

# PARAMETER SHARING AND MINIMUM CLASSIFICATION ERROR TRAINING OF MIXTURES OF FACTOR ANALYZERS FOR SPEAKER IDENTIFICATION

Hiro Yoshi Yamamoto<sup>†</sup>, Yoshihoko Nankaku<sup>†</sup>, Chiyomi Miyajima<sup>‡</sup>,  
Keiichi Tokuda<sup>†</sup>, and Tadashi Kitamura<sup>†</sup>

<sup>†</sup> Department of Computer Science and Engineering, Graduate School of Engineering  
Nagoya Institute of Technology, Gokiso-cho, Showa-ku, Nagoya 466-8555, Japan

<sup>‡</sup> Department of Media Science, Graduate School of Information Science  
Nagoya University, Furo-cho, Chikusa-ku, Nagoya 464-8603, Japan

## ABSTRACT

This paper investigates the parameter tying strategies of mixtures of factor analyzers (MFA) and discriminative training of MFA for speaker identification. The parameters of factor loading matrices or diagonal matrices are shared in different mixtures of MFA. The minimum classification error (MCE) training is applied to the MFA parameters to enhance the discrimination abilities. The results of text-independent speaker identification experiments show that MFA outperforms the conventional Gaussian mixture models (GMMs) with diagonal or full covariance matrices and achieve the best performance when sharing the diagonal matrices, resulting in a relative gain of 26% over the GMM with diagonal covariance matrices. The recognition performance is further improved by the MCE training with an additional 3% error reduction.

## 1. INTRODUCTION

Gaussian mixture models (GMMs) are widely used for text-independent speaker identification [1]. It is well known that GMM with full covariance matrices needs sufficient training data to guarantee the reliability of the estimated model parameters. Furthermore, GMM with diagonal covariance matrices requires a relatively large number of Gaussians to provide high recognition performance. In order to cope with the problem, mixtures of factor analyzers (MFA) [2] have been applied to speech as well as speaker recognition [3], [4]. MFA allows us to reduce the degree of freedom of the covariance matrices maintaining the recognition performance. Moreover, the reliability of the estimated parameters can be improved by sharing parameters in different mixture components of MFA.

In this paper, the parameter tying strategies of MFA are investigated for speaker identification. Factor loading matrices or diagonal matrices of MFA-based speaker models are shared in different mixture components assuming that all the mixture components in each MFA have the same number of factors. In this paper, the following three kinds of MFA with different parameter sharing structures are compared.

- 1) MFA without parameter sharing
- 2) MFA with shared diagonal matrices
- 3) MFA with shared factor loading matrices

In addition, minimum classification error (MCE) training is applied to MFA to improve the speaker recognition performance. The effectiveness of the MCE training for the parameter shared MFA is evaluated in a text-independent speaker identification task.

This paper is organized as follows. Sections 2 and 3 describe the general formulation of MFA and parameter tying strategies, respectively. Section 4 presents the MCE training of MFA, and

This work was partially supported by a Grant-in-Aid for Scientific Research from the Japan Society for the Promotion of Science, Encouragement of Young Scientists.

the experimental results are reported in Section 5. Finally, conclusions and future works are given in Section 6.

## 2. MIXTURES OF FACTOR ANALYZERS

### 2.1. Factor Analysis

Factor analysis (FA) is a statistical method for modeling the covariance structure of high dimensional data using a small number of latent variables. In FA, a  $d$ -dimensional speech feature vector  $\mathbf{x} = (x_1, x_2, \dots, x_d)^T$  is modeled using a  $q$ -dimensional vector  $\mathbf{z} = (z_1, z_2, \dots, z_q)^T$  and a  $d$ -dimensional observation noise  $\mathbf{n} = (n_1, n_2, \dots, n_d)^T$ :

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \mathbf{n}, \quad (1)$$

where  $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_q)$ ,  $\mathbf{w}_i = (w_{i1}, w_{i2}, \dots, w_{id})^T$  is a  $p \times d$  matrix known as a factor loading matrix, and  $\mathbf{z}$  is a latent variable assumed to be distributed according to a Gaussian density  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , i.e., zero-mean independent normals with unit variance. Each element of  $\mathbf{z}$  is referred to as "factor". The noise vector  $\mathbf{n}$  is distributed according to  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Psi})$ , where  $\boldsymbol{\mu}$  denotes a mean vector and  $\boldsymbol{\Psi}$  is a diagonal matrix.

