

# International Journal of Speech Technology

## A Nonlinear Autoregressive Model for Speaker Verification

--Manuscript Draft--

<b>Manuscript Number:</b>	IJST-D-13-00015R1
<b>Full Title:</b>	A Nonlinear Autoregressive Model for Speaker Verification
<b>Article Type:</b>	Manuscript
<b>Keywords:</b>	Gaussian mixture models; mixture autoregressive model; nonlinear statistical models; speaker verification
<b>Corresponding Author:</b>	Joseph Picone, Ph.D. Temple University Elkins Park, Pennsylvania UNITED STATES
<b>Corresponding Author Secondary Information:</b>	
<b>Corresponding Author's Institution:</b>	Temple University
<b>Corresponding Author's Secondary Institution:</b>	
<b>First Author:</b>	Sundararajan Srinivasan, PhD
<b>First Author Secondary Information:</b>	
<b>Order of Authors:</b>	Sundararajan Srinivasan, PhD
	Tao Ma, PhD
	Georgios Lazarou, PhD
	Joseph Picone, Ph.D.
<b>Order of Authors Secondary Information:</b>	
<b>Abstract:</b>	<p>Gaussian Mixture Models (GMM) have been the most popular approach in speaker recognition and verification for over two decades. The inefficiencies of this model for signals such as speech are well documented and include an inability to model temporal dependencies that result from nonlinearities in the speech signal. The resulting models are often complex and overdetermined, which leads to a lack of generalization. In this paper, we present a nonlinear mixture autoregressive model (MixAR) that attempts to directly model nonlinearities in the trajectories of the speech features. We apply this model to the problem of speaker verification. Experiments with synthetic data demonstrate the viability of the model. Evaluations on standard speech databases, including TIMIT, NTIMIT, and NIST-2001, demonstrate that MixAR, using only half the number of parameters and only static features, can achieve a lower equal error rate when compared to GMMs, particularly in the presence of previously unseen noise. Performance as a function of the duration of both the training and evaluation utterances is also analyzed.</p>

## General Comments:

In this paper, we present a nonlinear mixture autoregressive model (MixAR) that attempts to directly model nonlinearities in the trajectories of the speech features. Performance benefits are demonstrated across several databases including TIMIT, NTIMIT, and NIST-2001. We demonstrate that MixAR, using only half the number of parameters and only static features, can achieve a lower equal error rate when compared to GMMs, particularly in the presence of previously unseen noise. This is the first demonstration of improved performance on a significant speech recognition task for this approach, and demonstration of its efficacy across several databases makes the paper unique.

Below we address the specific concerns of the reviewers:

### Reviewer #2:

#### **FIGURE ONE IS INCOMPLETE (variables $w_i$ and $g_i$ )**

The relationship of  $w_i$  and  $g_i$  to the weights on the Gaussian mixtures is indirect. We have included an equation explaining this relationship, and redrawn the figure to explicitly show these dependencies. Hopefully it is clearer now.

### Reviewer #3:

**The authors present a nonlinear mixture autoregressive model and its application to the problem of speaker verification. The paper is well organized and the theoretical explanation is concise and technically sound. The authors conducted extensive experiments to illustrate the advantages of the presented model over the commonly used GMM model. However, in experiments the evaluation of the GMM model was not conducted on the state-of-the-art GMM system, what authors also mention in the summary. The authors should mention this fact already in the sections III and IV and should discuss the results emphasizing the performance limitations of the implemented GMM, which was not the state-of-the-art system. Therefore, further elaboration of the section III and especially IV is needed to adequately include this aspect in the discussion.**

We were careful not to use the term “state of the art” when describing GMMs. State of the art, of course, can mean many different things. The baseline system our work is based on uses the same number of mixture components per state. We understand not everyone does this today, but many practical systems are still constrained this way. We added some discussion of this at the beginning of Section III, and we reinforce this notion in the summary.

**The curves in the Figures 4, 5 and 6 are not readable (the curves cannot be distinguish) if the article is BW printed. To provide readable figures also in black/white layout of the article the authors may consider to change line types of the curves in Figures 4 - 6.**

We no longer have the ability to reproduce these figures. They are more legible in color than B&W. We will do our best to provide improved figures when the final color figures are requested.

**On page 3 last line the Section VVI should be changed to Section VI.**

Fixed.

**Summary Comments:**

We would like to thank the reviewers for their thoughtful and insightful feedback. We have done our best to address each of their concerns. We hope our responses have adequately addressed your concerns and we can proceed with publication.

## A NONLINEAR AUTOREGRESSIVE MODEL FOR SPEAKER VERIFICATION<sup>1</sup>

Sundararajan Srinivasan<sup>2</sup>, Tao Ma<sup>3</sup>, Georgios Lazarou<sup>4</sup> and Joseph Picone<sup>5</sup>

*Abstract*— Gaussian Mixture Models (GMM) have been the most popular approach in speaker recognition and verification for over two decades. The inefficiencies of this model for signals such as speech are well documented and include an inability to model temporal dependencies that result from nonlinearities in the speech signal. The resulting models are often complex and overdetermined, which leads to a lack of generalization. In this paper, we present a nonlinear mixture autoregressive model (MixAR) that attempts to directly model nonlinearities in the trajectories of the speech features. We apply this model to the problem of speaker verification. Experiments with synthetic data demonstrate the viability of the model. Evaluations on standard speech databases, including TIMIT, NTIMIT, and NIST-2001, demonstrate that MixAR, using only half the number of parameters and only static features, can achieve a lower equal error rate when compared to GMMs, particularly in the presence of previously unseen noise. Performance as a function of the duration of both the training and evaluation utterances is also analyzed.

*Keywords*— Gaussian mixture models, mixture autoregressive model, nonlinear statistical models, speaker verification

Manuscript submitted February 15, 2013.

1. This material is based upon work supported by the National Science Foundation under Grant No. IIS-0414450. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.
2. S. Srinivasan is with Nuance Communications Inc., 1198 East Arques Avenue Sunnyvale, CA 94085, USA (phone: 408-992-6243; email: sundararajan.srinivasan@gmail.com).
3. T. Ma is with Siri at Apple Inc., 2 Infinite Loop, mailstop 302-4APP, Cupertino, California 95014, USA (phone: 408-643-5909; email: tma@apple.com).
4. G. Lazarou is with The New York City Transit Authority, 30-74 38th Street, Apt 1A, Astoria, New York, New York, USA 11103 (phone: (662) 617-2064; email: glaz@ieee.org).
5. J. Picone is with the Department of Electrical and Computer Engineering at Temple University, 1947 North 12<sup>th</sup> Street, Philadelphia, Pennsylvania 19027 USA (phone: 215-204-4841; fax: 215-204-5960; email: joseph.picone@isip.piconepress.com).

## I. INTRODUCTION

The goal in speaker verification is to accept or reject an identity claim made by a speaker. This biometric is widely used in a variety of applications ranging from secured access and surveillance to multimodal verification. A challenge for statistical modeling in speaker verification is to accurately and efficiently represent the probability distribution of speaker features so that even similar sounding speakers can be distinguished. The majority of speaker recognition and verification systems today utilize Gaussian Mixture Models (GMMs) either entirely or as part of a hybrid model (Beigi, 2011).

There are two well-known drawbacks of the GMM model. The first involves statistical independence – there are obviously dependencies in the temporal evolution of both the static features and the derivatives of these features. Constructing a GMM from standard features decorrelated using only a diagonal covariance matrix does not adequately model these dependencies. Use of full covariance matrices results in models with an extremely large number of parameters and creates parameter estimation problems. Performance improvements with full covariance approaches (or constrained versions of these large matrices) have been minimal, and often increase the system’s sensitivity to mismatched channel conditions. A major overarching goal in our work is to improve performance when the training and evaluation conditions are mismatched, and the evaluation data contains previously unseen noise and channel conditions.

The second more serious drawback, which is the focus of this work, is the implicit assumption of linearity in the feature vector dynamics. The derivatives of the cepstral features, commonly used in speech processing as part of the MFCC feature vector representation (Chen & Bilmes, 2007) are only a linear approximation of the actual dynamics of the static features. However, a survey of studies on the subject shows that the speech signal contains significant nonlinear information, and using only derivative features to represent speech dynamics with GMM modeling is tantamount to discarding any nonlinear information present in the signal (May, 2008; Kokkinos & Maragos, 2005).

An obvious solution to this problem is to add features that can represent the nonlinear dynamic information. However, adding nonlinear invariants as features has not improved the robustness of

1  
2  
3  
4 recognition and verification technologies in harsh or mismatched environments. The reasons for these  
5 failures include (1) it is difficult to estimate invariants reliably from speech, resulting in parameter  
6 estimation algorithms that need to be extensively tuned; (2) these estimation algorithms typically require  
7 an acoustic event to have a long duration (Petry et al., 2002), and this gravely undermines the  
8 applicability of invariant features for a short-term stationary signal like speech; and (3) invariants only  
9 quantify the degree of nonlinearity and do not characterize the nature of the dynamics completely.

