

Speech Recognition Enhancement Via Robust CHMM Speech Background Discrimination

Waleed H. Abdulla and Prof. Nikola K. Kasabov
Knowledge Engineering Lab (KEL)
Information Science Department
University of Otago
New Zealand

1 Abstract

This paper describes a methodology for building robust Hidden Markov Model (HMM) based speech recognition system. The ultimate goal is to build a reliable large vocabulary isolated words speech recogniser. The model, that we are dealing with, is of continuous HMM type (CHMM). The topology selected is the left-right one as it is quite successful in speech recognition due to its consistency with the natural way of articulating the spoken words.

One important task here is to efficiently extract the spoken words from their background using 3 states CHMM and process them in isolation by another 9 states models. This is considered as a perceptual way of extracting the signal. This technique is substantially increasing the performance of the system and improving the incorporation of states' duration.

2 Introduction

A vitally important objective in implementing speech recognition system is the separation of the signal of essence from background environment as faithfully as possible. This operation has crucial effect on the overall performance of the recogniser. It is an issue to be tackled by the researchers from the early beginning of this field. The early milestone technique was using explicit features for speech non speech discrimination; such as speech signal energy and zero-crossings[1,2,7]. This technique is effective in case of low noise environment, but unreliable with the increasing noise and varied articulation manners such as breathing and clicks. The other approach was the pattern classification of voiced, unvoiced, and silence states[3, 4]. This technique implies some decision making to improve the performance of the system but it incurs a heavy computational load. Hybrid techniques were also suggested to alleviate the computational load while improving the performance. Wilpon et al. benchmarked a multispeaker digit recogniser to evaluate the effect of misaligned word boundaries on the recognition rate. The words and the reference patterns were manually extracted. The recognition rate was found to be 93%; which was the utmost value. Then misalignment procedure was practised with recognition error measured at each step. A similar experiment has been done on our system to see the recognition rate degradation due to different forced misalignments. Fig.(1) shows the contour plot of the spoken digits recognition performance under different start-

end constraints as tested on our system. The recognition rate down graded from 99%, in case of manually extracted words, to 75% due to the signal boundary misalignment. It can be noticed that the start points misalignment allowance is less than that of the end points. Recent techniques dealt with presilence and postsilence periods as pre and post states of Hidden Markov Models (HMM). During training phase the words are modelled without including the silence periods, while the silence periods are modelled as separate states. In recognition phase the pre and post silence states are concatenated to the initial and final states of the words' models. Then the maximum likelihood (or any other optimisation) procedure are followed to identify the tested words. Those HMM techniques, even they are effective, still need to concatenate the silence states during recognition phase that consequently increases the computational cost, especially with long silence periods and increasing number of models. Also, the different spikes that might be issued during silence periods such as lip flaps will be embedded in the final calculation of the model likelihood which in turn affect the performance. Other successful techniques are using neural networks (NN) to model the silence periods, but these need decision making steps to identify the positioning and the relevance of the detected silence periods[5].

This paper reveals a method that makes use of the early ideas of deleting the silence periods and recent ideas of modelling them with HMM. A method is implemented that shows superiority over the other techniques and translates it into great potential in recognition performance.

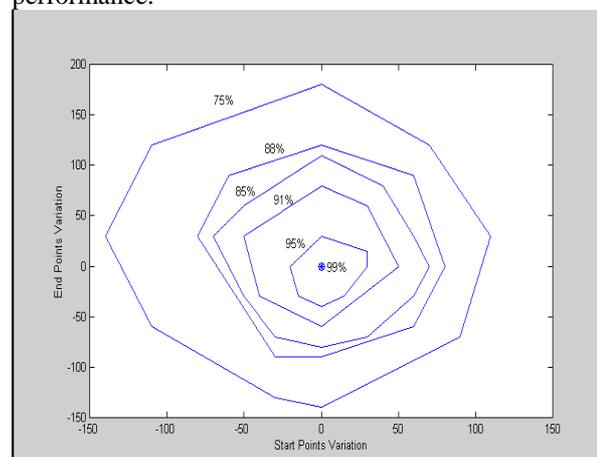


Fig.(1) Recognition performance as a function of start and end points detection.