The likelihood of an observation  $\mathbf{x}$  is given by

$$p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\boldsymbol{\mu} + \mathbf{W}\mathbf{z}, \boldsymbol{\Psi}) \quad (2)$$

because when  $\mathbf{z}$  is given, the product  $\mathbf{W}\mathbf{z}$  is a constant vector added to the observation noise vector  $\mathbf{n}$ . Therefore, distribution for  $\mathbf{x}$  is obtained by integrating out the latent variable  $\mathbf{z}$ :

$$\begin{aligned} p(\mathbf{x}) &= \int p(\mathbf{x} | \mathbf{z})p(\mathbf{z})d\mathbf{z} \\ &= \mathcal{N}(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \boldsymbol{\Psi}). \end{aligned} \quad (3)$$

### 2.2. Extension of FA to MFA

MFA is defined as mixtures of  $M$  factor analyzers. The likelihood of  $T$  independent feature vectors  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$  for the  $M$ -component MFA  $\boldsymbol{\theta} = \{c_m, \boldsymbol{\mu}_m, \mathbf{W}_m, \boldsymbol{\Psi}_m | m = 1, \dots, M\}$  is given by

$$p(\mathbf{X} | \boldsymbol{\theta}) = \prod_{t=1}^T \sum_{m=1}^M \int p_m(\mathbf{x}_t | \mathbf{z})p_m(\mathbf{z})c_m d\mathbf{z}, \quad (4)$$

where  $c_m$  denotes the weight of the  $m$ -th mixture component.

### 3. TYING STRATEGIES

Covariance matrices  $\Sigma_m$  of MFA consist of  $\mathbf{W}_m$  and  $\Psi_m$ . In this section, we present the various tying strategies of these parameters in MFA. We assume that all the mixture components have the same number of factors, and compare the following three kinds of MFA with different parameter sharing structures.

- 1) **Generic MFA:** MFA without parameter sharing, i.e., standard MFA.
- 2)  **$\Psi$ -shared MFA:** MFA with shared diagonal covariance matrices, i.e.,  $\Psi_1 = \Psi_2 = \dots = \Psi$ , where  $\mathbf{n}$  is assumed to be a sensor noise.
- 3)  **$\mathbf{W}$ -shared MFA:** MFA with shared factor loading matrices, i.e.,  $\mathbf{W}_1 = \mathbf{W}_2 = \dots = \mathbf{W}$ , where the weights of each factor in different mixtures are the same.

The maximum likelihood (ML) solution better suits the linear Gaussian model framework since the expectation maximization (EM) algorithm can be used. The EM steps for the MFA parameters  $\theta$  are summarized as follows.

#### 3.1. E-step

The E-step calculates the expectation of latent vector  $\mathbf{z}$  and the posterior of the  $m$ -th mixture component:

$$\langle \mathbf{z}_{tm} \rangle = E[\mathbf{z} | \mathbf{x}_t, m] = \beta_m (\mathbf{x}_t - \boldsymbol{\mu}_m), \quad (5)$$

$$\begin{aligned} \langle \mathbf{z} \mathbf{z}^T \rangle &= E[\mathbf{z} \mathbf{z}^T | \mathbf{x}_t, m] \\ &= \mathbf{I} - \beta_m \mathbf{W}_m + \langle \mathbf{z}_{tm} \rangle \langle \mathbf{z}_{tm} \rangle^T, \end{aligned} \quad (6)$$

$$h_{tm} = \frac{c_m \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_m, \Sigma_m)}{\sum_m c_m \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_m, \Sigma_m)}, \quad (7)$$

where  $\beta_m = \mathbf{W}_m^T \Sigma_m^{-1}$  and  $\Sigma_m = \mathbf{W}_m \mathbf{W}_m^T + \Psi_m$ .