10  
11 The primary goal of this work is to approach the information representation problem at the acoustic  
12 modeling level using a nonlinear mixture autoregressive model (MixAR) (Zeevi et al., 2000), thereby  
13 accounting for the nonlinear dynamics of speech in the base model and minimizing the dimensionality of  
14 the feature space. This model is shown in Figure 1. Previous work on mixture autoregressive modeling  
15 for speech has been in the context of hidden Markov models for speech recognition (Juang & Rabiner,  
16 1985). A more recent investigation of AR-HMMs (Ephraim & Roberts, 2005) used a switching  
17 autoregressive process to capture signal correlations during state transitions. Another model considered  
18 speech features as a GMM white noise process filtered through an autoregressive signal for speaker  
19 identification (Ayadi, 2008). Results on speech recognition showed that at best these models were only  
20 comparable to an MFCC-based HMM using a GMM observation model.

21  
22 A more sophisticated model (Wong & Li, 2000) considers a mixture of autoregressive filters (MAR)  
23 for the observation model. Our earlier work (Srinivasan et al., 2008) considered this model for phone  
24 classification. MixAR is a generalization of MAR, where the mixture weights are allowed to be  
25 time-varying and data-dependent. In this work, we apply the MixAR model to feature vectors in a speaker  
26 verification task, and demonstrate improved performance over the more limited MAR model.

27  
28 The rest of the paper is organized as follows: Section II formally defines the MixAR model, explains  
29 some of the relevant properties of this model and discusses the parameter estimation problem. Results of  
30 experiments using synthetic data are included in Section III and speaker verification experiments with real  
31 speech data are presented in Section IV. Experiments documenting variation in performance with the  
32 duration of training and evaluation utterances are also discussed in Section IV. Finally, in Section VVI we

present our conclusions and discuss future directions for this research.

## II. THE MIXAR MODEL

A mixture autoregressive process (MixAR) of order  $p$  with  $m$  components,  $X=\{x[n]\}$ , is defined as (Zeevi et al., 2000):

$$x[n] = \begin{cases} a_{1,0} + \sum_{i=1}^p a_{1,i} x[n-i] + \varepsilon_1[n] & \text{w.p. } W_1(x[n-1]) \\ a_{2,0} + \sum_{i=1}^p a_{2,i} x[n-i] + \varepsilon_2[n] & \text{w.p. } W_2(x[n-1]) \\ \vdots \\ a_{m,0} + \sum_{i=1}^p a_{m,i} x[n-i] + \varepsilon_m[n] & \text{w.p. } W_m(x[n-1]) \end{cases} \quad (1)$$

where  $\varepsilon_i$  is a zero-mean Gaussian random process with a variance of  $\sigma_j^2$ , “w.p.” denotes “with probability” and the gating weights,  $W_i$ , sum to 1. The linear prediction coefficients,  $\{a_i\}$ , represent the dynamic model, where  $a_{i,0}$  are the component means. The coefficients  $w_i$  and  $g_i$  are called gating coefficients, and are defined by the following relation:

$$W_i(x) = \frac{e^{w_i + g_i x}}{\sum_{j=1}^m e^{w_j + g_j x}} \quad (2)$$

It is apparent that an  $m$ -mixture MixAR process is the weighted sum of  $m$  Gaussian autoregressive processes, with the time-dependent weights depending on previous data and the gating coefficients.

One insightful way of viewing this model is as a process in which each data sample at any one point in time is generated from one of the component AR mixture processes chosen randomly according to its weight  $W_i$ , as depicted in Figure 1. One property of MixAR that is of particular relevance here is the ability of MixAR to model nonlinear time series (Zeevi et al., 2000; Wong & Li, 2000). Though the individual component AR processes are linear, the probabilistic mixing of these AR processes constitutes a nonlinear model. Even when the mixture weights are fixed, the model reduces to MAR, which is still

nonlinear. The addition of a gating system layer for weight generation increases the flexibility of the model even further, allowing us to model distributions as a function of past data.

Several other properties of MixAR, including a mathematically rigorous proof of the asymptotic performance of a MixAR model for stochastic processes, are derived in (Zeevi et al., 2000). Note that in the original formulation, both the gate and prediction orders were constrained to be equal. There are some practical implementation issues regarding parameter estimation for this model, and these are discussed next. In this paper, we restrict ourselves to MixAR models of order one to avoid some difficulties with parameter estimation.

#### A. Estimation of the Prediction and Variance Parameters

Similar to the well-known training procedure for GMM, maximum likelihood estimates for MixAR prediction and variance parameters can be calculated using the Expectation-Maximization (EM) algorithm (Dempster et al., 1977). Given an order,  $p$ , the parameter set for each of the  $m$  components of a MAR model consists of  $p+1$  predictor coefficients (including the mean), the error variance, and mixing weight, and is denoted as:

$$\theta_l = \{a_{l,0}, a_{l,1}, \dots, a_{l,p}, \sigma_l, w_l, g_l\}, \quad l = 1, 2, \dots, m. \quad (3)$$

To estimate these parameters, we first need an initial guess for these parameters and then we iterate with EM to successively refine the estimates.

An initialization strategy that we found to work reasonably well was to first train a GMM with the same number of mixtures and then set each component of the MixAR model to have the same mean, variance and weight as the GMM model. We initialize the predictor coefficients and the data-dependency gating coefficients,  $\{A_i\}$ , to zero. These initial parameters can be then refined recursively using an E-step (Zeevi et al., 2000; Wong & Li, 2000):

$$\gamma_l[n] = \frac{W_l p_l(x[n]|\theta)}{\sum_{k=1}^m W_k p_k(x[n]|\theta)}, \quad (4)$$



where

$$p_l(x[n]|\theta) \propto \frac{1}{\sigma_l} e^{-\frac{1}{2\sigma_l^2} (x[n] - a_{l,0} - \sum_{i=1}^m a_{l,i} x[n-i])^2} \quad (5)$$

represents the probability that a sample was generated from component  $l$  at time instant  $n$ .

The corresponding M-step is given by:

$$\hat{A}_l = R_l^{-1} r_l \quad (6)$$

$$\hat{\sigma}_l^2 = \frac{\sum_{n=p+1}^N \gamma_l[n] \left( x[n] - \hat{a}_{l,0} - \sum_{i=1}^m \hat{a}_{l,i} x[n-i] \right)^2}{\sum_{n=p+1}^N \gamma_l[n]} \quad (7)$$

$$R_l = \sum_{n=p+1}^N \gamma_l[n] X_{n-1} X_{n-1}^T \quad (8)$$

$$r_l = \sum_{n=p+1}^N \gamma_l[n] X_{n-1} x[n] \quad (9)$$

$$X_{n-1} = \begin{bmatrix} 1 \\ x[n-1] \\ x[n-2] \\ \vdots \\ x[n-p] \end{bmatrix}. \quad (10)$$

### B. Estimation of the Gating Coefficients

A complication arises with respect to the estimation of gating coefficients for MixAR. There is no closed-form solution for these, and hence a Newton gradient-ascent approach must be used:

$$\hat{w}_l = w_l + \beta \frac{\Delta Q}{\Delta w_l} \quad (11)$$

$$\hat{g}_l = g_l + \beta \frac{\Delta Q}{\Delta g_l}, \quad (12)$$

where  $Q$  denotes the log-likelihood of the MixAR model for the training data, and  $\beta$  and  $\Delta$  are design

parameters to be chosen empirically. The expression for computing Q is:

$$Q(\theta) = \sum_{n=1}^N \sum_{l=1}^m \gamma_l[n] \log(W_l[n]) + \sum_{n=1}^N \sum_{l=1}^m \gamma_l[n] \log(p_l(x[n]|\theta)) . \quad (13)$$

Due to this complication in the updates for the gate coefficients, the training procedure outlined above is not in the realm of strict EM algorithm but falls under a class of algorithms known as generalized EM algorithms (GEM) (McLachlan et al., 2008). For both EM and GEM algorithms, the E-step is similar. However, while an EM algorithm actually maximizes the expectation during each M-step, a GEM algorithm only guarantees that parameters that the model likelihood for the data is increased but does not guarantee that it is maximized at each M-step. This has the potential to increase the number of iterations required for training compared to a pure EM algorithm.

In a process analogous to the choice of an adaptation constant in adaptive filter theory, we can postulate that quick and smooth convergence of the GEM algorithm can be achieved by starting with a relatively high value for  $\beta$  and then reducing this value with successive iterations. In our experiments, we found that fixing  $\Delta = 0.01$  and running 10 iterations each with  $\beta = 0.9$ ,  $\beta = 0.5$ , and  $\beta = 0.2$  in succession provided a reasonably smooth and quick convergence. However, such convergence is not guaranteed in general and this poses a generalization problem for wider application of this model.