3 System Modelling

The main modelling here is based on Continuous Hidden Markov Modelling (CHMM) technique[6,7,8,9]. During training a model is constructed from a collection of words, including their silence periods. The model learns only the first and last states of words as they represent the pre and post silence periods. Throughout the recognition phase the candidate words are first aligned with this all words model and the speech samples belonging to the first and last states are removed. Then the extracted speech segment is aligned with the different words' models to select the most probable model to issue the spoken word.

The modelling of the system implies two steps: training and recognition.

3.1 Training

There are two types of models to train - word extraction model, and word recognition models.

3.1.1 Word Extraction Model

The data set here consists of words from many different words spoken in different background environments (not necessarily all the words of the recogniser), and normally selected in the range of 50-100 words.

The model parameters are trained for left-right CHMM with 3 states which represent presilence-speech-postsilence states respectively. Unimodal modelling is used in parameter estimation. Multimixtures could be used here but it does not show more difference in performance over the unimodal one. The observations are Mel scale coefficients of the speech signal frames with only 13 coefficients (12 mels plus one power coefficient). The delta coefficients are not included to make the model insensitive to the dynamic behaviour of the signal and then gives more stable background detection. The speech frames for building the model are selected to be 23 ms taken each 9 ms. This model is called an all-words model due to its way of training. The spoken signal is de-noised using wavelet technique as it is very efficient in redundancy removal and muting the external noise. The noise factor is crucial in deciding the beginning and the ending of the words. High level spiky noise signal could easily trigger a false new state. The underlying model for the noisy signal has the following form:

$$s(n) = f(n) + \sigma e(n)$$

where f is the signal and e is the noise sampled at time n , while σ is the noise amplitude factor.

The task of de-noising is to faithfully recovering the signal by suppressing the noise part. The wavelet de-noising relies in its efficient work on the signal decomposition into approximations and details tree structure. During decomposition the signal is filtered into approximations and details components. The approximations are the high-scale, low-frequency components of the signal. The details are the low-scale, high frequency components. The process of decomposition can be repeated for many levels depending on what we are looking for in the

signal[10],[11]. During any level of decomposition the small details can be removed without substantially affecting the main features of the signal. This property introduces the idea of thresholding which set to zero all the coefficients that are below certain threshold level. After thresholding the original signal can be reconstructed again as a clean signal(without noise). The key issue here is the decision of the type of wavelet, the level of decomposition, and the thresholding technique used.

In our approach the de-noising is done using the following steps:

3.1.1a – Signal Decomposition

The wavelet chosen for this task is symlet of form sym4 and decomposed up to level 8 [12]. This wavelet is a modification of the Daubechies family wavelets and it has more symmetry with great simplicity. The order is selected experimentally to achieve a compromise between best noise reduction and minimum decomposition level.

3.1.1b Thresholding

A fixed form thresholding is selected for each level from 1 to 8 and applied to the details coefficients to mute the noise.

3.1.1c Reconstruction

The signal is reconstructed based on the original approximation coefficients of level 8 and the modified detail coefficients of levels from 1 to 8.

3.1.2 Word Recognition Models

This is the process of building a model for each spoken word. In this stage 50-100 utterances of the same word are taken from different speakers to perform the data set . The training data are taken from Otago Speech Corpus, which is freely available from the Internet on:

<http://Kel.otago.ac.nz/hyspeech/corpus>

Fig.(2) shows the structure of the model used, it shows just five states for demonstration. The observation sequence in this case are the mel scales with 39 coefficients (12 mels and one power with their deltas and delta-delta coefficients). This makes the model more sensitive to the

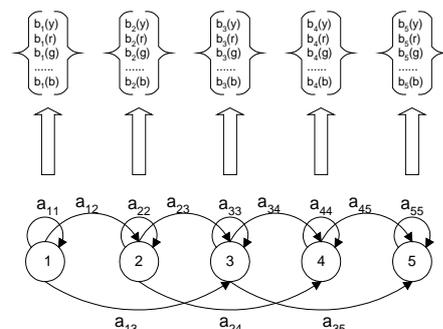


Fig.(2) CHMM left-right topology. Where a_{ij} is the state transition probability from state i to j , $b_k(o)$ is the observation probability function in state k .

dynamic behaviour of the signal which in this step is the main objective of modelling. Regarding the other parameters, the topology and the type is left-right CHMM as in Fig.(2), with 9 states and 12 mixtures. The speech frames in this case are of window length 23 ms taken each 9ms.