#### 3.2. M-step

The M-step is also very straightforward. The new model parameters  $\boldsymbol{\mu}'$ ,  $\mathbf{W}'$ ,  $\Psi'$ , and  $c'_m$  for the three kinds of MFA mentioned above can be obtained by the following re-estimation formulae.

##### 1) Generic MFA

The re-estimation formulae require some manipulation to obtain the new MFA parameters using the following convenient matrix operations.

$$\tilde{\mathbf{W}}_m = (\mathbf{W}_m \boldsymbol{\mu}_m) \quad (8)$$

$$\tilde{\mathbf{z}}_{tm} = \begin{pmatrix} \mathbf{z} \\ 1 \end{pmatrix} \quad (9)$$

The re-estimates of  $\tilde{\mathbf{W}}'_m$  and  $\Psi'_m$  are obtained by

$$\tilde{\mathbf{W}}'_m = \left( \sum_t h_{tm} \mathbf{x}_t \langle \tilde{\mathbf{z}}_{tm} \rangle^T \right) \cdot \left( \sum_t h_{tm} \langle \tilde{\mathbf{z}}_{tm} \rangle \langle \tilde{\mathbf{z}}_{tm} \rangle^T \right)^{-1}, \quad (10)$$

$$\Psi'_m = \frac{1}{\sum_t h_{tm}} \text{diag} \left\{ \sum_t h_{tm} \left( \mathbf{x}_t - \tilde{\mathbf{W}}'_m \langle \tilde{\mathbf{z}}_{tm} \rangle \right) \mathbf{x}_t^T \right\}, \quad (11)$$

where

$$\langle \tilde{\mathbf{z}}_{tm} \rangle = \begin{pmatrix} \langle \mathbf{z}_{tm} \rangle \\ 1 \end{pmatrix}, \quad (12)$$

$$\langle \tilde{\mathbf{z}} \tilde{\mathbf{z}}^T \rangle = \begin{pmatrix} \langle \mathbf{z} \mathbf{z}^T \rangle & \langle \mathbf{z}_{tm} \rangle \\ \langle \mathbf{z}_{tm} \rangle & 1 \end{pmatrix}, \quad (13)$$

and  $\text{diag}(\cdot)$  denotes setting the elements outside the main diagonal to zeros. The mixture weight  $c_m$  is re-estimated as follows.

$$c'_m = \frac{1}{T} \sum_{t=1}^T h_{tm} \quad (14)$$

##### 2) $\Psi$ -shared MFA

The re-estimation formulae for  $\Psi$ -shared MFA are the same as those for generic MFA except for the diagonal covariance matrix:

$$\Psi' = \frac{1}{T} \text{diag} \left\{ \sum_{t,m} h_{tm} \left( \mathbf{x}_t - \tilde{\mathbf{W}}'_m \langle \tilde{\mathbf{z}}_{tm} \rangle \right) \mathbf{x}_t^T \right\}. \quad (15)$$

##### 3) $\mathbf{W}$ -shared MFA

The new model parameters of  $\mathbf{W}$ -shared MFA is re-estimated as follows. The new factor loading matrix  $\mathbf{W}'$  is given by

$$\begin{aligned} \mathbf{W}'_{(k)} &= \left( \sum_{t,m} h_{tm} \Psi_{m(k)}^{-1} (\mathbf{x}_t - \boldsymbol{\mu}_{m(k)}) \langle \mathbf{z}_{tm} \rangle^T \right) \\ &\cdot \left( \sum_{t,m} h_{tm} \Psi_{m(k)}^{-1} \langle \mathbf{z} \mathbf{z}^T \rangle \right)^{-1} \end{aligned} \quad (16)$$

where  $\mathbf{W}_{(k)}$  is  $k$ -th row vector in the factor loading matrix  $\mathbf{W}$ . In the followings, the individual component parameters  $\boldsymbol{\mu}'_m$  and  $\Psi'_m$  can be re-estimated:

$$\boldsymbol{\mu}'_m = \frac{\sum_t h_{tm} (\mathbf{x}_t - \mathbf{W}' \langle \mathbf{z}_{tm} \rangle)}{\sum_t h_{tm}}, \quad (17)$$

$$\begin{aligned} \Psi'_m &= \frac{1}{\sum_t h_{tm}} \text{diag} \sum_t \left\{ h_{tm} (\mathbf{x}_t - \boldsymbol{\mu}'_m) (\mathbf{x}_t - \boldsymbol{\mu}'_m)^T \right. \\ &\quad \left. - h_{tm} \mathbf{W}' \left( 2 \langle \mathbf{z}_{tm} \rangle (\mathbf{x}_t - \boldsymbol{\mu}'_m)^T - \langle \mathbf{z} \mathbf{z}^T \rangle \mathbf{W}'^T \right) \right\}. \end{aligned} \quad (18)$$

## 4. MCE TRAINING FOR MFA SPEAKER MODEL

To enhance the discrimination abilities of MFA-based speaker models, MCE training based on the generalized probabilistic descent (GPD) method [5] is applied to the parameters of MFA [3].

### 4.1. Definition of Loss Function

For the MCE training, the misclassification measure of training data  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$  for speaker  $s$  is defined as

$$d_s(\mathbf{X}; \Theta) = -g_s(\mathbf{X}; \Theta) + \max_{y \neq s} g_y(\mathbf{X}; \Theta), \quad (19)$$

where  $\Theta = \{\theta_1, \theta_2, \dots, \theta_S\}$  denotes the speaker model parameter set of MFA, and  $g_s(\cdot; \cdot)$  is defined by the log likelihood of  $\mathbf{X}$  for speaker model  $\theta_s$ . Equation (19) is the approximation of the log likelihood ratio between the competing models and the correct one. The loss function is defined as a differentiable sigmoid function approximating the 0-1 step loss function:

$$l_s(\mathbf{X}; \theta) = (1 + \exp(-\gamma \cdot d_s))^{-1}, \quad (20)$$

where  $\gamma$  denotes the gradient of the sigmoid function. The goal of the discriminative training is to minimize the loss function based on the probabilistic descent method.

## 4.2. Parameter Adjustment of MFA

During the parameter adaptation in the MCE training, the constraints of the MFA parameters, e.g.,  $c_m > 0$ , should be satisfied. Hence, the MFA parameter set  $\Theta$  is transformed into a new model parameter set  $\tilde{\Theta}$ .

$$\tilde{\Theta} = \{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_S\}, \quad (21)$$

$$\tilde{\theta} = \{\tilde{c}_m, \tilde{\mu}_m, \mathbf{W}_m, \tilde{\Psi}_m \mid m = 1, 2, \dots, M\}, \quad (22)$$

where  $\tilde{c}_m = \log c_m$ ,  $\tilde{\mu}_{mi} = \frac{\mu_{mi}}{\Sigma_{mii}}$ ,  $\tilde{\Psi}_{mii} = \log \Psi_{mii}$ .  $\tilde{\Theta}$  is updated at each iteration  $r$  as

$$\tilde{\Theta}(r+1) = \tilde{\Theta}(r) - \varepsilon_r \nabla l_s(\mathbf{X}; \tilde{\theta}), \quad (23)$$

where  $\varepsilon_r$  is a monotonically decreasing learning step size at the  $r$ -th iteration. In this paper,  $\tilde{\Theta}$  is sequentially adjusted every time a training sample  $\mathbf{X}$  is given (i.e., sample-by-sample mode).