Fortunately, we can do better than guessing an appropriate value for  $\beta$ . We can use the secant method for root-finding and maximization (Dennis & Schnabel, 1996). The iterative formula for finding the maximum using Newton's method is:

$$\hat{x} = x + f'(x)/f''(x) . \quad (14)$$

In the secant method, the double derivative in the denominator is estimated numerically using the secant at the point. We estimate the scaling factor  $\beta$  as the inverse of double derivative of the log-likelihood with respect to the gate parameters:

$$\beta = 1 / (\Delta^2 Q / \Delta^2 w_l) . \quad (15)$$

During implementation, this scheme amounts to finding the value of Q at three different points ( $Q(w_l)$ ),

$Q(w_l + \Delta), Q(w_l - \Delta)$ ) for each gate coefficient  $w_l$  and then using the following update equation:

$$\hat{w}_l = w_l + \frac{Q(w_l + \Delta) - Q(w_l - \Delta)}{Q(w_l + \Delta) + Q(w_l - \Delta) - 2Q(w_l)}. \quad (16)$$

Similarly, the update equation for gate coefficients  $g_l$  is:

$$\hat{g}_l = g_l + \frac{Q(g_l + \Delta) - Q(g_l - \Delta)}{Q(g_l + \Delta) + Q(g_l - \Delta) - 2Q(g_l)}. \quad (17)$$

Using this method, we obtain reasonable convergence as shown in Figure 2. We have used three GEM iterations in the experiments described in the following sections.

### III. SYNTHETIC DATA EXPERIMENTS

In the experiments described in the next two sections, we used the ISIP public domain speech recognition software (Huang & Picone, 2002) to implement the MixAR model as well as integrate it into an existing speaker verification system. Our baseline system uses an HMM with Gaussian mixture models for the observation probabilities. We chose to use the same number of Gaussian mixtures per state, and used the standard mixture-splitting approach to generating the mixture components. Though this is not necessarily state of the art, it is sufficiently close to state of the art for the experiments described below.

#### A. Two-Way Classification of Synthetic Speech-Like Data

In this section we describe a set of pilot experiments designed to validate the basic properties of the MixAR model. We first selected two speakers from the 2001 NIST SRE Corpus (Greenberg & Martin, 2009) and generated the static features corresponding to a standard MFCC feature vector (12 mel-cestral coefficients). A 3-state HMM and a MixAR model, both with 4 Gaussian mixture components per state, were trained for each speaker. For each class (e.g., a speaker), two speech-like signals of 40,000 vectors were generated from these models – a linear speech-like signal ( $X_1$ ) was synthesized from the HMM model, and a nonlinear speech-like signal ( $X_2$ ) was generated from the MixAR model. To simulate a range of signals with varying degrees of nonlinearity, the two signals were mixed with a mixing coefficient  $\alpha$ :

$$X_{\alpha} = (1-\alpha)X_1 + \alpha X_2. \quad (18)$$

The first 20,000 vectors from each  $X_{\alpha}$  were used as a training set while the remaining vectors were split into 200 segments of 100 vectors each for evaluation.

The results of the classification experiments with this data are shown in Table 1. Since the data contains nonlinearities, we compared MixAR using static features to a GMM using static features plus first derivatives (static+ $\Delta$ ). Performance for GMMs with only static features was significantly worse except for the case of a purely linear signal ( $\alpha = 0$ ). We also allocated 8 mixture components for the GMM system since its feature vectors had additional components. The MixAR system had only 4 mixture components. Overall, the MixAR system had approximately 50% fewer parameters than the GMM system. We can see that when the amount of nonlinearity is insignificant, GMM performs as well as MixAR. However, as the amount of nonlinearity in the signal increases, MixAR performs significantly better. These results validated the basic model and provided motivation to do further testing on more realistic data.

### *B. Speaker Verification With Synthetic Data*

Since our goal is to study speaker verification, we next used the development database from the 1-speaker detection task of the 2001 NIST SRE Corpus (Greenberg & Martin, 2009). Though the use of development database is not standard practice in most published baselines, it is small enough to allow us to quickly generate results using only modest computing power and yet large enough to provide a reliable estimate of the performance. All 60 speakers in the training set were used. Each training utterance was about 2 minutes long. Static (12 MFCCs + energy), delta (26 MFCCs) and delta-delta (39 MFCCs) features were extracted.

Two types of clean data were synthesized. For the first type, a 10-state HMM with 4-Gaussians per state was trained for each utterance. For the second type, a 32-mixture MixAR model of prediction order 1 was trained for each utterance. For each of the models trained, new training data of about 30,000 frames per speaker and evaluation data of 20 utterances with about 200 frames for each utterance per

1  
2  
3  
4 speaker were generated according to that model.  
5

6 Similarly, two types of noisy data were generated. First, the clean training utterances from the  
7 development data were corrupted with car noise to achieve an SNR of 5 dB. This approach followed a  
8 methodology previously used to generate the Aurora database (Parihar et al., 2004). The remaining steps  
9 to yield the two types of noisy data were the same as those for the clean case. The goal of creating data in  
10 this way was to simulate 4 different test conditions: clean+linear, clean+nonlinear, noisy+linear and  
11 noisy+nonlinear.  
12  
13  
14  
15  
16  
17  
18  
19

20 Using the synthetic training data, both GMMs and prediction order-1 MixAR models were trained for  
21 each speaker under each condition. Next, the corresponding synthesized evaluation data were used for  
22 evaluating speaker verification performance. For the clean case, there was little difference in performance  
23 between GMM and MixAR. For noisy evaluation data at 5 dB SNR, there was not much variation in  
24 performance between GMM and MixAR for HMM-generated data.  
25  
26  
27  
28  
29  
30

31 However, for the data generated from the nonlinear MixAR model and with the addition of noise, the  
32 MixAR model showed a significant improvement in performance using far fewer parameters. This is  
33 evident from the DET plot shown in Figure 3. These results provide support to the hypothesis that when  
34 there are significant nonlinearities in the signal, using this information makes the nonlinear model much  
35 more robust in the presence of noise.  
36  
37  
38  
39  
40  
41

#### 42 **IV. SPEAKER VERIFICATION EXPERIMENTS**

43

##### 44 *A. Evaluation Using the NIST-2001 Development Database*

45

46 We applied the MixAR model to the 1-speaker detection task in the 2001 NIST SRE Corpus. All 60  
47 speakers in the development set were used for training and all 78 utterances were used for evaluation.  
48 Each training utterance was about 2 minutes long, while the durations of the test utterances varied but did  
49 not exceed 60 seconds. A standard 39-dimensional MFCC feature vector was used.  
50  
51  
52  
53  
54

55 Performance was evaluated with and without delta features and energy for a fixed number of  
56 mixtures. The results are tabulated in Table 2. For GMM, a substantial improvement is obtained using the  
57 delta features and marginal improvements were obtained using delta-delta features. For MixAR, the use  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 of any delta features provided no measurable improvements. This clearly indicates that MixAR can  
5  
6 extract all necessary information from only the static features.  
7

8  
9 MixAR and GMM performance was then evaluated as a function of the number mixtures. The  
10  
11 detection error trade-off (DET) curves are shown in Figure 4. The EER results are shown in Table 3. Also  
12  
13 indicated in parenthesis is the number of parameters for each case. From this table it is clear that MixAR  
14  
15 can achieve about the same performance using 2x fewer parameters than GMM. This reduction in the  
16  
17 number of parameters points to the efficiency of MixAR in capturing the dynamic information.  
18

19  
20 Moreover, even when considering the best-case scenario for GMM with a large number of parameters  
21  
22 (8 mixtures with static as well as delta and delta-delta coefficients), there is a 10.6% relative reduction in  
23  
24 EER with MixAR. These results appear to strongly indicate that there is nonlinear evolution information  
25  
26 in speech features that the GMM model cannot capture using linear derivatives alone and that MixAR can  
27  
28 effectively employ this information for achieving better speaker verification.  
29

### 30 31 *B. Evaluation on Unseen Noise Conditions*

32  
33 To evaluate the robustness of MixAR compared to GMM on unseen noise conditions, several noise  
34  
35 conditions were simulated with the TIMIT database (Garofolo et al., 1993) by adding synthesized noise  
36  
37 from three different noise sources: white, car, and babble. Three SNR levels were used: 10, 5 and 0 dB (in  
38  
39 addition to the clean set). The core test partition of the database containing 168 speakers was used. The  
40  
41 three types of noise sources were chosen to represent commonly occurring types of noise.  
42  
43

44  
45 The matrix of experimental results is shown in Table 4. From this table, it is clear that while unseen  
46  
47 noise conditions degrades performance for both models, MixAR performs relatively better than GMM  
48  
49 and also uses 2.5x fewer parameters. The DET curves shown for the different noise conditions in Figure 5  
50  
51 also support the conclusion that MixAR performs better than GMM.  
52

53  
54 A related problem to unseen noise conditions is variation in acoustic channel. NTIMIT is a database  
55  
56 that was created by transmitting TIMIT utterances over different analog telephone channels (Jankowski et  
57  
58 al., 1990). We studied speaker verification performance on NTIMIT by splitting the data for each speaker  
59  
60 into 8 utterances for training and the remaining 2 utterances for evaluation. The DET performance curves  
61  
62  
63  
64  
65

1  
2  
3  
4 for the 8-mixture MixAR using only static MFCCs (with 480 parameters) and for the 16-mixture GMM  
5  
6 (with 1168 parameters) using both static and delta features are shown in Figure 6. The corresponding  
7  
8 EERs are shown in Table 5. From this analysis it is clear that MixAR using 2.5x fewer parameters  
9  
10 achieves the same or a higher level of performance as a GMM.  
11  
12

