# Glottal Excitation Feature based Gender Identification System using Ergodic HMM

R. Rajeshwara Rao
DVR College of Engineering & Technology
Department of CSE, Hyderabad, AP, India

A. Prasad
Vignan University, Department of MCA
Guntur, AP, India

## ABSTRACT

In this paper, through different experimental studies it is demonstrated that the time varying glottal excitation component of speech can be exploited for text independent gender recognition studies. Linear prediction (LP) residual is used as a representation of excitation information in speech. The gender-specific information in the excitation of voiced speech is captured using the Hidden Markov Models (HMMs). The decrease in the error during training and recognizing   genders during testing phase close to 100 % accuracy demonstrates that the excitation component of speech contains gender-specific information and is indeed being effectively captured by continuous Ergodic HMM. A gender recognition study using gender specific features for different HMM states, mixture components,   size of testing data on the performance of the gender recognition is   evaluated. We demonstrate the gender recognition studies on TIMIT database.

**Keywords-** *Gender, Hidden Markov Model (HMM); LPC; MFCC*

## 1. INTRODUCTION

With the development of more and more identification systems to identity a person, there is a need for the development of a system which can provide personal identification task such as gender identification automatically without any human interface. Gender identification using voice of a person is comparatively easier than that from other approaches. There exist several algorithms for automatic gender identification but none of them has found to be 100% accurate.

In Gender identification based on the voice of a speaker consists of detecting if a speech signal is uttered by a male or a female. Automatically detecting the gender of a speaker has several potential applications. In the context of Automatic Speech Recognition, gender dependent models are more accurate than gender independent ones [1] [2]. Hence, gender recognition is needed prior to the of speaker recognition, gender dependent model. In the context of speaker recognition, gender detection can improve the performance by limiting the search space to speakers from the same gender. Also, in the context of content based multimedia indexing the speaker's gender is a cue used in the annotation. Therefore, automatic gender detection can be a tool in  a  content-based multimedia indexing system.

Much information can be inferred form a speech, such as sequences of words, gender, age, dialect, emotion, and even level of education, height or weight etc. Gender is an important characteristic of a speech. Automatically detecting the gender of a speaker has several potential applications such as (1) sorting telephone calls by gender (e.g. for gender sensitive surveys), (2) as part of an automatic speech recognition system to enhance speaker adaptation, and (3) as part of automatic speaker recognition systems. In the past, many methods of gender classification have been proposed. For parameters selections, some methods used gender dependent features such as pitch and formants [3] [4].

Speech is composite signal which has information about the message, gender, the speaker identity and the language [5][6]. It is difficult to isolate the speaker specific features alone from the signal.  The speaker characteristics present in the signal can be attributed to the anatomical and the behavioural aspects of the speech production mechanism.  The representation of the behavioural characteristics is a difficult task, and usually requires large amount of data.  Automatic speaker recognition systems rely mainly on features derived from the physiological characteristics of the speaker.

Speech is produced as sequence of sounds. Hence the state of vocal folds, shape and size of various articulators, change over time to reflect the sound being produced. To produce a particular sound the articulators have to be positioned in a particular way. When different speakers try to produce same sound, through their vocal tracts are positioned in a similar manner, the actual vocal tract shapers will be different due to differences in the anatomical structure of the vocal tract. System features represent the structure of vocal tract. The movements of vocal folds vary from one speaker to another. The manner and speed in which the vocal folds close also varies across speakers. Hence different voices are produced. Source features represent these variations in the vibrations of the vocal folds.

The theory of Linear Prediction (LP) is closely linked to modelling of the vocal tract system, and relies upon the fact that a particular speech sample may be predicted by a linear combination of previous samples. The number of previous samples used for prediction is known as the order of the prediction. The weights applied to each of the previous speech samples are known as Linear Prediction Coefficients (LPC). They are calculated so as to minimize the prediction error. As a byproduct of the LP analysis, reflection coefficients and log area coefficients are also obtained [7].

A study into the use of LPC for speaker recognition was carried out by Atal [8]. These coefficients are highly correlated, and the use of all prediction coefficients may not be necessary for speaker recognition task [9]. Sambur [10] used a method called orthogonal linear prediction.  It is shown that only a small subset of the resulting orthogonal coefficients exhibits significant variation over the duration of an utterance. It is also shown that reflection coefficients are as good as the other feature sets. Naik et. al., [11] used principal spectral components derived from linear prediction coefficients for speaker verification task. Hence a detailed exploration to know the speaker-specific excitation information

present in the residual of speech is needed and hence the motivation for the present work.