The gradient of (23) is obtained as follows.

$$\nabla_{\tilde{\theta}_y} l_s(\mathbf{X}; \tilde{\theta}) = \frac{\partial l_s}{\partial d_s} \frac{\partial d_s}{\partial g_y} \cdot \nabla_{\tilde{\theta}_y} g_y(\mathbf{X}; \tilde{\theta}), \quad (24)$$

where  $\frac{\partial l_s}{\partial d_s}$ ,  $\frac{\partial d_s}{\partial g_y}$ ,  $\nabla_{\tilde{\theta}_y} g_y(\mathbf{X}; \tilde{\theta})$  are given by

$$\frac{\partial l_s}{\partial d_s} = \gamma l_s(1 - l_s), \quad \frac{\partial d_s}{\partial g_y} = \begin{cases} -1, & y = s \\ 1, & y \neq s \end{cases}, \quad (25)$$

$$\nabla_{\tilde{\theta}_y} g_y(\mathbf{X}; \tilde{\theta}) = \frac{1}{T} \sum_{t=1}^T \frac{1}{b_y(\mathbf{x}_t)} \nabla_{\tilde{\theta}_y} b_y(\mathbf{x}_t). \quad (26)$$

For the three kinds of MFA, the gradient of  $b_y(\mathbf{x}_t)$  with respect to each element in  $\tilde{\theta}_y$  is obtained by the following formulae, where the subscript  $y$  is dropped for the simplicity of notation.

### 1) generic MFA

For the generic MFA, the gradients are obtained as follows.

$$\frac{\partial b(\mathbf{x}_t)}{\partial \tilde{c}_m} = f_m, \quad \frac{\partial b(\mathbf{x}_t)}{\partial \tilde{\mu}_{mi}} = f_m \delta_{mi} \Sigma_{mii}, \quad (27)$$

$$\frac{\partial b(\mathbf{x}_t)}{\partial W_{mij}} = -f_m \left\{ (\Sigma_m^{-1} \mathbf{W}_m)_{ij} - \delta_{mi} [\delta_m^T \mathbf{W}_m]_j \right\}, \quad (28)$$

$$\frac{\partial b(\mathbf{x}_t)}{\partial \tilde{\Psi}_{mii}} = -\frac{1}{2} f_m \left\{ \Sigma_{mii}^{-1} - \delta_m^2 \right\} \Psi_{mii}, \quad (29)$$

where  $f_m = c_m \mathcal{N}(\mathbf{x}_t \mid \mu_m, \Sigma_m)$ ,  $\delta_m = \Sigma_m^{-1}(\mathbf{x}_t - \mu_m)$ , and  $[\cdot]_i$  denotes the  $i$ -th vector element.

### 2) $\Psi$ -shared MFA

In the case of  $\Psi$ -shared MFA, the gradients with respect to mixture weights, mean vectors and factor loading matrices are obtained by follows: (27) and (28), respectively, and only (29) is changed as

$$\frac{\partial b(\mathbf{x}_t)}{\partial \tilde{\Psi}_{ii}} = \sum_{m=1}^M \frac{\partial b(\mathbf{x}_t)}{\partial \tilde{\Psi}_{mii}}. \quad (30)$$

### 3) $W$ -shared MFA

The gradients in (27) and (29) apply to the  $W$ -shared MFA case, and (28) is changed as follows.

$$\frac{\partial b(\mathbf{x}_t)}{\partial W_{ij}} = \sum_{m=1}^M \frac{\partial b(\mathbf{x}_t)}{\partial W_{mij}} \quad (31)$$

## 5. EXPERIMENTAL EVALUATION

### 5.1. Database and Experimental Conditions

A text-independent speaker identification experiment was conducted for 80 speakers (40 males and 40 females) in the ATR Japanese speech database. 216 words were used for training each speaker model, and 520 words were used for testing. The number of tests was 41600 in total. The speech data was down-sampled from 20kHz to 10kHz, windowed at a 10-ms frame rate using a 25-ms Blackman window, and parameterized into 12 mel-cepstral coefficients excluding zero-th coefficients with a mel-cepstral analysis technique.

GMM parameters were initialized using an LBG codebook. Mixture weights and mean vectors of MFA were also initialized using the LBG codebook, and factor loading matrices were initialized with random values considering full covariance. Diagonal covariance matrices were initialized using diagonal elements of full covariance matrices  $\Sigma$  [2]. The number of mixture components was changed from 4 to 64, and the number of factors was changed from 2 to 10.