### 13 *C. Effect of Utterance Duration on Speaker Verification Performance*

14

15 Even if MixAR could do better under the conditions we have tested so far, it is possible that MixAR  
16  
17 requires more training data than GMM for reliable parameter estimation. This could be a particular  
18  
19 concern considering that MixAR attempts to learn nonlinear dynamic information, and nonlinear  
20  
21 dynamics are notoriously difficult to characterize from short lengths of data. For example, it is known that  
22  
23 estimates of Lyapunov exponents can be unreliable when the length of data is short (Banbrook et al.,  
24  
25 1997). One particular concern with insufficient training data is the problem of overfitting. It is therefore  
26  
27 necessary to study performance as a function of the amount of training data.  
28  
29  
30

31 Towards this end, we conducted experiments with varying training utterance durations keeping the  
32  
33 evaluation utterance duration a constant. Utterances corresponding to five durations – 120, 90, 60, 30 and  
34  
35 15s – were extracted from training data for each of the 60 speakers from the training part of the NIST  
36  
37 2001 development database. All evaluation data for the 78 speakers with durations ranging mostly  
38  
39 between 20 and 40s were used. The NIST-2001 database is particularly suited here because the training  
40  
41 data is clean and the evaluation data is corrupted by different kinds of noise. This means that models that  
42  
43 are overtrained will perform poorly on the evaluation utterances.  
44  
45

46 The number of mixtures for MixAR was fixed at 8. The number of mixtures for GMM was also fixed  
47  
48 at 8 to alleviate the problem of overfitting. The results of the experiment are reported in Table 6. There is  
49  
50 a 43.9% increase in EER for GMM when the training utterance duration is reduced from 120s to 15s. On  
51  
52 the other hand, the corresponding increase in EER for MixAR is only 26.56%. Thus, MixAR does not  
53  
54 necessarily require longer training utterance durations than GMM. Utterance duration is an important  
55  
56 practical constraint on voice biometrics.  
57  
58

59 It is also well known that increasing the evaluation utterance duration improves speaker verification  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 performance. We tested this hypothesis by examining evaluation utterances of five different durations –  
5  
6 30, 15, 10, 5, and 3s. These were extracted from each of the 78 test utterances in the NIST-2001  
7  
8 development database. All training data from all 60 speakers were used. The results of this experiment are  
9  
10 reported in Table 7. For GMM, there is an increase in EER of 31.2% as the evaluation utterance duration  
11  
12 reduces from 30s to 3s. The corresponding reduction for MixAR is 33.3%. Hence, both systems are  
13  
14 equally sensitive to the duration of the evaluation data.  
15  
16

## 17 **V. SUMMARY**

18  
19  
20 In this work, we have applied a nonlinear mixture autoregressive model (MixAR) to several speaker  
21  
22 verification tasks. Our experiments with synthetic as well as real speech data show that the MixAR model  
23  
24 outperforms GMM under several noise conditions, particularly the case where the type of noise in the  
25  
26 evaluation data was not observed in the training database. Equally important, MixAR did not require delta  
27  
28 features and used 2.5x fewer parameters to achieve comparable or better performance as that of GMM.  
29  
30 The dynamic modeling capability of MixAR is effective at capturing and exploiting speech dynamics  
31  
32 better than GMM.  
33  
34

35  
36 This work is part of our ongoing interest in exploiting nonlinear statistical models in speech  
37  
38 recognition. It has only begun to explore the potential of the nonlinear MixAR model in speaker  
39  
40 verification. Most modern systems incorporate more advanced adapted GMMs (Reynolds & Campbell,  
41  
42 2008) or discriminative model approaches (Li & Kinnunen, 2010). Before an unbiased comparison can be  
43  
44 made to these systems, a framework for integrating these approaches into MixAR model must be  
45  
46 developed. We also plan to apply the MixAR model to large vocabulary speech recognition tasks.  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



## VI. REFERENCES

- 1  
2  
3  
4  
5  
6  
7 Ayadi, M. (2008). *Autoregressive models for text independent speaker identification in noisy environments*. University of Waterloo.
- 8  
9  
10 Banbrook, M., Ushaw, G., & McLaughlin, S. (1997). How to extract Lyapunov exponents from short and  
11 noisy time series. *IEEE Transactions on Signal Processing*, 45(5), 1378–1382.
- 12  
13 Beigi, H. (2011). *Fundamentals of Speaker Recognition* (p. 942). Upper Saddle River, New Jersey, USA:  
14 Springer.
- 15  
16 Chen, C.-P., & Bilmes, J. A. (2007). MVA Processing of Speech Features. *IEEE Transactions on Audio,  
17 Speech and Language Processing*, 15(1), 257–270.
- 18  
19 Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum Likelihood from Incomplete Data via  
20 the EM Algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1), 1–38.
- 21  
22  
23 Dennis, J., & Schnabel, R. (1996). *Numerical Methods for Unconstrained Optimization and Nonlinear  
24 Equations* (p. 394). Englewood Cliffs, New Jersey, USA: Prentice Hall.
- 25  
26 Ephraim, Y., & Roberts, W. (2005). Revisiting autoregressive hidden Markov modeling of speech  
27 signals. *IEEE Signal Processing Letters*, 12(2), 166–169.
- 28  
29  
30 Garofolo, J., Lamel, L., Fisher, W., Fiscus, J., Pallet, D., Dahlgren, N., & Zue, V. (1993). TIMIT  
31 Acoustic-Phonetic Continuous Speech Corpus. *The Linguistic Data Consortium Catalog*.  
32 Philadelphia, Pennsylvania, USA: The Linguistic Data Consortium.
- 33  
34 Greenberg, C. S., & Martin, A. F. (2009). NIST speaker recognition evaluations 1996-2008. *Proceedings  
35 of SPIE (Stereoscopic Displays and Applications XX)* (p. 732411). Orlando, Florida, USA.
- 36  
37 Huang, K., & Picone, J. (2002). Internet-Accessible Speech Recognition Technology. *Proceedings of the  
38 IEEE Midwest Symposium on Circuits and Systems* (pp. III–73 – III–76). Tulsa, Oklahoma, USA.
- 39  
40 Jankowski, C., Kalyanswamy, A., Basson, S., & Spitz, J. (1990). NTIMIT: a phonetically balanced,  
41 continuous speech, telephone bandwidth speech database. *IEEE International Conference on  
42 Acoustics Speech and Signal Processing* (pp. 109–112 vol.1). Albuquerque, New Mexico, USA.
- 43  
44 Juang, B.-H., & Rabiner, L. (1985). Mixture autoregressive hidden Markov models for speech signals.  
45 *IEEE Transactions on Acoustics, Speech and Signal Processing*, 33(6), 1404–1413.
- 46  
47  
48 Kinnunen, T., & Li, H. (2010). An overview of text-independent speaker recognition: From features to  
49 supervectors. *Speech Communication*, 52(1), 12–40.
- 50  
51 Kokkinos, I., & Maragos, P. (2005). Nonlinear speech analysis using models for chaotic systems. *IEEE  
52 Transactions on Speech and Audio Processing*, 13(6), 1098–1109.
- 53  
54  
55 McLachlan, Geoffrey, & Thriyambakam, K. (2008). *The EM Algorithm and Extensions* (p. 400).  
56 Hoboken, New Jersey, USA: Wiley-Interscience.
- 57  
58 May, D. (2008). *Nonlinear dynamic invariants for continuous speech recognition*. Mississippi State  
59 University.
- 60  
61  
62  
63  
64  
65

- 1  
2  
3  
4 Parihar, N., Picone, J., Pearce, D., & Hirsch, H.-G. (2004). Performance Analysis of the Aurora Large  
5 Vocabulary Baseline System. *Proceedings of the European Signal Processing Conference* (pp. 553–  
6 556). Vienna, Austria.  
7  
8  
9 Petry, A., Augusto, D., & Barone, C. (2002). Speaker Identification using Nonlinear Dynamical Features.  
10 *Chaos, Solitons and Fractals*, 13(2), 221–231.  
11  
12 Reynolds, D., & Campbell, W. (2008). Text-Independent Speaker Recognition. *Springer Handbook of*  
13 *Speech Processing* (1st ed., p. 1176). Berlin, Germany: Springer.  
14  
15 Srinivasan, S., Ma, T., May, D., Lazarou, G., & Picone, J. (2008). Nonlinear Mixture Autoregressive  
16 Hidden Markov Models For Speech Recognition. *Proceedings of the International Conference on*  
17 *Spoken Language Processing* (pp. 960–963). Brisbane, Australia.  
18  
19 Zeevi, A., Meir, R., & Adler, R. (2000). *Nonlinear Models for Time Series using Mixtures of*  
20 *Autoregressive Models* (p. 25). Haifa, Israel. Retrieved from  
21 <http://ie.technion.ac.il/~radler/mixar.pdf>.  
22  
23  
24 Wong, C. S., & Li, W. K. (2000). On a Mixture Autoregressive Model. *Journal of the Royal Statistical*  
25 *Society. Series B (Statistical Methodology)*, 62(1), 95–115.  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## VII. LIST OF FIGURES

Figure 1. An overview of the (a) GMM and (b) MixAR approaches. The MixAR model is a weighted sum of Gaussian autoregressive models with time-dependent weights.

Figure 2. Performance of (Generalized) EM using the secant method as a function of the number of iterations for an 8-mixture MixAR model is shown (speaker 4516 from the NIST-2001 database).

Figure 3. DET curves are shown for a simulated speaker verification task. MixAR performance in the presence of noise exceeds GMM performance.