The rest of the paper is organized as follows: In Section II we examine the gender characteristics of the LP residual, and discuss issues involved in extracting the speaker-specific information from the residual. In Section III we discuss feature extraction using Mel Ceptral coefficients to capture the speaker specific information from the residual. Section IV describes Gaussian Mixture Model for Gender Recognition. Section V describes the database used in the study and Section V1 describes performance evaluation of Gender identification sytem. The proposed gender recognition system, based on the LP residual, may not require large amounts of data.

## 2. GENDER CHARACTERISTICS IN THE LP RESIDUAL

Speech signals, as any other real world signals, are produced by exciting a system with source. A simple block diagram representation of the speech production mechanism is shown in the Fig.1. Vibrations of the vocal folds, powered by air coming from the lungs during exhalation, are the sound source for speech. Hence, as can be from Fig. 1, the glottal excitation forms the source, and the vocal tract forms the system. One of the most powerful speech analysis technique is the method of linear predictive analysis. The philosophy of linear prediction is intimately related to the basic speech production model. The Linear Predictive Coding (LPC) analysis approach performs spectral analysis on short segments of speech with an all-pole modelling constraint [12]. Since speech can be modelled as the output of linear, time-varying system excited by a source, LPC analysis captures the vocal tract system information in terms of coefficients of the filter representing the vocal tract mechanism. Hence, analysis of speech signal by LP results in two components, namely the synthesis filter on one hand and the residual on the other hand. In brief, the LP residual signal is generated as a by product of the LPC analysis, and the computation of the residual signal is given below.
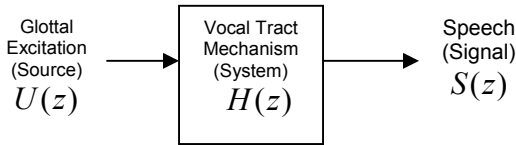


**Fig. 1: Source and System representation of speech production mechanism.**

If the input signal is represented by $u_n$ and the output signal by $s_n$, then the transfer of the system can be expressed as,

$$H(z) = \frac{S(z)}{U(z)} \quad (1)$$

Where $s(z)$ and $u(z)$ are z-transforms of $s_n$ and $u_n$ respectively.

Consider the case where we have output signal and the system and have to compute the input signal. The above equation can be expressed as $S(z) = H(z)U(z)$

$$U(z) = \frac{S(z)}{H(z)} \quad (2)$$

$$U(z) = \frac{1}{H(z)} S(z) \quad (3)$$

$$U(z) = A(z)S(z) \quad (4)$$

Where $A(z) = \frac{1}{H(z)}$ is the inverse filter representation of the vocal tract system.

Linear prediction models the output $s_n$ as the linear function of past outputs and present and past inputs. Since prediction is done by a linear function, the name linear prediction. Assuming an all-pole for the vocal tract, the signal $s_n$ can be expressed as linear combination of past values and some input $u_n$ as shown below.

$$Sn = - \sum_{k=1}^{p} akSn-k + GUn \quad (5)$$

Where G is a gain factor.

Now assuming that the input $u_n$ is unknown, the signal $s_n$ can be predicted only approximately from a linear weighted sum of past samples. Let this approximation of $s_n$ be , where

$$\widetilde{S}_n = - \sum_{k=1}^{p} a_k s_{n-k} \quad (6)$$

Then the error between the actual value Sn and predicted value is given by $en = Sn - \widetilde{S}n$ [13]. This error is nothing but LP residual of signal is shown in Fig 2.
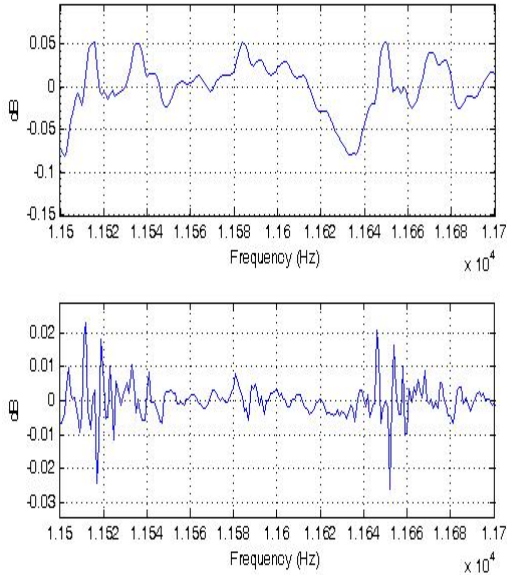
By logarithm of amplitude of mel spectrum and applying reverse Fourier transformation we achieve frame cepstrum:

$$mel - cepstrum(frame) = FFT^{-1}\big[mel(\log | FFT(frame)|)\big] \quad (9)$$

The FFT-base cepstral coefficients are computed by taking IFFT of the log magnitude spectrum of the Speech signal. The mel-warped cepstrum is obtained by inserting a intermediate step of transforming the frequency scale to place less emphasis on higher frequencies before taking the IFFT [7][15][16].

## 3. HIDDEN MARKOV MODEL FOR GENDER IDENTIFICATION

The Hidden Markov Models (HMM) is a doubly embedded stochastic process where the underlying stochastic process is not directly observable. HMMs can be used as probabilistic speaker models for both text-dependent and independent speaker recognition. An HMM not only models the underlying speech sounds but also the temporal sequencing of the sounds. This temporal modelling is advantageous for text-dependent tasks. For text-dependent speaker recognition task, HMM-based methods have achieved significantly better recognition[17][18][19].

Since the stressful cues contained in an utterance cannot be assumed as specific sequential events in the signal, an ergodic or fully connected HMM structure becomes more appropriate than LeftToRight (LTR) structure because every state in the ergodic structure can be reached in a single step from every other state. An ergodic or fully connected that is derived from ergodic or fully connected HMM has been used in this work. The transition matrix, A, of this structure can be written in terms of the bij coefficients (positive coefficients) as,

$$A = \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{bmatrix}$$

In training phase, an HMM for each speaker is obtained by estimating the parameters of model using feature vectors from the training data. The parameters of HMM are:
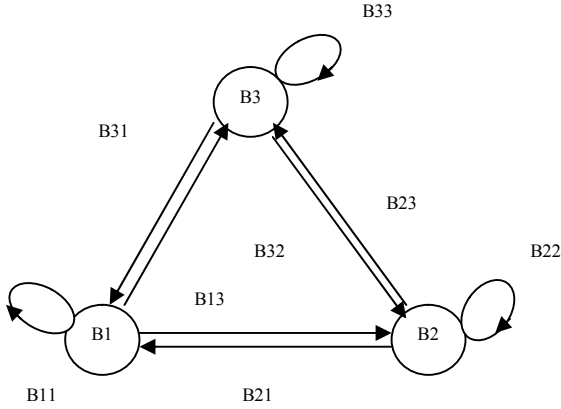


**Fig 2: Actual signal and its LP residual**

## Feature Extraction Of Lp Residual Signal

MFCC is the best known and most popular, and this features has been used for gender identification. MFCC's are based on the known variation of the human ear's critical bandwidths with frequency. The MFCC technique makes use of two types of filter, namely, linearly spaced filters and logarithmically spaced filters. To capture the phonetically important characteristics of speech, signal is expressed in the Mel frequency scale. This scale has a linear frequency spacing below 1000 Hz and a logarithmic spacing above 1000 Hz. Normal speech waveform may vary from time to time depending on the physical condition of speakers' vocal cord. Rather than the speech waveforms themselves, MFFCs are less susceptible to the said variations [6].

### 2.1.1 Motivation to use Melfrequency Cepstral coefficients(MFCC)

Since our interest is in capturing global features which correspond to source excitation, the low frequency or pitch components are to be emphasized. To fulfil this requirement it is felt that MFCC are most suitable as they emphasize low frequency and de-emphasize high frequencies.

### 2.1.2 MFCC

In this phase the digital speech signal is partitioning into segments (frames) with fixed length 10-30 ms from which the features are extracted due to their spectral qualities. Spectrum is achieved with fast Fourier transformation [14]. Then an arrangement of frequency range to mel scale follows according to relation

$$f_{mel} = 2595 \log\left(1 + \frac{f_{Hz}}{700}\right) \quad (8)$$

**Fig 3: Three-state ergodic HMM**

## Experimental Setup

The system has been implemented in Matlab7 on Windows XP platform. We have trained the HMMs using Gaussian Components as 2, 4, 8, and 16 by varying the HMM states from 2 to 4, and for training speech duration of 30 sec. Testing is performed using different test speech durations such as 1 sec., 2 sec., and 3 sec..