### 5.2. Results

Figures 1–3 compare the identification error rates between the three kinds of MFA and the conventional GMMs with full or diagonal covariance matrices (full-GMM and diag-GMM). All speaker models in Figs.1–3 were trained with 216 words based on ML-estimation. Figure 1 compares the results of generic MFA with the conventional GMMs, where the number of factors  $q$  is changed as 2, 4, 6, 8, and 10. The horizontal axis corresponds to the number of model parameters in a logarithmic scale. Generic MFA show better performance with a smaller number of factors. However, the error rate of generic MFA with a larger number of factors is close to that of full-GMM, because the model structure of generic MFA with  $q = 12$  is almost the same as the full-GMM. Figure 2 shows the results of  $\Psi$ -shared MFA.  $\Psi$ -shared MFA achieved a significant improvement over the conventional GMMs with larger number of mixtures. In the case of 64-mixture models, error reductions of 19% ( $q = 2$ ) and 26% ( $q = 6$ ) over diag-GMM were obtained. Figure 3 shows the results of  $W$ -shared MFA. The performance of  $W$ -shared MFA is almost equivalent to that of diag-GMM, because the model structure of  $W$ -shared MFA is similar to that of diag-GMM, and has the lowest flexibility among the three kinds of sharing structures.

Figure 4 shows the results of the three kinds of MFA with the number of factors  $q = 2$  and diag-GMM, where the amount of training data was changed as 27, 54, and 216 words. We can see that the MFA-based speaker models show relatively high performance with such a small number of factors, and all the MFA-based speaker models outperform the conventional GMM with any amount of training data. Among the three kinds of MFA,  $\Psi$ -shared MFA achieves the best performance and a significant difference is found with smaller amounts of training data.

Finally, MCE training is applied to the MFA-based speaker models. Figure 5 compares the performance of  $\Psi$ -shared MFA with six factors before and after the MCE training. We can see that the performance is further improved by the MCE training and the 2.73% error rate was reduced to 2.65% with a 3% error reduction.

## 6. CONCLUSIONS

This paper has investigated the parameter tying strategies of MFA for speaker identification and MCE training has been applied to the parameter shared MFA. Sharing diagonal covariance matrices provided the best performance leading to a relative gain of 26% over the GMM with diagonal covariance matrices. The MCE training has further improved the recognition performance.

Our future works include the application of other variations of MFA to speaker identification [6] and automatic determination of the optimal number of mixture components and factors using the variational Bayesian approach [7].

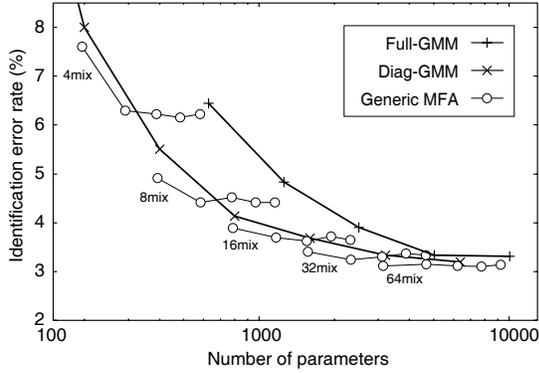


Fig. 1. Comparison between generic MFA and conventional GMMs with diagonal or full covariance matrices.

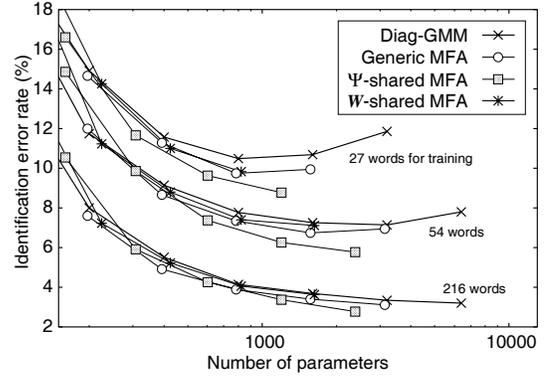


Fig. 4. Comparison between diag-GMM and three kinds MFA ( $q = 2$ ) with increasing the number of mixtures, using 27 words for training (upper), 54 words (middle), 216 words (lower) for training.