The probability density function (pdf) of certain observations O being in a state is considered to be of Gaussian Distribution .

The general form of $b_i(O)$ is:

$$b_i(O) = \sum_{m=1}^M c_{im} \mathfrak{N}(O; \mu_{im}, U_{im}), \quad 1 \leq i \leq N$$

where:

c_{jm} : is the m-th mixture gain coefficient in state i .

\mathfrak{N} : is the pdf distribution which is considered to be Gaussian in our case.

μ_{im} : is the mean of the m-th mixture in state i .

U_{im} : is the covariance of the m-th mixture in state i .

O : is the observations of feature dimension d .

M : is the number of mixtures used.

N : is the number of states.

In the optimisation procedure to find the best mixture distribution during training step, the Vector Quantization (VQ) technique is applied on the unimodal data set of each state. The observations belonging to each state are extracted by Viterbi

$$c_{im} = \frac{\text{number of observations being in state } i \text{ and mixture } m}{\text{total number of observations in state } i}$$

Algorithm during training and then optimised by using the maximum likelihood method. This representation results in a good modelling of the data. The following section explains this technique in more details. During this training step a model for each word is built in addition to one model for word extraction.

The state duration factor is incorporated through using heuristic technique, which boosts the performance to the same level as the correct theoretical duration inclusion with very low computational and storage costs. The state duration probability function $p_j(\tau)$ is estimated during the model training and defined as:

$p_j(\tau)$: is the probability of being in state j for τ duration, normalised to the length of observations.

The duration probability density function is considered to be Gaussian with 5 mixtures.

3.1.3 Mixture Density Components Estimation using Maximum Likelihood (ML):

The ML estimation is an optimisation technique that can be used efficiently in estimating the different component of multimixture models. We are not going through the mathematical derivations of the ML but we only describe the method used in our task.

The following definitions are used further in the paper:

$b_i(O_t)$: probability of being in state i given observation sequence O_t . It is considered of Gaussian distribution .

c_{im} : probability of being in state i with mixture m (gain coefficient).

$b_{im}(O_t)$: probability of being in state i with mixture m and given O_t .

$\Phi(w_{im}|O_t)$: probability function of being in a mixture class w_{im} given O_t in state i .

T_i : total number of observations in state i .

T_{im} : number of observations in state i with mixture m .

N : number of states.

M : number of mixtures in each state.

Now we are ready to implement the algorithm through applying the following steps:

1 – Take several versions of observations of certain word, say digit "zero", spoken several times by many speakers.

2 – Apply standard CHMM using unimodal representation ; then via Viterbi algorithm detect the states of each version of the training spoken word.

3 – Put the whole observations belonging to each state from all the versions of the spoken word into separate cells. Now we have N cells and each one represents the population of certain state derived from several observation sequences of the same word.

4 – Apply vector quantization technique to split the population of each cell into M mixtures and getting w_M classes within each state.

5 – Use any of the well known statistical methods to find the mean μ_{im} and the covariance U_{im} of each class. The gain factor c_{im} can be calculated by:

6 – Calculate $\Phi(w_{im}|O_t)$ from the following formula:

$$\Phi(w_{im} | O_t) = c_{im} \cdot \frac{b_{im}(O_t)}{b_i(O_t)}$$

7 – Find the next estimate of \hat{c}_{im} , $\hat{\mu}_{im}$, and \hat{U}_{im} from the formulas given by ML :

$$\hat{c}_{im} = \frac{1}{T_i} \sum_{t=1}^{T_i} \Phi(w_{im} | O_t)$$

$$\hat{\mu}_{im} = \frac{1}{T_{im}} \sum_{t=1}^{T_i} \Phi(w_{im} | O_t) \cdot O_t$$

$$\hat{U}_{im} = \frac{1}{T_{im}} \sum_{t=1}^{T_i} \Phi(w_{im} | O_t) \cdot (O_t - \hat{\mu}_{im})(O_t - \hat{\mu}_{im})'$$