## 3.2 PERFORMANCE EVALAUATION

### 3.2.1 Gender Recognition Performance for varying Mixture Components

Determining the optimal number of mixture components needed to model a gender adequately is an important task. There is no theoretical way to estimate the number of mixture components to model a gender. The experiment is carried out for a 2-state HMM, 3-state-HMM and 4-state HMM for varying number of Gaussian components such as 2, 4, 8 and 16 for each state to evaluate the performance of the gender recognition. Here the model is trained with 30 seconds of speech duration. The system is tested with 1 second of test speech length and the performance is shown in the Fig. 4, Fig. 5 and Fig. 6.. The percentage of recognition for different Gaussian components such as 2, 4, 8, and 16 seems to be uniformly increasing. The minimum number of Gaussian components to achieve good recognition performance seems to be 4 and thereafter the recognition performance is minimal. The number of Gaussian components may vary for different experimental setups such as number of states, amount of training data, and amount of test data to achieve maximum recognition performance.

### 3.2.2 Gender Recognition Performance for varying Test Duration

In this experiment each gender model is tested with varying test durations such as 1 sec., 3 sec., and 5 sec.. The model is trained with 30 sec.. As shown in the Fig. 4 to Fig. 6. The recognition performance of the HMM drastically increases for the test speech length 1 sec. to 3 sec. Increasing the test length from 3 sec. to 5 sec. improve the performance with small improvement. This suggest that at least 3 sec. of test speech data is required to maintain better gender recognition performance. If there is no enough test speech data, the selection of the HMM states, number of Gaussian components on each state becomes more important to get good gender recognition performance

### 3.2.3 Average Gender Recognition Performance for varying HMM states

In this experiment the gender recognition performance is calculated for varying number of states from 2 to 4. The model is trained with number of Gaussian components as 4 on each state, training speech duration of 30 seconds. The model is tested with test duration of 1 sec., 3 sec., and 5 sec.. and the performance is shown in Fig. 7. For the above experimental setup to get good speaker recognition performance the minimum number of states required are 3. Above this minimal number of states the performance seems to be slightly increasing.

State transition probability distribution: It is represented by A=$\left[a_{ij}\right]$ where

$$aij = P(q_{t+1} = j \mid q_t = i) \quad 1 \le i, j \le N$$

defines the probability of transition from state i to j at time t.

Observation symbol probability distribution: It is given by $B = b_j(k)$, in which

$$b_j(k) = P(O_t = V_k \mid q_t = j) \quad 1 \le k \le M$$

defines the symbol distribution in state j, j=1, 2, … N

The initial state distribution: It is given by $\prod = \left[\pi\right]$, where

$$\pi_i = P(q_1 = i) \quad 1 \le i \le N$$

Here, N is the total number of states, and $q_t$ is the state at time t. M is the number of distinct observation symbols per state, and $O_t$ is the observation symbol at time t.

The model parameters can be collectively represented as $\lambda = \left\{A_i, B_i, \pi_i\right\}$ for i = 1, … M. Each speaker in a speaker identification system can be represented by a HMM and is referred to by the speaker's respective models λ.

In the testing phase, $P(O \mid \lambda)$ for each model is calculated [19], where $O = \left(o_1 o_2 o_3 \ldots o_T\right)$ is the sequence of the test feature vectors. The goal is to find the probability, given the model, that the test utterance belongs to that particular model. The speaker model that gives the highest score is declared as the identified speaker.

# 4. CONCLUSION

In this work we have demonstrated the importance of information in the excitation component of speech ( pitch ) for gender recognition task. Linear prediction residual is used to represent the excitation information. Performance of the recognition experiments shows that Ergodic HMM can capture speaker-specific excitation information from the LP residual. Performance of the system for different Mixture components shows that the optimal mixture components are 4 for speech signals sampled at 16 kHz. The recognition performance depends on the training speech length selected for training to capture the speaker-specific excitation information. Larger the training length, the better is the performance, although smaller number reduces computational complexity.