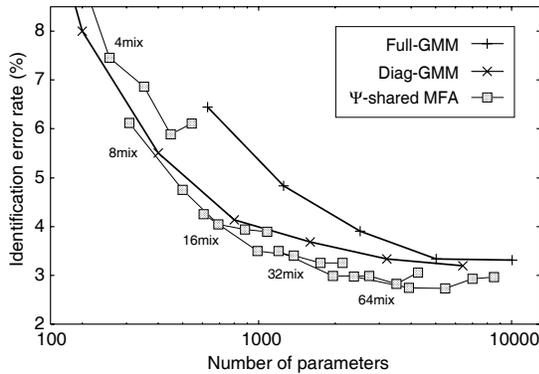


Fig. 2. Comparison between  $\Psi$ -shared MFA and conventional GMMs with diagonal or full covariance matrices.

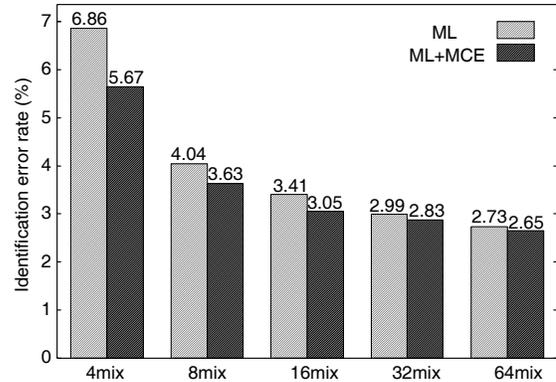


Fig. 5. Comparison of  $\Psi$ -shared MFA before and after MCE training ( $q = 6$ ).

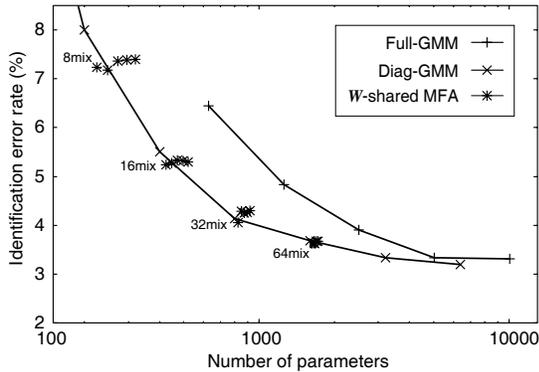


Fig. 3. Comparison between  $W$ -shared MFA and conventional GMMs with diagonal or full covariance matrices.

## 7. REFERENCES

- [1] D. A. Reynolds and R. C. Rose, "Robust text-independent speaker identification using Gaussian mixture speaker models," *IEEE Trans. on Speech and Audio Processing*, vol.3, no.1, pp.72–83, Jan. 1995.
- [2] Z. Ghahramani and G. E. Hinton, "The EM algorithm for mixtures of factor analyzers," *Tech. Rep. Univ. of Toronto, CRGTR-96-1*, May 1996.

- [3] L. K. Saul and M. G. Rahim, "Maximum likelihood and minimum classification error factor analysis for automatic speech recognition," *IEEE Trans. on Speech and Audio Processing*, vol.8, no.2, pp.115–125, Mar. 2000.
- [4] P. Ding, Y. Liu, and B. Xu, "Factor analyzed Gaussian mixture models for speaker identification," *Proc. of ICSLP-2002*, pp.1341–1344, Sept. 2002.
- [5] B.-H. Juang and S. Katagiri, "Discriminative learning for minimum error classification," *IEEE Trans. Signal Process.*, vol.40, no.12, pp.3043–3054, Dec. 1992.
- [6] A.-V. I. Rosti and M. J. F. Gales, "Generalised linear Gaussian models," *Tech. Rep. Cambridge Univ., CUED/F-INFENG/TR.420*, Nov. 2001.
- [7] Z. Ghahramani and M. J. Beal, "Variational inference for Bayesian mixtures of factor analysers," *Neural Information Processing Systems 12*, 1999.