecaman suara dalam bahasa Melayu. Dalam model ini, SOM yang berdimensi dua digunakan sebagai fungsi pemetaan turutan untuk menukar turutan vector akustik bagi isyarat suara kepada matrik binari. Hal ini bertujuan untuk mengurangkan dimensi bagi vektor suara. SOM menyimpan nod pemenang bagi suara dalam bentuk matrik di mana nod pemenang diskalakan kepada nilai "1" dan yang lain diskalakan kepada nilai "0". Hal ini membentuk satu matrik binari yang mewakili kandungan suara tersebut. Kemudian, MLP mengelaskan matrik binari tersebut kepada kelas masing-masing. Eksperimen dijalankan terhadap model tradisional (MLP) and model hybrid (SOM dan MLP) dalam pengecaman digit ("satu" to "sembilan") dan pengecaman perkataan 2-suku (30 perkataan yang dipilih). Eksperimen ini bertujuan untuk mendapat ketepatan pengecaman dengan menggunakan nilai parameter yang berbeza (dimensi cepstral, dimensi SOM, bilangan nod tersembunyi dan kadar pembelajaran). Kedua-dua model ini juga diuji dengan menggunakan dua teknik pengelasan: pengelasan mengikut suku perkataan dan perkataan. Perbandingan dan perbincangan telah dibuat berdasarkan ketepatan pengecaman masing-masing. Keputusan eksperimen menunjukkan bahawa model kami mencapai ketepatan yang lebih tinggi.

## TABLE OF CONTENTS

<b>CHAPTER</b>	<b>TITLE</b>	<b>PAGE</b>
	<b>DECLARATION</b>	ii
	<b>DEDICATION</b>	iii
	<b>ACKNOWLEDGEMENT</b>	iv
	<b>ABSTRACT</b>	v
	<b>ABSTRAK</b>	vi
	<b>TABLE OF CONTENTS</b>	vii
	<b>LIST OF TABLES</b>	xiv
	<b>LIST OF FIGURES</b>	xviii
	<b>LIST OF ABBREVIATIONS</b>	xxiii
	<b>LIST OF SYMBOLS</b>	xxv
	<b>LIST OF APPENDICES</b>	xxvi
<b>1</b>	<b>INTRODUCTION</b>	
	1.1 Introduction	1
	1.2 Background of Study	2
	1.3 Problem Statements	4
	1.4 Aim of the Research	5
	1.5 Objectives of the Research	5
	1.6 Scopes of the Research	5
	1.7 Justification	6
	1.8 Thesis Outline	8

**2****REVIEW OF SPEECH RECOGNITION AND  
NEURAL NETWORK**

2.1	Fundamental of Speech Recognition	10
2.2	Linear Predictive Coding (LPC)	11
2.3	Speech Recognition Approaches	16
2.3.1	Dynamic Time Warping (DTW)	16
2.3.2	Hidden Markov Model (HMM)	17
2.3.3	Artificial Neural Network (ANN)	18
2.4	Comparison between Speech Recognition Approaches	20
2.5	Review of Artificial Neural Networks	21
2.5.1	Processing Units	21
2.5.2	Connections	22
2.5.3	Computation	22
2.5.4	Training	23
2.6	Types of Neural Networks	24
2.6.1	Supervised Learning	24
2.6.2	Semi-Supervised Learning	25
2.6.3	Unsupervised Learning	25
2.6.4	Hybrid Networks	26
2.7	Related Research	27
2.7.1	Phoneme/Subword Classification	27
2.7.2	Word Classification	29
2.7.3	Classification Using Hybrid Neural Network Approach	31
2.8	Summary	32

**3****SPEECH DATASET DESIGN**

3.1	Human Speech Production Mechanism	33
3.2	Malay Morphology	35
3.2.1	Primary Word	35
3.2.2	Derivative Word	38
3.2.3	Compound Word	39

3.2.4	Reduplicative Word	39
3.3	Malay Speech Dataset Design	39
3.3.1	Selection of Malay Speech	
	Target Sounds	40
3.3.2	Acquisition of Malay Speech Dataset	44
3.4	Summary	46

## 4

### FEATURE EXTRACTION AND CLASSIFICATION ALGORITHM

4.1	The Architecture of Speech Recognition	
	System	47
4.2	Feature Extractor (FE)	48
4.2.1	Speech Sampling	49
4.2.2	Frame Blocking	50
4.2.3	Pre-emphasis	51
4.2.4	Windowing	51
4.2.5	Autocorrelation Analysis	52
4.2.6	LPC Analysis	52
4.2.7	Cepstrum Analysis	53
4.2.8	Endpoint Detection	54
4.2.9	Parameter Weighting	55
4.3	Self-Organizing Map (SOM)	55
4.3.1	SOM Architecture	57
4.3.2	Learning Algorithm	58
4.4.3	Dimensionality Reduction	63
4.4	Multilayer Perceptron (MLP)	65
4.4.1	MLP Architecture	65
4.4.2	Activation Function	66
4.4.3	Error-Backpropagation	67
4.4.4	Improving Error-Backpropagation	69
4.4.5	Implementation of Error- Backpropagation	73
4.5	Summary	74