$$\hat{b}_{im}(O_t) = \sum_{m=1}^M \hat{c}_{im} \mathfrak{N}(O; \hat{\mu}_{im}, \hat{U}_{im}), \quad 1 \leq i \leq N$$

$$\hat{b}_i(O_t) = \sum_{m=1}^M \hat{c}_{im} \hat{b}_{im}(O_t)$$

8 – Compute the next estimate of Φ by using the formula:

$$\hat{\Phi}(w_{im} | O_t) = \frac{\hat{c}_{im} \hat{b}_{im}(O_t)}{\sum_{n=1}^M \hat{c}_{in} \hat{b}_{in}(O_t)}$$

9-IF $|\Phi(w_{im} | O_i) - \hat{\Phi}(w_{im} | O_i)| \leq \epsilon$ THEN END
 ELSE Make the new value of $\Phi(w_{im} | O_i)$ equal
 the newly predicted one.
 $\Phi(w_{im} | O_i) = \hat{\Phi}(w_{im} | O_i)$

GO TO STEP 7.

Here ϵ is a very small threshold value.

3.2 Recognition

This step comprises two operations:

- 1- The input unknown utterance is submitted to the all-words model, word extraction model, to extract efficiently the spoken word from the background.
- 2- The extracted word from the previous operation is submitted to all the other models. The model that scores maximum log likelihood $\log[P(O/\lambda)]$ is representing the submitted input, where $P(O/\lambda)$ is the probability of observation O given a model λ .

The duration factor is incorporated through an efficient formula which results in improved performance.

During recognition, the states' duration are calculated from the backtracking procedure in Viterbi Algorithm. Then, the log likelihood value is incremented by the log of the duration probability value as shown below:

$$\log[\hat{P}(q, O | \lambda)] = \log[P(q, O | \lambda)] + \eta \sum_{j=1}^N \log[p_j(\tau_j)]$$

where: η is a scaling factor;

τ_j is the normalised duration of being in state j as detected by Viterbi Algorithm.

4 Results

The data sets used in training the CHMM to built the all-words models, for silence detection, has direct effect on the system performance. To build a robust word recognition model, different effects must be included in the silence periods of the data sets. The training pre- and post silence periods include: microphone clicks, sound artefacts, lip slaps, heavy breathing. The collected words should include, in their start and end phones, the most problematic phones such as weak fricatives, weak plosives and nasals. The best performance is achieved by inclusion as many effects as possible from different speakers.

The system is benchmarked against several other techniques and shows tractable results in determining the actual start and end points (boundaries) of the tested words.

The following figures show clearly how the system works in extracting the spoken word of digit "eight" from its background as compared with the explicit one. This example is chosen as it suffers from bad extraction results using the other techniques due to low energy ending fricative state as well as fading end preceded by short silence before the letter "t". The background noise level is taken to be comparable to the end level of the tested word to make it as difficult to detect as it could be. Fig(3) shows the boundary of the signal using an explicit well known technique

suggested by Rabiner et al [1]. The result is troublesome for the recogniser. Fig.(4) shows the same spoken utterance as detected by all-words model.

The time signal and spectrogram are displayed to show the correspondence between the signal and the states, which indicate how precise the model is. The states are detected through the backtracking phase in Viterbi Algorithm. The extracted signal (after removing the samples belong to states 1 and 3) is submitted to the trained models of the words to be recognised.

The recognition rate using the technique described in this paper is scoring above 98% when tested by twenty four persons, speaking the digits words in four different accents.