The objective in this paper was mainly to demonstrate the significance of the speaker-specific excitation information ( pitch) present in the linear prediction residual for gender recognition. We have not made any attempt to optimize the parameters of the model used for feature extraction, and also the decision making stage. Therefore the performance of speaker recognition may be improved by optimizing the various design parameters

# 5. REFERENCES

[1] Alex Acero and Xuedong Huang, Speaker and Gender Normalization for Continuous-Density Hidden Markov Models, in Proc. of the Int. Conf. on Acoustics, Speech, and Signal , IEEE, May 1996

[2] C. Neti and Salim Roukos. Phone-specific gender-dependent models for continuous speech recognition, Automatic Speech Recognition and Understanding Workshop (ASRU97), Santa Barbara, CA, 1997.

[3] R. Vergin, A. Farhat and D.O'Shaughnessy, "Robust gender-dependent acoustic-phonetic modeling in continuous speech recognition based on a new automatic male/female classification", Proc. Of IEEE Int. Conf. on Spoken Language (ICSLP), pp. 1081, Oct. 1996.

[4] S. Slomka and S. Sridharan, "Automatic gender identification optimized for language independence", Proc. Of IEEE TENCON'97, pp. 145-148,Dec. 1997.

[5] O'Shaughnessy, D., 1987. Speech Communication: Human and Machine. Addison-Wesley, New York.

[6] Rabiner, L.R., Juang, B.H., 1993. Fundamentals of Speech Recognition. Prentice-Hall, Englewood Cli s, NJ.

[7] Makhoul, J., 1975. Linear prediction: a tutorial review. Proc. IEEE 63, 561–580.

[8] B.S. Atal, "Effectiveness of linear prediction characteristics of the speech wave for automatic speaker identification and verification" J. Acoust. Soc. Ameri., vol. 55, pp.1304-1312, Jun. 1974.K. Elissa, "Title of paper if known," unpublished

[9] A.E. Rosenberg and M. Sambur, "New techniques for automatic speaker verification.", vol. 23, no.2, pp.169-175, 1975.

[10] M. R. Sambur, "Speaker recognition using orthogonal linear prediction," IEEE Trans. Acoust. Speech, Signal Processing, vol. 24, pp.283-289, Aug. 1976

[11] J. Naik and G. R. Doddington, " high performance speaker verification using principal spectral components", in proc.

IEEE Int. Conf. Acoust. Speech, Singal Processing, pp. 881-884, 1986.

[12] Furui, S., 1997. Recent advances in speaker recognition. Pattern Recognition Lett. 18, 859–872.

[13] S.R.Mahadeva Prassana, Cheedella S. Gupta, B. Yegnanarayana. Extraction of speaker-specific excitation information from linear prediction residual of speech. Speech Communications Vol.48 (2006) pp.1243-1261.

[14] Dempster, A., Laird, N., and Rubin, D., "Maximum likelihood from incomplete data via the EM algorithm," Journal of the Royal Statistical Society, vol. 39, pp. 1-38, 1977.

[15] Molau, S., Pitz, M., Schluter, R., and Ney, H., "Computing Mel-frequency cepstral coefficients on the power spectrum," Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP), vol. 1, pp. 73-76, May. 2001.

[16] Picone, J. W., "Signal modeling techniques in speech recognition," Proceedings of IEEE, vol. 81, no. 9, pp. 1215-1247, Sep. 1993.

[17] M. Forsyth and M. Jack, Discriminating semi-continuous HMM for speaker verification,‖ in proc. IEEE Int. Conf. Acoust. Speech, Signal Processing, vol.1, pp. 313-316, 1994.

[18] M. Forsyth, Discriminating observation probability (DOP) HMM for speaker verification, ‖ Speech Communicaiton, vol. 17, pp.117-129, 1995.

[19] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm", J. Royal Statist. Soc. Ser. B. (methodological), vol. 39, pp. 1-38, 1977

[20] K.N. Stevens, Acoustic Phonetics. Cambridge, England: The MIT Press, 1999

| No. of HMM states | No. of mixture components | Recognition rate (%) | | |
|---|---|---|---|---|
| | | Testing speech length | | |
| | | 1 sec. | 3 sec. | 5 sec. |
| 2 | 2 | 92 | 95 | 96 |
| | 4 | 96 | 98 | 99 |
| | 8 | 95 | 98 | 99 |
| | 16 | 94 | 97 | 96 |
| 3 | 2 | 96 | 98 | 99 |
| | 4 | 98 | 100 | 100 |
| | 8 | 97 | 98 | 98 |
| | 16 | 93 | 95 | 95 |
| 4 | 2 | 95 | 95 | 97 |
| | 4 | 98 | 99 | 99 |
| | 8 | 98 | 97 | 97 |
| | 16 | 95 | 96 | 97 |