<b>5</b>	<b>SYSTEM DESIGN AND IMPLEMENTATION</b>	
5.1	Introduction	75
5.2	Implementation of Speech Processing	76
5.2.1	Feature Extraction using LPC	76
5.2.2	Endpoint Detection	80
5.3	Implementation of Self-Organizing Map	91
5.4	Implementation of Multilayer Perceptron	97
5.4.1	MLP Architecture for Digit Recognition	97
5.4.2	MLP Architecture for Word Recognition	98
5.4.3	Implementation of MLP	99
5.5	Experiment Setup	107
<b>6</b>	<b>RESULTS AND DISCUSSION</b>	
6.1	Introduction	109
6.2	Testing of Digit Recognition	111
6.2.1	Testing Results for Conventional System	111
6.2.1.1	Experiment 1: Optimal Cepstral Order (CO)	111
6.2.1.2	Experiment 2: Optimal Hidden Node Number (HNN)	112
6.2.1.3	Experiment 3: Optimal Learning Rate (LR)	114
6.2.2	Results for Proposed System Testing	115
6.2.2.1	Experiment 1: Optimal Cepstral Order (CO)	115
6.2.2.2	Experiment 2: Optimal Dimension of SOM (DSOM)	116
6.2.2.3	Experiment 3: Optimal Hidden Node Number (HNN)	117

6.2.2.4	Experiment 4: Optimal Learning Rate (LR)	119
6.2.3	Discussion for Digit Recognition Testing	120
6.2.3.1	Comparison of Performance for DRCS and DRPS (CO)	120
6.2.3.2	Comparison of Performance for DRCS and DRPS (HNN)	121
6.2.3.3	Comparison of Performance for DRCS and DRPS (LR)	123
6.2.3.4	Discussion on Performance of DRPS according to DSOM	124
6.2.3.5	Summary for Digit Recognition Testing	125
6.3	Testing of Word Recognition	126
6.3.1	Results for Conventional System Testing (Syllable Classification)	126
6.3.1.1	Experiment 1: Optimal Cepstral Order (CO)	126
6.3.1.2	Experiment 2: Optimal Hidden Node Number (HNN)	127
6.3.1.3	Experiment 3: Optimal Learning Rate (LR)	128
6.3.2	Results for Conventional System Testing (Word Classification)	130
6.3.2.1	Experiment 1: Optimal Cepstral Order (CO)	130
6.3.2.2	Experiment 2: Optimal Hidden Node Number (HNN)	131
6.3.2.3	Experiment 3: Optimal Learning Rate (LR)	132
6.3.3	Results for Proposed System Testing (Syllable Classification)	133

6.3.3.1 Experiment 1: Optimal	
Cepstral Order (CO)	133
6.3.3.2 Experiment 2: Optimal	
Dimension of SOM (DSOM)	135
6.3.3.3 Experiment 3: Optimal	
Hidden Node Number (HNN)	136
6.3.3.4 Experiment 4: Optimal	
Learning Rate (LR)	137
6.3.4 Results for Proposed System	
Testing (Word Classification)	138
6.3.4.1 Experiment 1: Optimal	
Cepstral Order (CO)	138
6.3.4.2 Experiment 2: Optimal	
Dimension of SOM (DSOM)	140
6.3.4.3 Experiment 3: Optimal	
Hidden Node Number (HNN)	141
6.3.4.4 Experiment 4: Optimal	
Learning Rate (LR)	142
6.3.5 Discussion for Word Recognition	
Testing	143
6.3.5.1 Comparison of Performance	
for WRCS and WRPS	
according to CO	143
6.3.5.2 Comparison of Performance	
for WRCS and WRPS	
according to HNN	146
6.3.5.3 Comparison of Performance	
for WRCS and WRPS	
according to LR	149
6.3.5.4 Comparison of Performance	
of WRPS according to DSOM	152

6.3.5.5 Comparison of Performance for WRCS and WRPS according to Type of Classification	155
6.3.5.6 Summary of Discussion for Word Recognition	156
6.4 Summary	157
<b>7 CONCLUSION AND SUGGESTION</b>	
7.1 Conclusion	158
7.2 Directions for Future Research	159
<b>REFERENCES</b>	162
<b>PUBLICATIONS</b>	170
<b>Appendices A – V</b>	171 - 192

## LIST OF TABLES

<b>TABLE NO.</b>	<b>TITLE</b>	<b>PAGE</b>
1.1	Comparison of different speech recognition systems.	7
2.1	The comparison between different speech recognition approaches.	20
2.2	The performance comparison between different speech recognition approaches.	20
3.1	Structure of words with one syllable.	36
3.2	Structure of words with two syllables.	37
3.3	Structure of words with three syllables or more.	38
3.4	15 selected syllables in order to form two-syllable words as target sounds.	41
3.5	Two-syllable Malay words combined using 15 selected syllables.	42
3.6	30 selected Malay two-syllable words as the speech target sounds.	43
3.7	10 selected digit words as the speech target sounds for digit recognition.	44
3.8	Specification of dataset for word recognition	45
3.9	Specification of dataset for digit recognition	45
5.1	The setting of the target values for MLP in digit recognition.	98
5.2	The setting of the target values for MLP (syllable classification).	100