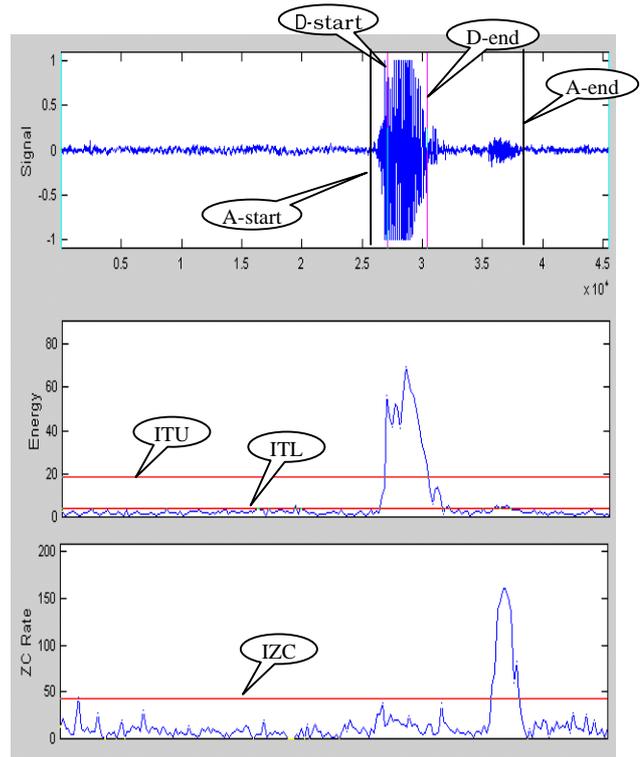


Fig.(3) End Points Detection of Spoken Digit "eight" Using Energy and Zero-Crossings Technique.

D-start and D-end are the detected start and end points of the signal, A-start and A-end are their corresponding actual start/end points. ITU and ITL are the upper and lower thresholds of the signal energy. IZC is the zero-crossings threshold.

The potential of extracting the speech signals from their backgrounds using silence states detection is not limited by words spoken in isolated mode. It can be applied on a whole spoken sentence. It shows precise discrimination capability in isolating the spoken sentence from the background environment. The inter-silence periods within the sentences are considered speech like periods as it appears clearly from Fig.(5). This figure shows the capability of discriminating the spoken sentence of connected digits " 5-6-7-8 "

A speech recognition system is implemented based on the techniques described in this paper and it will be available soon on the Internet to be tested by the researchers and to visualise the state assignment of different word and the efficient way of isolating the spoken signal from the carrying ambient. Fig.(6) shows the GUI interface of the system and one example of

assigning the different states to extracted spoken digit(0).

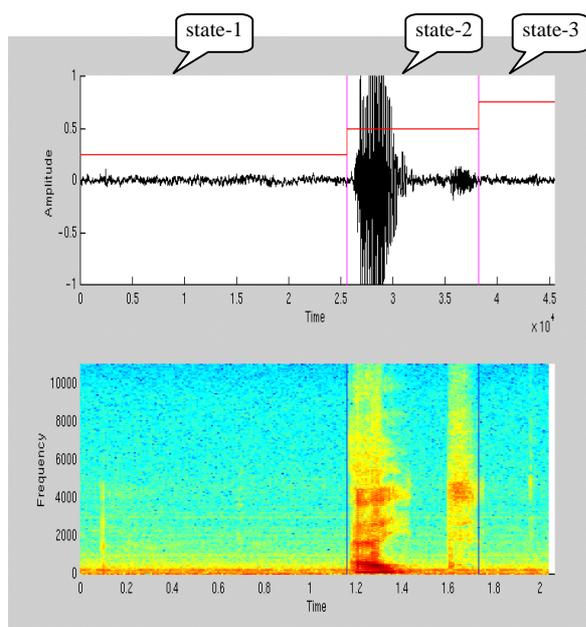


Fig.(4) Presilence-Speech_Postsilence States Discrimination.

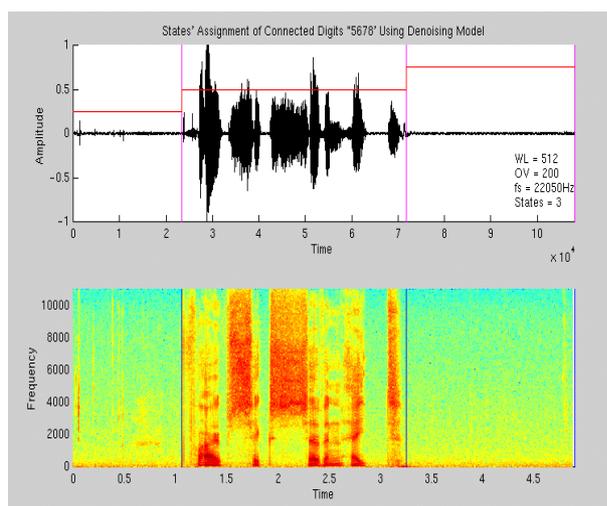


Fig.(5) Speech Signal Discrimination of the Connected Digits "5678"