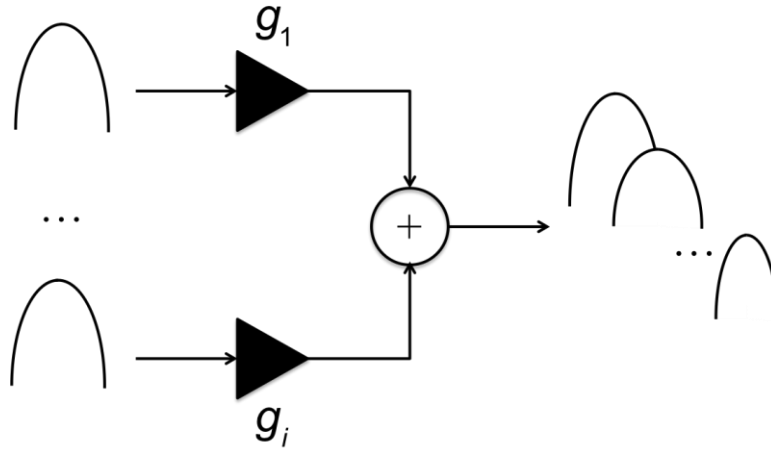
Figure 4. A DET curve is shown for a 1-speaker detection task based on the NIST-2001 development database. MixAR with 4 mixture components and only static features performs better than a GMM with 16 mixture components and static+delta features.

Figure 5. DET curves for GMM and MixAR models are shown for noisy TIMIT data with three types of additive noise: a) white, b) babble and c) car noise. A variety of SNRs are used.

Figure 6. DET curves for GMM and MixAR models on TIMIT and NTIMIT are shown. MixAR performance exceeds GMM performance while using a fewer number of parameters.

## VIII. FIGURES

(a)



(b)

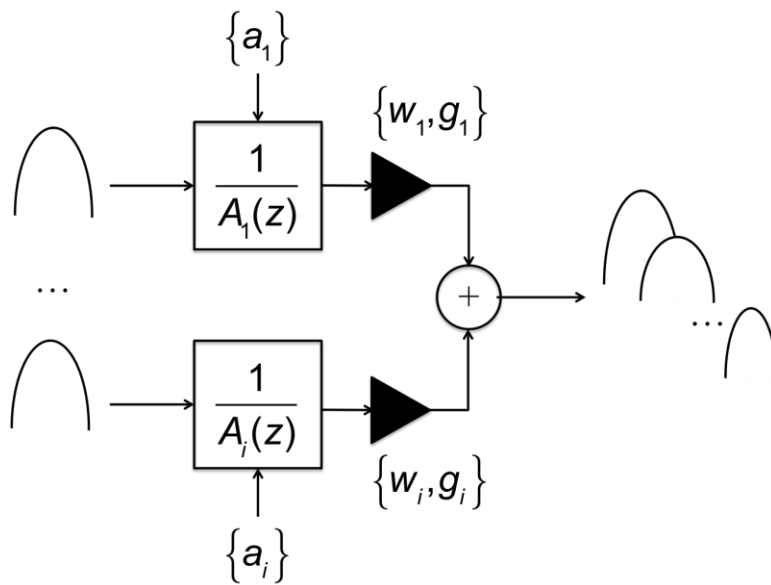


Figure 1. An overview of the (a) GMM and (b) MixAR approaches. The MixAR model is a weighted sum of Gaussian autoregressive models with time-dependent weights.

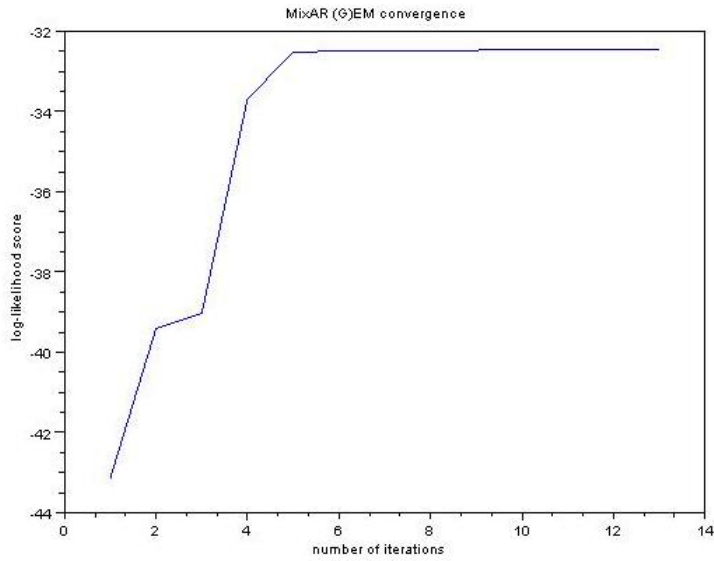


Figure 2. Performance of (Generalized) EM using the secant method as a function of the number of iterations for an 8-mixture MixAR model is shown (speaker 4516 from the NIST-2001 database).

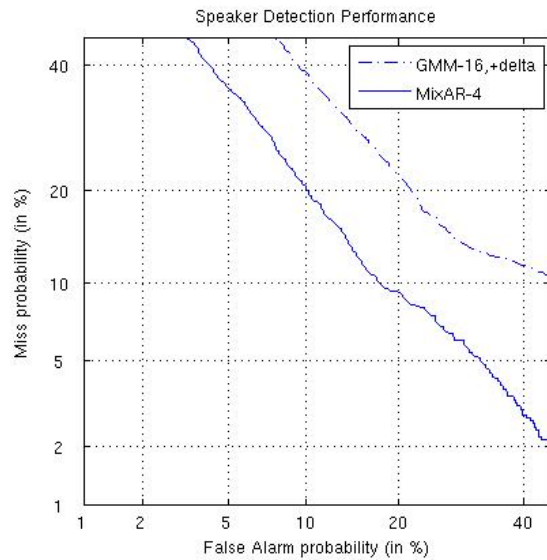


Figure 3. DET curves are shown for a simulated speaker verification task. MixAR performance in the presence of noise exceeds GMM performance.

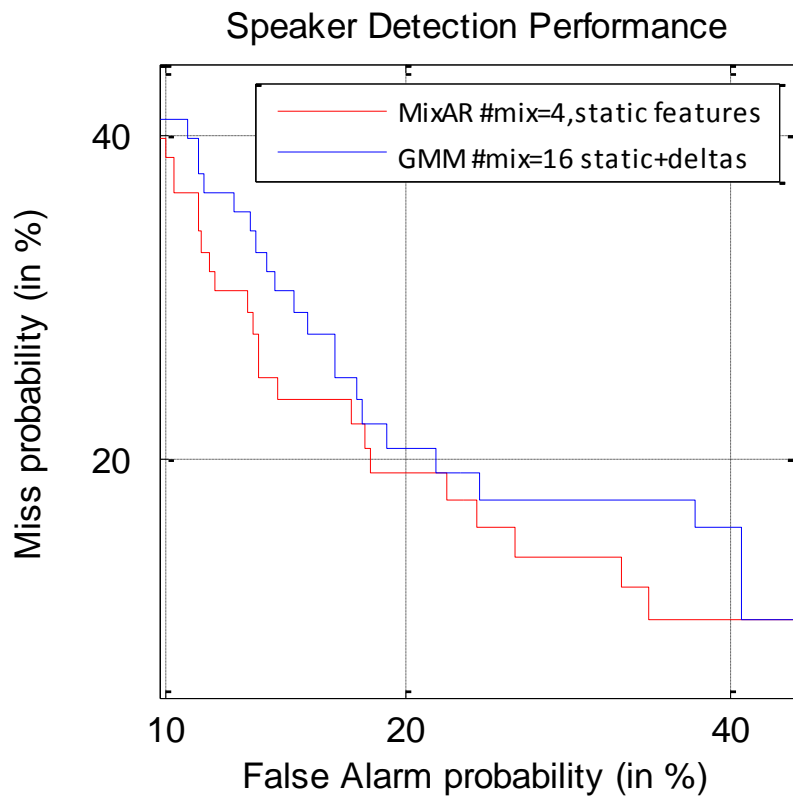


Figure 4. A DET curve is shown for a 1-speaker detection task based on the NIST-2001 development database. MixAR with 4 mixture components and only static features performs better than a GMM with 16 mixture components and static+delta features.