5.3	The setting of the target values for MLP (word classification).	101
6.1	Recognition accuracy for different CO for Experiment 1 (DRCS)	111
6.2	Recognition accuracy for different HNN for Experiment 2 (DRCS)	113
6.3	Recognition accuracy for different LR for Experiment 3 (DRCS)	114
6.4	Recognition accuracy for different CO for Experiment 1 (DRPS)	115
6.5	Recognition accuracy for different DSOM for Experiment 2 (DRPS)	116
6.6	Recognition accuracy for different HNN for Experiment 3 (DRPS)	118
6.7	Recognition accuracy for different LR for Experiment 4 (DRPS)	119
6.8	Comparison of performance for DRCS and DRPS according to CO	120
6.9	Comparison of performance for DRCS and DRPS according to HNN	122
6.10	Comparison of performance for DRCS and DRPS according to LR	123
6.11	The optimal parameters and the architecture for DRPS	125
6.12	Recognition accuracy for different CO for Experiment 1 (WRCS(S))	126
6.13	Recognition accuracy for different HNN for Experiment 2 (DRCS)	127
6.14	Recognition accuracy for different LR for Experiment 3 (WRCS(S))	129
6.15	Recognition accuracy for different CO for Experiment 1 (WRCS(W))	130
6.16	Recognition accuracy for different HNN for Experiment 2 (WRCS(W))	131

6.17	Recognition accuracy for different LR for Experiment 3 (WRCS(W))	132
6.18	Recognition accuracy for different CO for Experiment 1 (WRPS(S))	134
6.19	Recognition accuracy for different DSOM for Experiment 2 (WRPS(S))	135
6.20	Recognition accuracy for different HNN for Experiment 3 (WRPS(S))	136
6.21	Recognition accuracy for different LR for Experiment 4 (WRPS(S))	137
6.22	Recognition accuracy for different CO for Experiment 1 (WRPS(W))	139
6.23	Recognition accuracy for different DSOM for Experiment 2 (WRPS(W))	140
6.24	Recognition accuracy for different HNN for Experiment 3 (WRPS(W))	141
6.25	Recognition accuracy for different LR for Experiment 4 (WRPS(W))	142
6.26	Comparison of performance for WRCS(S) and WRPS(S) according to CO.	144
6.27	Comparison of performance for WRCS(W) and WRPS(W) according to CO.	145
6.28	Comparison of performance for WRCS(S) and WRPS(S) according to HNN.	146
6.29	Comparison of performance for WRCS(W) and WRPS(W) according to HNN.	147
6.30	Comparison of performance for WRCS(S) and WRPS(S) according to LR.	149
6.31	Comparison of performance for WRCS(W) and WRPS(W) according to LR.	150
6.32	Comparison of performance for WRPS according to DSOM.	152
6.33	Results of testing for WRCS and WRPS according to type of classification	155

6.34	The optimal parameters and the architecture for WRPS(S).	156
------	--	-----



## LIST OF FIGURES

FIGURE NO.	TITLE	PAGE
1.1	Feature map with neurons (circles) which is labeled with the symbols of the phonemes to which they “learned” to give the best responses.	3
1.2	The sequence of the responses obtained from the trained feature map when the Finnish word <i>humppila</i> was uttered.	3
2.1	Basic model of speech recognition system.	11
2.2	The current speech sample is predicted as a linear combination of past $p$ samples. ( $n$ = total number of speech sample).	12
2.3	Dynamic Time Warping (DTW)	17
2.4	A basic architecture of Multilayer Perceptron (MLP)	19
2.5	A basic neuron processing unit	22
2.6	Neural network topologies: (a) Unstructured, (b) Layered, (c) Recurrent and (d) Modular	23
2.7	Perceptrons: (a) Single-layer Perceptron (b) Multilayer Perceptron	25
2.8	Decision regions formed by a 2-layer Perceptron using backpropagation training and vowel formant data.	28
3.1	The vocal tract	34
3.2	Structure of one-syllable word “ <i>Ya</i> ” and “ <i>Stor</i> ”.	36
3.3	Structure of two-syllable word “ <i>Guru</i> ” and “ <i>Jemput</i> ”.	37
4.1	Proposed speech recognition model	47