5 Conclusions

The technique of efficiently extracting the speech signal from the background environment described in this paper is a development of our previous system[13]. It has better performance in dealing with noisy environment as well as less number of parameters for modelling the all-words model. It uses only 3 states to represent the presilence- speech- postsilence periods, instead of 7 states used in the previous modelling. It is as if aggregating the 5 inter states of the previous model into a single state. This can be achieved by using wavelet methodology in removing the noise from the signal before modelling the 3 states.

The technique presented in this paper offers an efficient way of extracting the speech signals from

their backgrounds. It shows superiority over the other known techniques of end points detection as it is a perceptual way which takes into consideration the input signal as well as the background status in taking the decision of signal boundary. This does not incur further or more computational cost as it might appear from the first look. The word extraction model will save at least 1/3 of the computations as the extracted signal has shorter duration than the original one (signal plus silence periods).

The CHMM could be applied on the input signal without performing word extraction but the computation in this case will be more as the duration of the signal will be longer. Accordingly, the number of states will be more to compensate for the background states. The known end points detection methods degrade the performance of the system specially in the case of low energy segments at the beginning and/or at the end of the speech signal. The all-words model which is used to extract the words from background environment is flexible and could be easily adapted to any environment just by presenting the new environment during training.

The precise signal /background separation leads to high recognition rate. The post inclusion of the normalised states' duration in the log probability equation using the way described in section 3.2 adds further reinforcement to the performance of the system to raise the recognition rate to more than 98% with multi-speaker digits data set. The ultimate recognition rate of our system is 99% as shown in Fig.(1) and it happens when the words are presented to the recogniser after manually extracted from their backgrounds.

6 Acknowledgement

This work is partially supported by grant UOO808, University of Otago, funded by the New Zealand Foundation for Research, Science, and Technology.

6 References

- [1] Lawrence R. Rabiner and M. R. Sambur. " An Algorithm for Determining the End Points of Isolated Utterances", Bell System Technical Journal (BSTJ), no. 54, pp 297-315, Feb 1975.
- [2] Lori F. Lamel, Lawrence R. Rabiner, Aaron E. Rosenberg, and Jay G. Wilpon. " An Improved End Points Detector for Isolated word Recognition", IEEE ASSP , vol. 29, no. 4, pp 777-785, Aug 1981.
- [3] Bishnu S. Atal and Lawrence R. Rabiner. " A pattern Recognition Approach to Voiced-Unvoiced-Silence Classification with Application to Speech Recognition", IEEE ASSP vol. 24, no. 4, pp 201-212, June 1976.
- [4] J. G. Wilpon, L. R. Rabiner, and T. B. Martin, "An Improved Word -Detection Algorithm for Telephone

Quality Speech Incorporating both Syntactic and Semantic Constraints," AT & T Tech. J. vol. 63, no. 3, pp 479-498, March 1984.

[5] N. K. Kasabov, "Foundations of Neural Networks, Fuzzy Systems, and Knowledge Engineering", Cambridge, MIT Press, 1996, Chap. 5.

[6] L. R. Rabiner and Biing-Hwang Juang, "Fundamentals of Speech Recognition", Printice Hall Signal Processing Series, Alan V. Oppenheim, Series Editor, 1993, Chap. 4.

[7] J. R. Deller, J. G. Proakis, and J. H. Hansen, "Discrete-Time Processing of Speech Signals", New York: Macmillan Publishing, 1993, Chap. 4, 12.

[8] W. H. Abdulla and N. K. Kasabov, "The Concepts of Hidden Markov Model in Speech Recognition", Technical Report TR99/09, University of Otago, July 1999.

[9] L. R. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition" Proceedings of the IEEE, Vol. 77, No. 2, pp. 257-286, Feb. 1989.

[10] A. Graps, "An Introduction to Wavelets", IEEE Computational Science and Engineering, Vol. 2, No. 2, 1995.

[11] M. Misiti, Y. Misiti, G. Oppenheim, J.M. Poggi, "Wavelet Toolbox", Math Works Inc., 1996.