**Table 1:Gender recognition performance with TIMIT database   ( MALE -100  +  FEMALE-100 )**
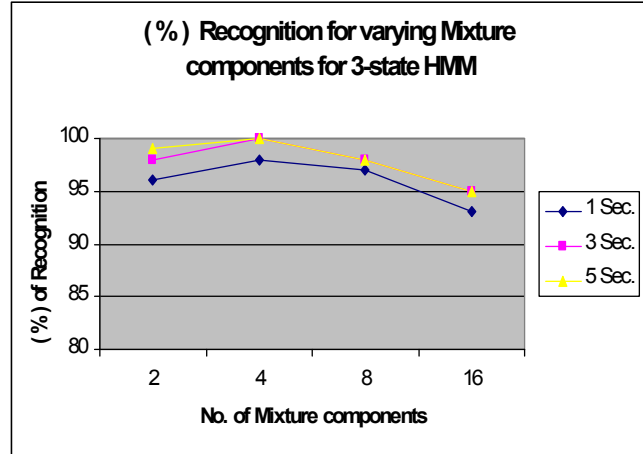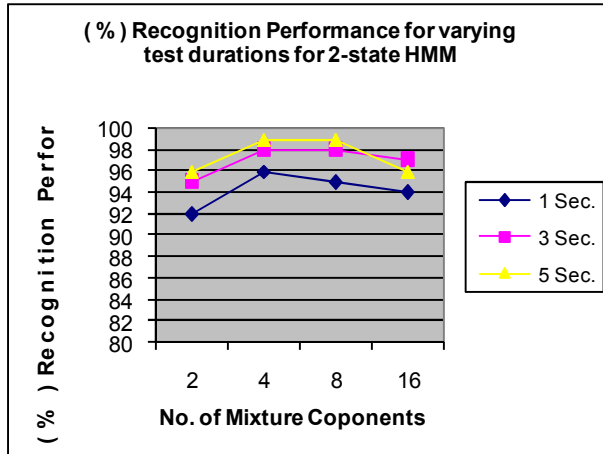
Fig. 4: Gender Recognition Performance for varying Mixtures



Fig. 5: Gender Recognition Performance for varying Train Duration



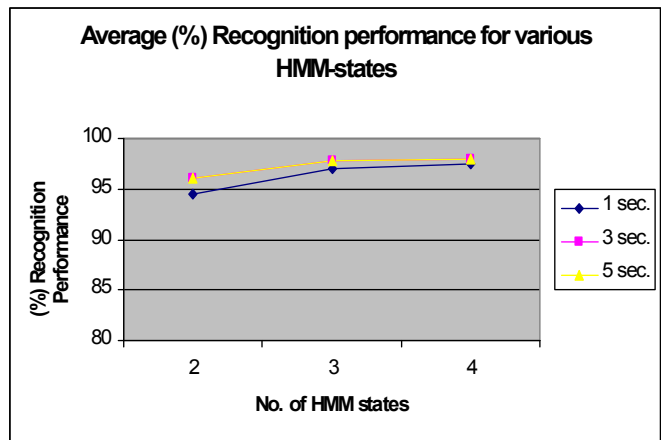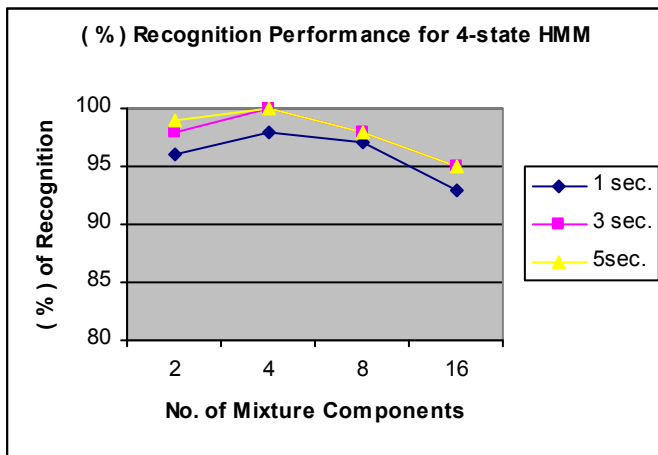Fig. 6: Gender Recognition Performance for 4 state HMM



Fig. 7: Gender ( %) Recognition Performance for varying HMM states