4.2	Feature Extractor (FE) schematic diagram	48
4.3	Figure 4.3: Speech signal for the word <i>kosong01.wav</i> sampled at 16 kHz with a precision of 16 bits.	49
4.4	Blocking of speech waveform into overlapping frames with $N$ analysis frame length and $M$ shifting length.	50
4.5	Cepstral coefficient of <i>BU.cep</i>	54
4.6	SOM transforms feature vectors generated by speech processing into binary matrix which performs dimensionality reduction.	56
4.7	The 2-D SOM architecture	57
4.8	Flow chart of SOM learning algorithm	61
4.9	Trained feature map after 1,250,000 iterations.	62
4.10	Dimensionality reduction performed by SOM.	63
4.11(a)	The 12 x 12 mapping of binary matrix of <i>/bu/</i> syllable.	64
4.11(b)	Binary matrix of <i>/bu/</i> which is fed as input for MLP.	64
4.12	A three-layer Multilayer Perceptron	66
4.13	The determination of hidden node number using Geometric Pyramid Rule (GPR).	71
4.14	Flow chart of error-backpropagation algorithm	73
5.1	The implementation of speech recognition system	75
5.2(a)	The detected boundaries of <i>sembilan04.wav</i> using rms energy in Level 1 of Initial endpoint detection	90
5.2(b)	The detected boundaries of <i>sembilan04.wav</i> using zero crossing rate in Level 2 of Initial endpoint detection	90
5.2(c)	The actual boundaries of <i>sembilan04.wav</i> using Euclidean distance of cepstrum in Level 3 of Actual endpoint detection	90
5.3	The architecture of Self-Organizing Map (SOM)	91
5.4	MLP with 10 output. The 10 output nodes correspond to 10 Malay digit words respectively.	97

5.5	MLP with 15 output nodes. The 15 output nodes correspond to 15 Malay syllables respectively.	99
5.6	MLP with 30 output nodes. The 30 output nodes correspond to 30 Malay two-syllable words respectively.	99
5.7	System architecture for conventional model (single-network)	108
5.8	System architecture for proposed model (hybrid network)	108
5.9	Training and testing of the digit recognition system	109
5.10	Training and testing of the word recognition system	109
6.1	Presentation and discussion of the results of the tests in table and graph form in stages.	110
6.2	Recognition accuracy for different CO for Experiment 1 (DRCS)	112
6.3	Recognition accuracy for different HNN for Experiment 2 (DRCS)	113
6.4	Recognition accuracy for different LR for Experiment 3 (DRCS)	114
6.5	Recognition accuracy for different CO for Experiment 1 (DRPS)	116
6.6	Recognition accuracy for different DSOM for Experiment 2 (DRPS)	117
6.7	Recognition accuracy for different HNN for Experiment 3 (DRPS)	118
6.8	Recognition accuracy for different LR for Experiment 4 (DRPS)	119
6.9	Analysis of comparison of performance for DRCS and DRPS according to CO.	121
6.10	Analysis of comparison of performance for DRCS and DRPS according to HNN.	122
6.11	Analysis of comparison of performance for DRCS and DRPS according to LR.	124

6.12	Recognition accuracy for different CO for Experiment 1 (WRCS(S))	127
6.13	Recognition accuracy for different HNN for Experiment 2 (WRCS(S))	128
6.14	Recognition accuracy for different LR for Experiment 3 (WRCS(S))	129
6.15	Recognition accuracy for different CO for Experiment 1 (WRCS(W))	130
6.16	Recognition accuracy for different HNN for Experiment 2 (WRCS(W))	132
6.17	Recognition accuracy for different LR for Experiment 3 (WRCS(W))	133
6.18	Recognition accuracy for different CO for Experiment 1 (WRPS(S))	134
6.19	Recognition accuracy for different DSOM for Experiment 2 (WRPS(S))	135
6.20	Recognition accuracy for different HNN for Experiment 3 (WRPS(S))	137
6.21	Recognition accuracy for different LR for Experiment 4 (WRPS(S))	138
6.22	Recognition accuracy for different CO for Experiment 1 (WRPS(W))	139
6.23	Recognition accuracy for different DSOM for Experiment 2 (WRPS(W))	140
6.24	Recognition accuracy for different HNN for Experiment 3 (WRPS(W))	142
6.25	Recognition accuracy for different LR for Experiment 4 (WRPS(W))	143
6.26	Comparison of performance for WRCS(S) and WRPS(S) according to CO.	144
6.27	Comparison of performance for WRCS(W) and WRPS(W) according to CO.	145
6.28	Comparison of performance for WRCS(S) and WRPS(S) according to HNN.	147

6.29	Comparison of performance for WRCS(W) and WRPS(W) according to HNN.	148
6.30	Comparison of performance for WRCS(S) and WRPS(S) according to LR.	150
6.31	Comparison of performance for WRCS(W) and WRPS(W) according to LR.	151
6.32	Comparison of performance for WRPS according to DSOM.	153
6.33(a)	Matrix mapping of “ <i>buku</i> ” word where the arrows show the direction of the sequence of phonemes.	154
6.33(b)	Matrix mapping of “ <i>kubu</i> ” word where the arrows show the direction of the sequence of phonemes.	154
6.34	Analysis of comparison of performance for WRCS and WRPS according to syllable classification and word classification.	155