[12] I. Daubechies, "Ten Lectures on Wavelets", SIAM, pp.254-257, 1992.

[13] W.H. Abdulla, N.K. Kasabov "Two Pass Hidden Markov Model for Speech Recognition Systems", to be published in Proceedings of the ICICS'99, Singapore, Dec. 1999.

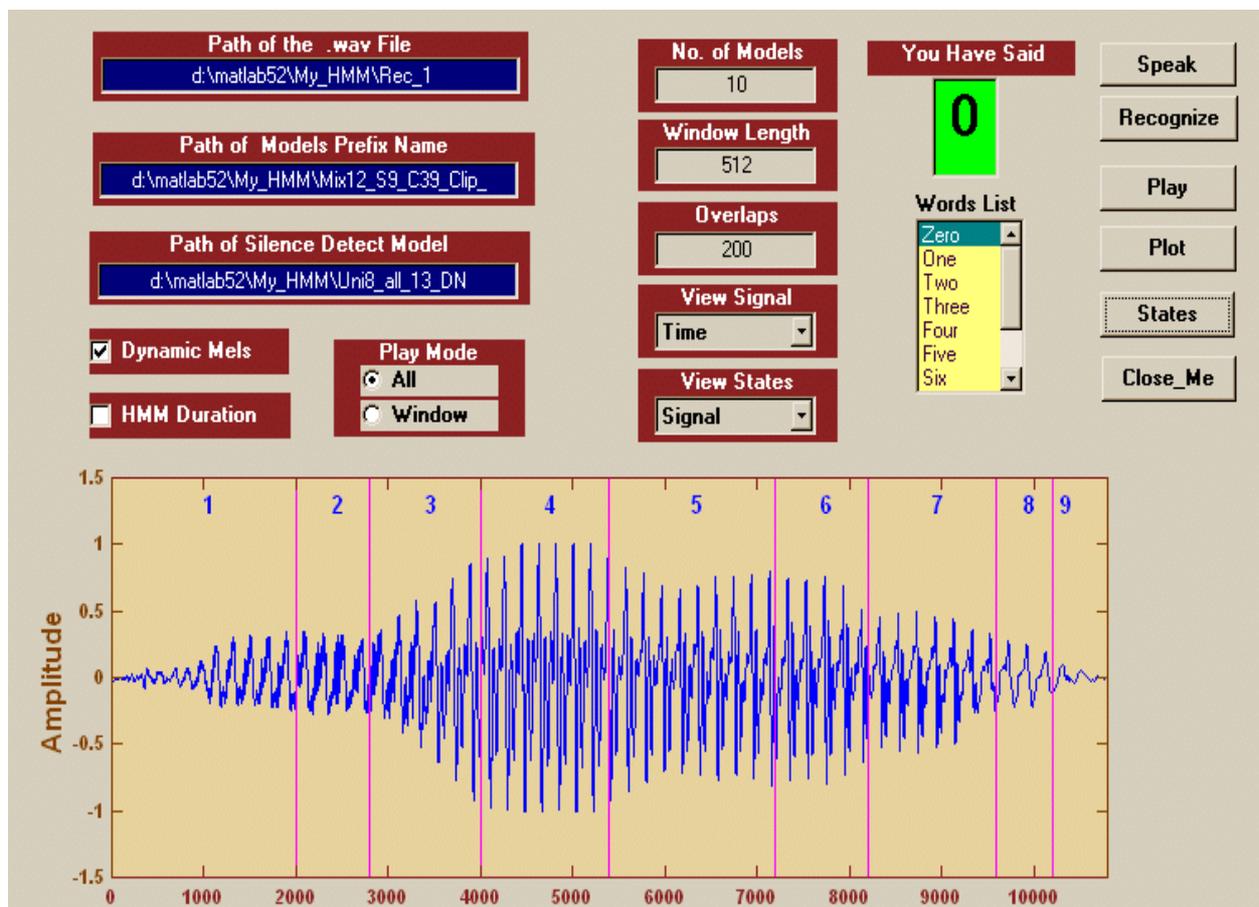


Fig.(6) Graphical User Interface of the Speech Recognition System.