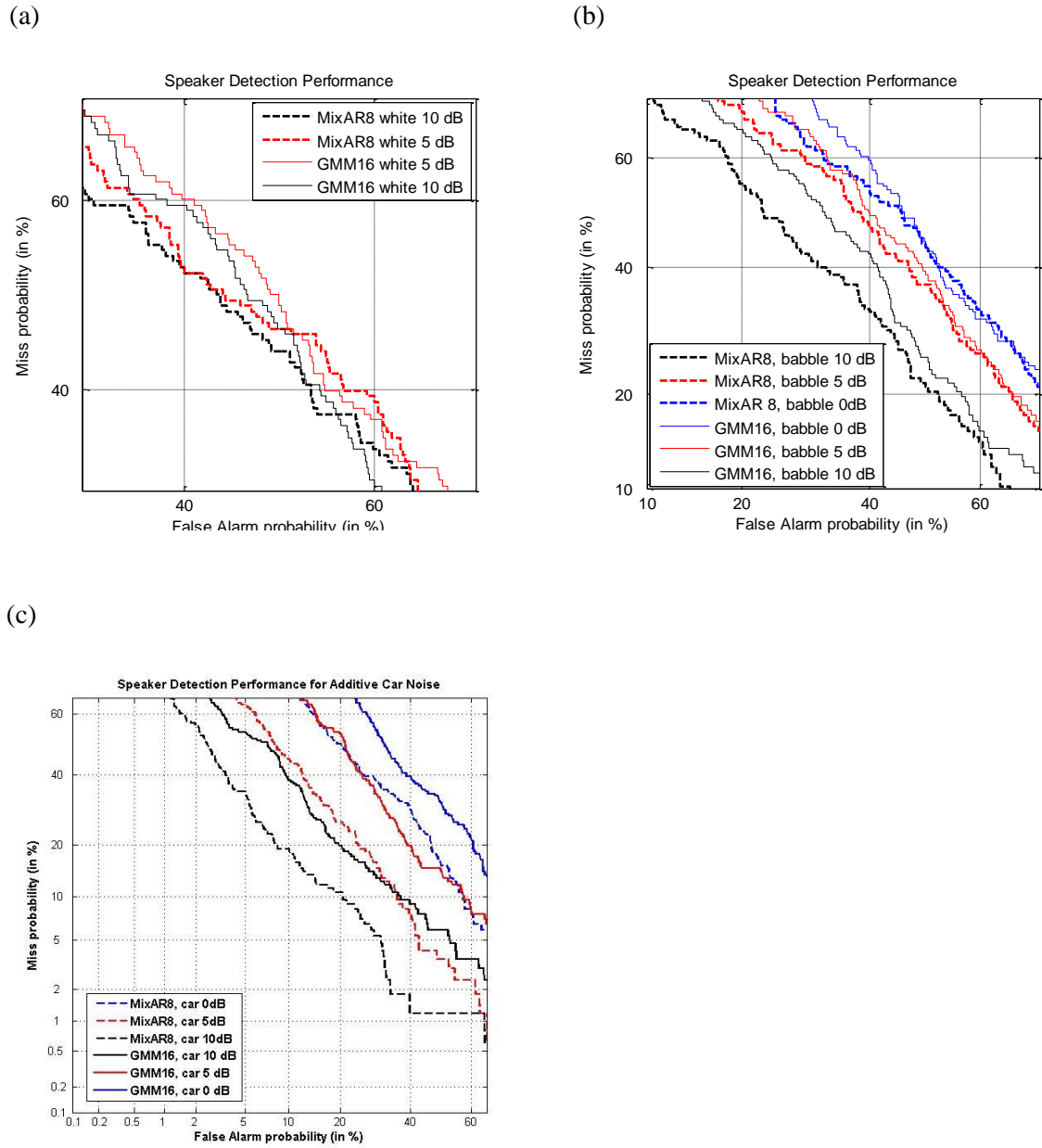


Figure 5. DET curves for GMM and MixAR models are shown for noisy TIMIT data with three types of additive noise: a) white, b) babble and c) car noise. A variety of SNRs are used.

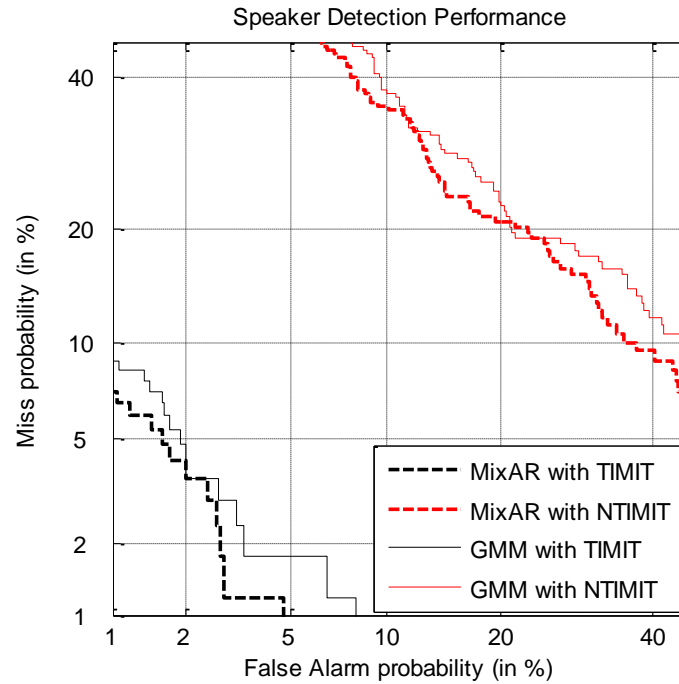


Figure 6. DET curves for GMM and MixAR models on TIMIT and NTIMIT are shown. MixAR performance exceeds GMM performance while using a fewer number of parameters.



## IX. LIST OF TABLES

Table 1. A comparison of classification error rates is shown for a GMM system using static+delta features and a MixAR system operating only on the static features is shown on synthetic data. The number of parameters for each system is shown in parentheses. The GMM system, which uses static and delta features, performs significantly worse than the MixAR system as the nonlinearity in the data increases.

Table 2. Speaker verification EERs are shown for MixAR and GMM for a variety of feature vector combinations. MixAR does not need delta features since the model itself encodes temporal dynamics.

Table 3. EERs are shown as a function of the number of mixtures. MixAR performs slightly better with almost half the number of parameters.

Table 4. EERs are shown for a variety of noise conditions.

Table 5. EERs are shown for TIMIT (clean data) and NTIMIT (noisy data).

Table 6. Performance is analyzed as a function of the duration of the training data utterances. The evaluation utterance durations were held constant and varied between 20 and 40 seconds.

Table 7. Performance is analyzed as a function of the duration of the evaluation data is shown. The training utterance duration was fixed and averaged around 120s.

## X. TABLES

$\alpha$	GMM (8 mix.) Static+ $\Delta$	MixAR (4-mix.) Static
0.00	1.50 (288)	1.50 (240)
0.25	3.25 (576)	3.50 (240)
0.50	10.25 (576)	6.25 (240)
0.75	24.75 (576)	9.75 (240)
1.00	26.75 (576)	13.75 (240)

Table 1. A comparison of classification error rates is shown for a GMM system using static+delta features and a MixAR system operating only on the static features is shown on synthetic data. The number of parameters for each system is shown in parentheses. The GMM system, which uses static and delta features, performs significantly worse than the MixAR system as the nonlinearity in the data increases.

Features	GMM 16-mix.	MixAR 8-mix.
Static(12)	22.1	19.1
Static+E(13)	33.1	41.1
Static+ $\Delta$ (24)	20.6	20.4
Static+ $\Delta$ + $\Delta\Delta$ (36)	20.5	20.5

Table 3. Speaker verification EERs are shown for MixAR and GMM for a variety of feature vector combinations. MixAR does not need delta features since the model itself encodes temporal dynamics.

No. Mixtures	GMM Static+ $\Delta$ + $\Delta\Delta$	MixAR Static Only
2	23.1 (216)	24.1(120)
4	21.7 (432)	19.2(240)
8	20.5 (864)	19.1(480)
16	20.5 (1728)	19.2(960)

Table 2. EERs are shown as a function of the number of mixtures. MixAR performs slightly better with almost half the number of parameters.

	SNR (dB)	Clean	Car Noise	White Noise	Babble Noise
	<b>GMM (1168)</b>		<b>2.4</b>		
	10 dB		19.7	48.7	40.6
	5 dB		31.2	50.0	44.7
	0 dB		39.3	49.8	48.2
<b>MixAR (480)</b>		<b>1.8</b>			
	10 dB		13.7	47.0	36.9
	5 dB		23.2	47.6	42.8
	0dB		33.9	48.5	47.6

Table 4. EERs are shown for a variety of noise conditions.

Database	GMM Static+ $\Delta$ + $\Delta\Delta$ (1728)	MixAR Static Only (480)
TIMIT	2.4	1.8
NTIMIT	21.0	20.9

Table 6. EERs are shown for TIMIT (clean data) and NTIMIT (noisy data).

	Training Utterance Duration	EER
<b>GMM (864)</b>	120	20.5
	90	20.4
	60	20.4
	30	24.4
	15	29.5
<b>MixAR (480)</b>	120	19.2
	90	21.5
	60	21.8
	30	21.8
	15	24.3

Table 5. Performance is analyzed as a function of the duration of the training data utterances. The evaluation utterance durations were held constant and varied between 20 and 40 seconds.

	Evaluation Utterance Duration	EER
<b>GMM (864)</b>	30	20.5
	15	21.8
	10	21.5
	5	24.4
	3	26.9
<b>MixAR (480)</b>	30	19.2
	15	23.4
	10	23.1
	5	25.6
	3	25.6

Table 7. Performance is analyzed as a function of the duration of the evaluation data is shown. The training utterance duration was fixed and averaged around 120s.

## I. LIST OF FIGURES

Figure 1. An overview of the (a) GMM and (b) MixAR approaches. The MixAR model is a weighted sum of Gaussian autoregressive models with time-dependent weights.