## LIST OF ABBREVIATIONS

AI	-	Artificial Intelligence
ANN	-	Artificial Neural Network
BMU	-	Best Matching Unit
BP	-	Back-Propagation
CO	-	Cepstral Order
CS	-	Conventional System
DR	-	Digit Recognition
DRCS	-	Digit Recognition Conventional System
DRPS	-	Digit Recognition Proposed System
DSOM	-	Dimension of Self-Organizing Map
DTW	-	Dynamic Time Warping
FE	-	Feature Extractor
GPR	-	Geometric Pyramid Rule
HMM	-	Hidden Markov Model
HNN	-	Hidden Node Number
KSOM	-	Kohonen Self-Organization Network
LP	-	Linear Prediction
LPC	-	Linear Predictive Coding
LR	-	Linear Rate
LVQ	-	Learning Vector Quantization
MLP	-	Multilayer Perceptron
PARCOR	-	Partial-Correlation
PC	-	Personal Computer
PS	-	Proposed System
SAMSOM	-	Structure Adaptive Multilayer Self-Organizing Map
SLP	-	Single-layer Perceptron
SOM	-	Self-Organizing Map

TDNN	-	Time-Delay Neural Network
VQ	-	Vector Quantization
WPF	-	Winning Probability Function
WR	-	Word Recognition
WRCS	-	Word Recognition Conventional System
WRCS(S)	-	Word Recognition Conventional System using Syllable Classification
WRCS(W)	-	Word Recognition Conventional System using Word Classification
WRPS	-	Word Recognition Proposed System
WRPS(S)	-	Word Recognition Proposed System using Syllable Classification
WRPS(W)	-	Word Recognition Proposed System using Word Classification

## LIST OF SYMBOLS

$s$	-	Speech sample
$\hat{s}$	-	Predicted speech sample
$a$	-	Predicted coefficient
$e$	-	Prediction error
$E$	-	Mean squared error (LPC)
$E$	-	Energy power (Endpoint detection)
$Z$	-	Zero-crossing
$T$	-	Threshold (Endpoint detection)
$D$	-	Weighted Euclidean distance
$R$	-	Autocorrelation function
$w$	-	Hamming window
$p$	-	The order of the LPC analysis
$k$	-	PARCOR coefficients
$c$	-	Cepstral coefficients
$X$	-	Input nodes
$Y$	-	Output nodes
$H$	-	Hidden nodes
$M$	-	Weights
$B$	-	Bias
$\sigma$	-	Width of lattice (SOM)
$\lambda$	-	Time constant (SOM)
$\alpha$	-	Learning rate (SOM)
$\Theta$	-	The amount of influence a node's distance from the BMU (SOM)
$\eta$	-	Learning rate (MLP)
$\delta$	-	Error information term



**LIST OF APPENDICES**

<b>APPENDIX</b>	<b>TITLE</b>	<b>PAGE</b>
A	Specification of test on optimal Cepstral Order for DRCS	171
B	Specification of test on optimal Hidden Node Number for DRCS	172
C	Specification of test on optimal Learning Rate for DRCS	173
D	Specification of test on optimal Cepstral Order for DRPS	174
E	Specification of test on optimal Dimension of SOM for DRPS	175
F	Specification of test on optimal Hidden Node Number for DRPS	176
G	Specification of test on optimal Learning Rate for DRPS	177
H	Specification of test on optimal Cepstral Order for WRCS(S)	178
I	Specification of test on optimal Hidden Node Number for WRCS(S)	179
J	Specification of test on optimal Learning Rate for WRCS(S)	180
K	Specification of test on optimal Cepstral Order for WRCS(W)	181
L	Specification of test on optimal Hidden Node Number for WRCS(W)	182

M	Specification of test on optimal Learning Rate for WRCS(W)	183
N	Specification of test on optimal Cepstral Order for WRPS(S)	184
O	Specification of test on optimal Dimension of SOM for WRPS(S)	185
P	Specification of test on optimal Hidden Node Number for WRPS(S)	186
Q	Specification of test on optimal Learning Rate for WRPS(S)	187
R	Specification of test on optimal Cepstral Order for WRPS(W)	188
S	Specification of test on optimal Dimension of SOM for WRPS(W)	189
T	Specification of test on optimal Hidden Node Number for WRPS(W)	190
U	Specification of test on optimal Learning Rate for WRPS(W)	191
V	Convergences file (dua12.cep) which shows the rms error in each epoch.	192

## CHAPTER 1

### INTRODUCTION

#### 1.1 Introduction

By 1990, many researchers had demonstrated the value of neural networks for important task like phoneme recognition and spoken digit recognition. However, it is still unclear whether connectionist techniques would scale up to large speech recognition tasks. There is a large variety in the speech recognition technology and it is important to understand the differences between the technologies. Speech recognition system can be classified according to the type of speech, size of the vocabulary, the basic units and the speaker independence. The position of a speech recognition system in these dimensions determines which algorithm can or has to be used. Speech recognition has been another proving ground for neural networks. Some researchers achieved good results in such basic tasks as voiced/unvoiced discrimination (Watrous, 1988), phoneme recognition (Waibel *et al.*, 1989), and spoken digit recognition (Peeling and Moore, 1987). However, research in finding a good neural network model for robust speech recognition still has a wide potential to be developed.

Why does the speech recognition problem attract researchers? If an efficient speech recognizer is produced, a very natural human-machine interface would be obtained. By natural means something that is intuitive and easy to be used by a person, a method that does not require special tools or machines but only the natural capabilities that every human possesses. Such a system could be used by any person who is able to speak and will allow an even broader use of machines, specifically computers.

## 1.2 Background of Study

Neural network classifier has been compared with other pattern recognition classifiers and is explored as an alternative to other speech recognition techniques. Lippman (1989) has proposed a static model which is employed as an input pattern of Multilayer Perceptron (MLP) network. The conventional neural network (Pont *et al.*, 1996; Ahkuputra *et al.*, 1998; Choubassi *et al.*, 2003) defines a network as consisting of a few basic layers (input, hidden and output) in a Multilayer Perceptron type of topology. Then a training algorithm such as backpropagation is applied to develop the interconnection weights. This conventional model or system has also been used in a variety of pattern recognition and control applications that are not effectively handled by other AI paradigms.

However, there are some difficulties in using MLP alone. The most major difficulty is that, increasing the number of connections not only increases the training time but also makes it more probable to fall in a poor local minima. It also necessitates more data for training. Perceptron as well as Multilayer Perceptron (MLP) usually needs input pattern of fixed length (Lippman, 1989). This is the reason why the MLP has difficulties when dealing with temporal information (essential speech information or feature extracted during speech processing). Since the word has to be recognized as a whole, the word boundaries are often located automatically by endpoint detector and the noise is removed outside of the boundaries. The word patterns have to be also warped using some pre-defined paths in order to obtain fixed length word patterns.

Since the early eighties, researchers have been using neural networks in the speech recognition problem. One of the first attempts was Kohonen's electronic typewriter (Kohonen, 1992). It uses the clustering and classification characteristics of the Self-Organizing Map (SOM) to obtain an ordered feature map from a sequence of feature vectors which is shown in Figure 1.1. The training was divided into two stages, where the first stage was used to obtain the SOM. Speech feature vectors were fed into the SOM until it converged. The second training stage consisted in labeling the SOM, as example, each neuron of the feature map was assigned a phoneme label. Once the labeling process was completed, the training process

ended. Then, unclassified speech was fed into the system, which was then translated it into a sequence of labels. Figure 1.2 shows the sequence of the responses obtained from the trained feature map when the Finnish word *humppila* was uttered. This way, the feature extractor plus the SOM behaved like a transducer, transforming a sequence of speech samples into a sequence of labels. Then, the sequence of labels was processed by some AI scheme (Grammatical Transformation Rules) in order to obtain words from it.

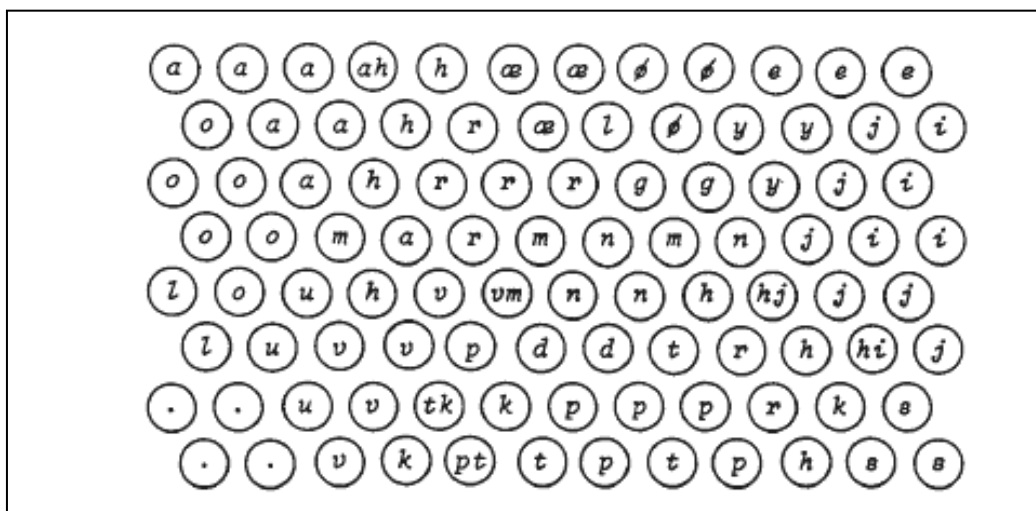


Figure 1.1: Feature map with neurons (circles) which is labeled with the symbols of the phonemes to which they “learned” to give the best responses.