Figure 2. Performance of (Generalized) EM using the secant method as a function of the number of iterations for an 8-mixture MixAR model is shown (speaker *4516* from the NIST-2001 database).

Figure 3. DET curves are shown for a simulated speaker verification task. MixAR performance in the presence of noise exceeds GMM performance.

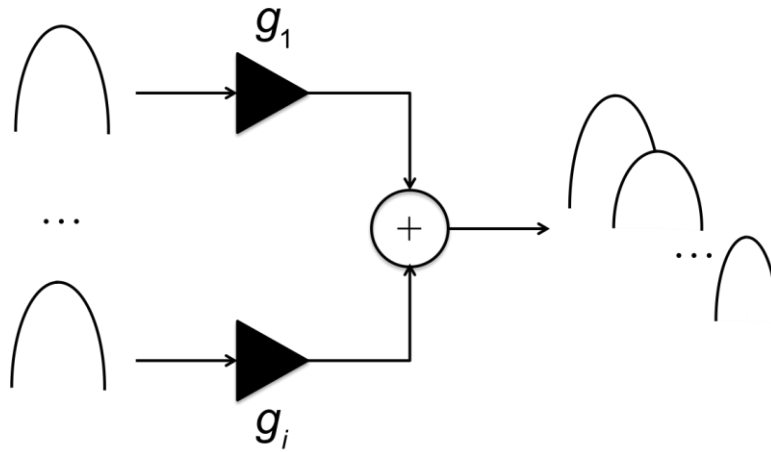
Figure 4. A DET curve is shown for a 1-speaker detection task based on the NIST-2001 development database. MixAR with 4 mixture components and only static features performs better than a GMM with 16 mixture components and static+delta features.

Figure 5. DET curves for GMM and MixAR models are shown for noisy TIMIT data with three types of additive noise: a) white, b) babble and c) car noise. A variety of SNRs are used.

Figure 6. DET curves for GMM and MixAR models on TIMIT and NTIMIT are shown. MixAR performance exceeds GMM performance while using a fewer number of parameters.

## II. FIGURES

(a)



(b)

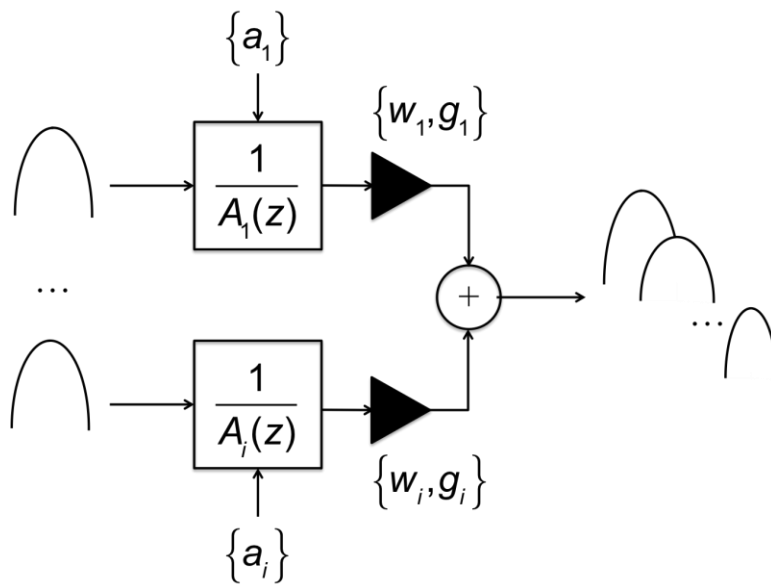


Figure 1. An overview of the (a) GMM and (b) MixAR approaches. The MixAR model is a weighted sum of Gaussian autoregressive models with time-dependent weights.

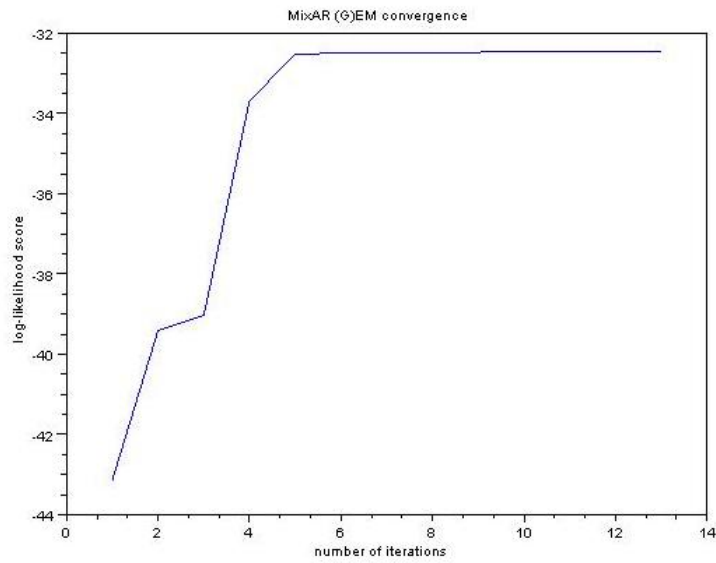


Figure 2. Performance of (Generalized) EM using the secant method as a function of the number of iterations for an 8-mixture MixAR model is shown (speaker 4516 from the NIST-2001 database).

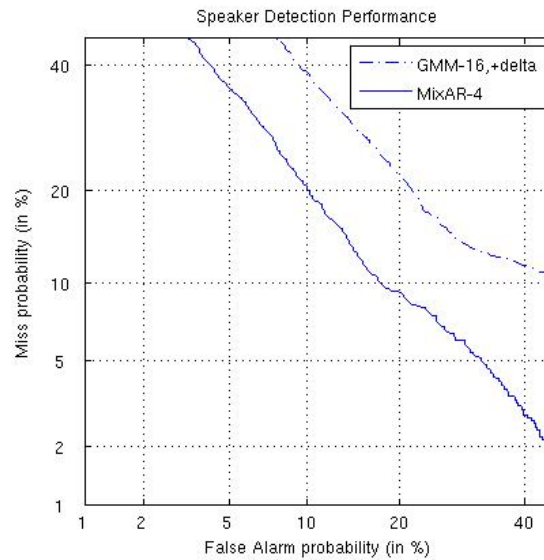


Figure 3. DET curves are shown for a simulated speaker verification task. MixAR performance in the presence of noise exceeds GMM performance.

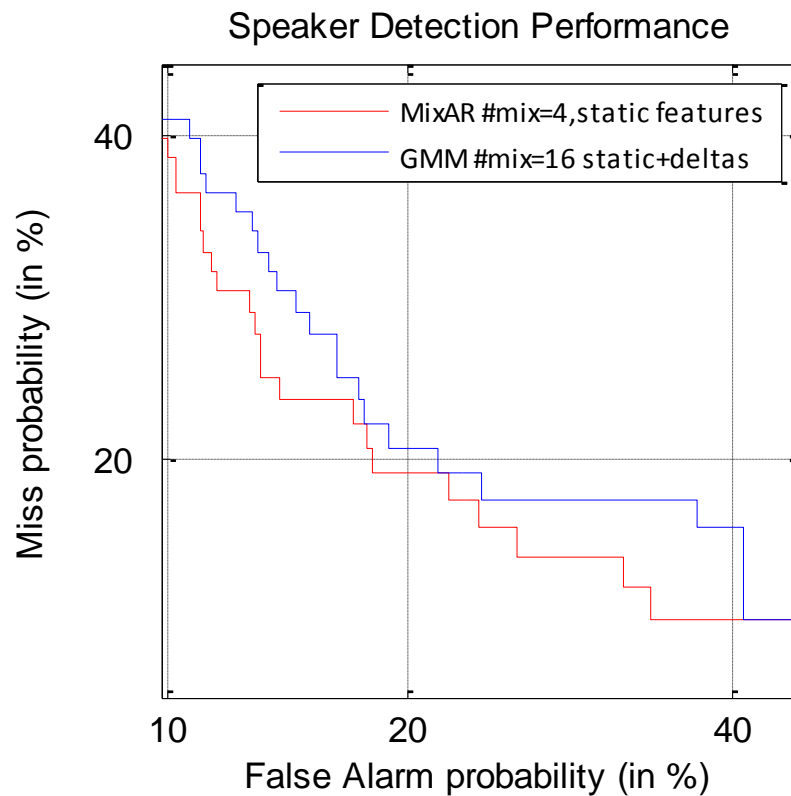


Figure 4. A DET curve is shown for a 1-speaker detection task based on the NIST-2001 development database. MixAR with 4 mixture components and only static features performs better than a GMM with 16 mixture components and static+delta features.