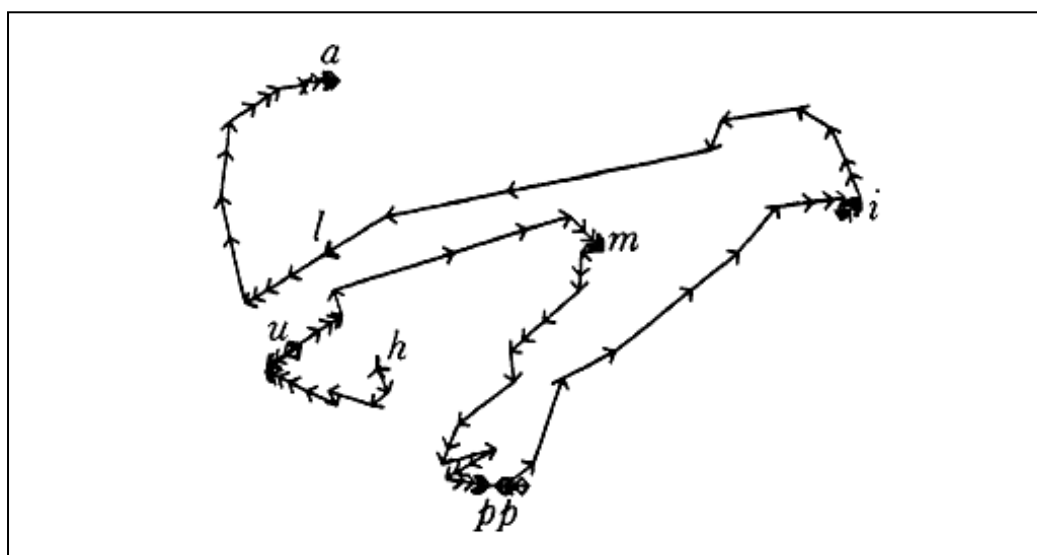


Figure 1.2: The sequence of the responses obtained from the trained feature map when the Finnish word *humppila* was uttered.

Usage of an unsupervised learning neural network as well as SOM seems to be wise. The SOM constructs a topology preserving mapping from the high-dimensional space onto map units (neurons) in such a way that relative distances between data points are preserved. The way SOM performs dimensionality reduction is by producing a map of usually 2 dimensions which plot the similarities of the data by grouping similar data items together. Because of its characteristic which is able to form an ordered feature map, the SOM is found to be suitable for dimensionality reduction of speech feature. Forming a binary matrix to feed to the MLP makes the training and classification simpler and better. Such a hybrid system consists of two neural-based models, a SOM and a MLP. The hybrid system mostly tries to overcome the problem of the temporal variation of utterances where the utterances for same word by same speaker may be different in duration and speech rate).

### 1.3 Problem Statements

According to the background of study, here are the problem statements:

- i. Various approaches have been introduced for Malay speech recognition in order to produce an accurate and robust system for Malay speech recognition. However, there are only a few approaches which have achieved excellent performance for Malay speech recognition (Ting *et al.*, 2001a, 2001b and 2001c; Md Sah Haji Salam *et al.*, 2001). Thus, research in speech recognition for Malay language still has a wide potential to be developed.
- ii. Multilayer Perceptron (MLP) has difficulties when dealing with temporal information. Since the word has to be recognized as a whole, the word patterns have to be warped using some pre-defined paths in order to obtain fixed length word patterns (Tebelskis, 1995; Gavat *et al.*, 1998). Thus, an efficient model is needed to improve this drawback.
- iii. Self-Organizing Map (SOM) is considered as a suitable and effective approach for both clustering and dimensionality reduction. However, is SOM an efficient neural network to be applied in MLP-based speech recognition in order to reduce the dimensionality of feature vector?

#### 1.4 Aim of the Research

The aim of the research is to investigate how hybrid neural network can be applied or utilized in speech recognition area and propose a hybrid model by combining Self-Organizing Map (SOM) and Multilayer Perceptron (MLP) for Malay speech recognition in order to achieve a better performance compared to conventional model (single network).

#### 1.5 Objectives of the Research

- i. Studying the effectiveness of various types of neural network models in terms of speech recognition.
- ii. Developing a hybrid model/approach by combining SOM and MLP in speech recognition for Malay language.
- iii. Developing a prototype of Malay speech recognition which contains three main components namely speech processing, SOM and MLP.
- iv. Conducting experiments to determine the optimal values for the parameters (cepstral order, dimension of SOM, hidden node number, learning rate) of the system in order to obtain the optimal performance.
- v. Comparing the performance between conventional model (single network) and proposed model (SOM and MLP) based to the recognition accuracy to prove the improvement achieved by the proposed model. The recognition accuracy is based on the calculation of percentage below:

$$\text{Recognition Accuracy (\%)} = \frac{\text{Total of Correct Recognized Word}}{\text{Total of Sample Word}}$$

#### 1.6 Scopes of the Research

The scope of the research clearly defines the specific field of the study. The discussion of the study and research is confined to the scope.