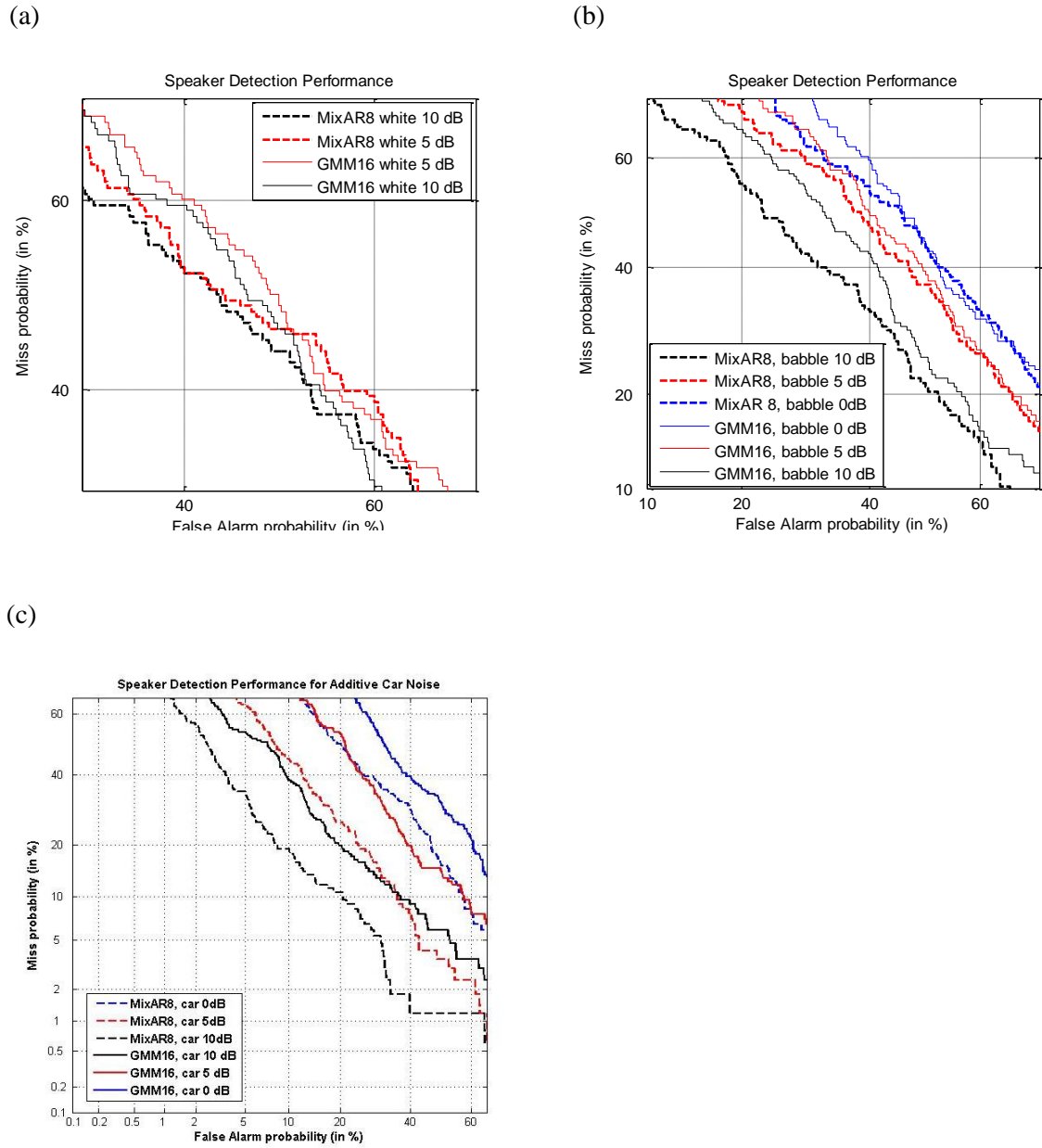


Figure 5. DET curves for GMM and MixAR models are shown for noisy TIMIT data with three types of additive noise: a) white, b) babble and c) car noise. A variety of SNRs are used.



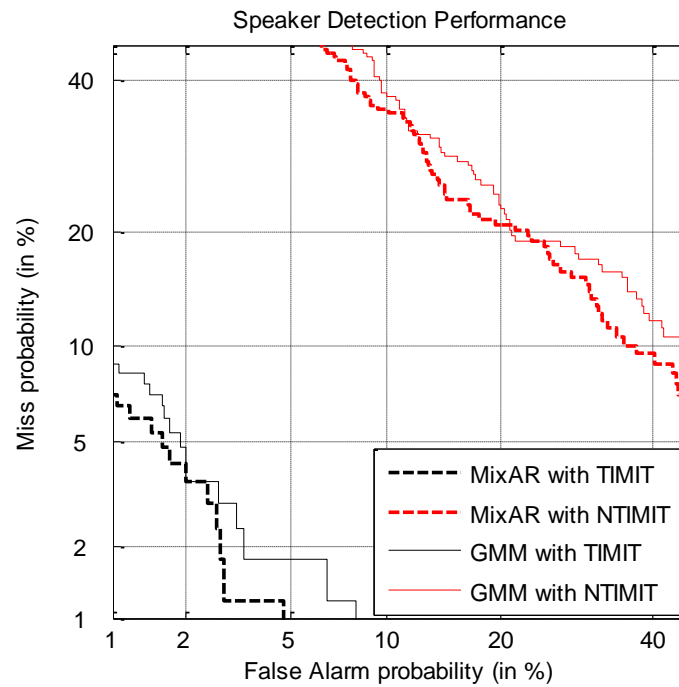


Figure 6. DET curves for GMM and MixAR models on TIMIT and NTIMIT are shown. MixAR performance exceeds GMM performance while using a fewer number of parameters.

## I. LIST OF TABLES

Table 1. A comparison of classification error rates is shown for a GMM system using static+delta features and a MixAR system operating only on the static features is shown on synthetic data. The number of parameters for each system is shown in parentheses. The GMM system, which uses static and delta features, performs significantly worse than the MixAR system as the nonlinearity in the data increases.

Table 2. Speaker verification EERs are shown for MixAR and GMM for a variety of feature vector combinations. MixAR does not need delta features since the model itself encodes temporal dynamics.

Table 3. EERs are shown as a function of the number of mixtures. MixAR performs slightly better with almost half the number of parameters.

Table 4. EERs are shown for a variety of noise conditions.

Table 5. EERs are shown for TIMIT (clean data) and NTIMIT (noisy data).

Table 6. Performance is analyzed as a function of the duration of the training data utterances. The evaluation utterance durations were held constant and varied between 20 and 40 seconds.

Table 7. Performance is analyzed as a function of the duration of the evaluation data is shown. The training utterance duration was fixed and averaged around 120s.

## II. TABLES

$\alpha$	GMM (8 mix.) Static+ $\Delta$	MixAR (4-mix.) Static
0.00	1.50 (288)	1.50 (240)
0.25	3.25 (576)	3.50 (240)
0.50	10.25 (576)	6.25 (240)
0.75	24.75 (576)	9.75 (240)
1.00	26.75 (576)	13.75 (240)

Table 1. A comparison of classification error rates is shown for a GMM system using static+delta features and a MixAR system operating only on the static features is shown on synthetic data. The number of parameters for each system is shown in parentheses. The GMM system, which uses static and delta features, performs significantly worse than the MixAR system as the nonlinearity in the data increases.

Features	GMM 16-mix.	MixAR 8-mix.
Static(12)	22.1	19.1
Static+E(13)	33.1	41.1
Static+ $\Delta$ (24)	20.6	20.4
Static+ $\Delta$ + $\Delta\Delta$ (36)	20.5	20.5

Table 3. Speaker verification EERs are shown for MixAR and GMM for a variety of feature vector combinations. MixAR does not need delta features since the model itself encodes temporal dynamics.

No. Mixtures	GMM Static+ $\Delta$ + $\Delta\Delta$	MixAR Static Only
2	23.1 (216)	24.1(120)
4	21.7 (432)	19.2(240)
8	20.5 (864)	19.1(480)
16	20.5 (1728)	19.2(960)

Table 2. EERs are shown as a function of the number of mixtures. MixAR performs slightly better with almost half the number of parameters.

	SNR (dB)	Clean	Car Noise	White Noise	Babble Noise
			<b>2.4</b>		
<b>GMM (1168)</b>	10 dB		19.7	48.7	40.6
	5 dB		31.2	50.0	44.7
	0 dB		39.3	49.8	48.2
			<b>1.8</b>		
<b>MixAR (480)</b>	10 dB		13.7	47.0	36.9
	5 dB		23.2	47.6	42.8
	0dB		33.9	48.5	47.6

Table 4. EERs are shown for a variety of noise conditions.

<b>Database</b>	<b>GMM Static+<math>\Delta</math>+<math>\Delta\Delta</math> (1728)</b>	<b>MixAR Static Only (480)</b>
TIMIT	2.4	1.8
NTIMIT	21.0	20.9

Table 6. EERs are shown for TIMIT (clean data) and NTIMIT (noisy data).

	<b>Training Utterance Duration</b>	<b>EER</b>
<b>GMM (864)</b>	120	20.5
	90	20.4
	60	20.4
	30	24.4
	15	29.5
<b>MixAR (480)</b>	120	19.2
	90	21.5
	60	21.8
	30	21.8
	15	24.3

Table 5. Performance is analyzed as a function of the duration of the training data utterances. The evaluation utterance durations were held constant and varied between 20 and 40 seconds.

	<b>Evaluation Utterance Duration</b>	<b>EER</b>
<b>GMM (864)</b>	30	20.5
	15	21.8
	10	21.5
	5	24.4
	3	26.9
<b>MixAR (480)</b>	30	19.2
	15	23.4
	10	23.1
	5	25.6
	3	25.6

Table 7. Performance is analyzed as a function of the duration of the evaluation data is shown. The training utterance duration was fixed and averaged around 120s.