- i. There are two datasets created where one is used for digit recognition and another one is used for word recognition. The former consists of 10 Malay digits and the latter consists of 30 selected two-syllable Malay words. Speech samples are collected in a noise-free environment using unidirectional microphone.
- ii. Human speakers comprise of 3 males and 3 females. The system supports speaker-independent capability. The age of the speakers ranges between 18 – 25 years old.
- iii. Linear Predictive Coding (LPC) is used as the feature extraction method. The method is to extract the speech feature from the speech data. The LPC coefficients are determined using autocorrelation method. The extracted LPC coefficients are then converted to cepstral coefficients.
- iv. Self-Organizing Map (SOM) and Multilayer Perceptron (MLP) is applied in the proposed system. SOM acts as a feature extractor which converts the higher-dimensional feature vector into lower-dimensional binary vector. Then MLP takes the binary vectors as its input for training and classification.

## 1.7 Justification

Many researchers have worked in automatic speech recognition for almost few decades. In the eighties, speech recognition research was characterized by a shift in technology from template-based approaches (Hewett, 1989; Aradilla *et al.*, 2005) to statistical-based approaches (Gold, 1988; Huang, 1992; Siva, 2000; Zbancioc and Costin, 2003) and connectionist approaches (Watrous, 1988; Hochberg *et al.*, 1994). Instead of Hidden Markov Model (HMM), the use of neural networks has become another idea in speech recognition problems. Anderson (1999) has made a comparison between statistical-based and template-based approaches. Today's research focuses on a broader definition of speech recognition. It is not only concerned with recognizing the word content but also prosody (Shih *et al.*, 2001) and personal signature.



Despite all of the advances in the speech recognition area, the problem is far from being completely solved. A number of commercial products are currently sold in the commercial market. Products that recognize the speech of a person within the scope of a credit card phone system, command recognizers that permit voice control of different types of machines, “electronic typewriters” that can recognize continuous speech and manage several tens of thousands word vocabularies, and so on. However, although these applications may seem impressive, they are still computationally intensive, and in order to make their usage widespread more efficient algorithms must be developed. Summing up, there is still room for a lot of improvement and research.

Currently there are many speech recognition applications released, whether as a commercial or free software. The technology behind speech output has changed over times and the performance of speech recognition system is also increasing. Early system used discrete speech; *Dragon Dictate* is the only discrete speech system still available commercially today. On the other hand, the main continuous speech systems currently available for PC are *Dragon Naturally Speaking* and *IBM ViaVoice*. Table 1.1 shows the comparison of different speech recognition systems with the prototype to be built in this research. This comparison is important as it gives an insight of the current trend of speech recognition technology.

Table 1.1: Comparison of different speech recognition systems

<b>Software</b> <b>Feature</b>	<b>Dragon Dictate</b>	<b>IBM Voice</b>	<b>Naturally Speaking 7</b>	<b>Microsoft Office XP SR</b>	<b>Prototype To Be Built</b>
<b>Discrete Speech Recognition</b>	√	X	X	X	√
<b>Continuous Speech Recognition</b>	X	√	√	√	X
<b>Speaker Dependent</b>	√	√	√	√	√
<b>Speaker Independent</b>	X	X	X	X	√
<b>Speech-to-Text</b>	√	√	√	√	√
<b>Active Vocabulary Size (Words)</b>	30,000 – 60,000	22,000 – 64,000	300,000	Finite	30 – 100
<b>Language</b>	English	English	English	English	Malay

In this research, the speech recognition problem is transformed into simplified binary matrix recognition problem. The binary matrices are generated and simplified while preserving most of the useful information by means of a SOM. Then, word recognition turns into a problem of binary matrix recognition in a smaller dimensional feature space and this performs dimensionality reduction. Besides, the comparison between the single-network recognizer and hybrid-network recognizer conducted here sheds new light on future directions of research in the field. It is important to understand that it is not the purpose of this work to develop a full-scale speech recognizer but only to test proposed hybrid model and explore its usefulness in providing more efficient solutions in speech recognition.

## **1.8 Thesis Outline**

The first few chapters of this thesis provide some essential background and a summary of related work in speech recognition and neural networks.

Chapter 2 reviews the field of speech recognition, neural network and also the intersection of these two fields, summarizing both past and present approaches to speech recognition using neural networks.

Chapter 3 introduces the speech dataset design.

Chapter 4 presents the algorithms of the proposed system: speech feature extraction (Speech processing and SOM) and classification (MLP).

Chapter 5 presents the implementation of the proposed system: Speech processing, Self-Organizing Map and Multilayer Perceptron. The essential parts of the source code are shown and explained in detail.

Chapter 6 presents the experimental tests on both of the systems: conventional system and the proposed system. The tests are conducted using digit dataset for digit recognition and word dataset for word recognition. For word recognition, two classification approaches are applied such as syllable and word

classification. The tests are conducted on speaker-independent system with different values of the parameters in order to obtain optimal performance according to the recognition accuracy. Discussion and comparison of the experimental results are also included in this chapter.

Chapter 7 presents the conclusions and future works of the thesis.