



# Voice conversion based on matrix variate Gaussian mixture model using multiple frame features

Yi Yang, Hidetsugu Uchida, Daisuke Saito, Nobuaki Minematsu

The University of Tokyo, Japan

{yang, uchida, dsk\_saito, mine}@gavo.t.u-tokyo.ac.jp

## Abstract

This paper presents a novel voice conversion method based on matrix variate Gaussian mixture model (MV-GMM) using features of multiple frames. In voice conversion studies, approaches based on Gaussian mixture models (GMM) are still widely utilized because of their flexibility and easiness in handling. They treat the joint probability density function (PDF) of feature vectors from source and target speakers as that of joint vectors of the two vectors. Addition of dynamic features to the feature vectors in GMM-based approaches achieves certain performance improvements because the correlation between multiple frames is taken into account. Recently, a voice conversion framework based on MV-GMM, in which the joint PDF is modeled in a matrix variate space, has been proposed and it is able to precisely model both the characteristics of the feature spaces and the relation between the source and target speakers. In this paper, in order to additionally model the correlation between multiple frames in the framework more consistently, MV-GMM is constructed in a matrix variate space containing the features of neighboring frames. Experimental results show that a certain performance improvement in both objective and subjective evaluations is observed.

**Index Terms:** voice conversion, Gaussian mixture model, matrix variate Gaussian mixture model, multiple frame features

## 1. Introduction

Voice conversion (VC), specifically speaker conversion discussed in this paper, is a technique to modify non-linguistic information — in this case speaker characteristics — while maintaining linguistic information unchanged. In speaker conversion, a statistical mapping function is constructed using pairs of features extracted from utterances of source and target speakers. Gaussian mixture model (GMM) [1] is widely used to construct mapping functions, as well as Neural Network (NN) [2] and Nonnegative Matrix Factorization (NMF) [3].

Founded on flexibilities of GMM-based approaches, several expanded models are proposed, such as a method in which dynamic features are taken into account [4]. In this method, relationship between features of adjacent frames are considered, and a conversion function is designed to maximize the likelihood of a time sequence of the target features given that of the source ones. Since dynamics of speech are captured, conversion performance of the method is improved well.

Generally in GMM-based VC, both the features extracted from the source and target speakers are concatenated and they are represented as a joint vector. In addition, when dynamic features are taken into account, they are also included in the joint vector. Finally, the characteristics of the static feature spaces, those of the dynamic ones, and the correlation of the source and

target speakers are mixed together in the joint vector space. In order to model VC functions precisely, these information should be properly treated.

Voice conversion based on matrix variate Gaussian mixture models (MV-GMM), in which joint features from the source and target are represented as matrices, has been proposed [5]. In this model, a separable structure is derived to covariances of matrix features. Two types of covariance matrices, i.e. row and column matrices capture the characteristics of the feature spaces and the correlation between the source and target, respectively. Since they induce an effective training algorithm for the model, MV-GMM constructs a precise conversion function.

Basically both GMM and MV-GMM are mixture models of normal distributions. Hence, similar extensions to GMM-based VC can be applied to the MV-GMM approach. In this paper, in order to additionally model the correlation between adjacent frames in MV-GMM, voice conversion based on MV-GMM using multiple frame features is proposed. In the proposed method, it is expected that the separable structure in MV-GMM enables the model to capture the correlation between adjacent frames more precisely than the joint vector approach.

## 2. GMM-based VC with joint vectors

In this section, GMM-based voice conversion is briefly described [1]. A feature vector of time index  $t$  from an utterance of a source speaker is defined as  $\mathbf{x}_t = [x_1, x_2, \dots, x_n]^\top$ , while  $\mathbf{y}_t = [y_1, y_2, \dots, y_n]^\top$  represents that from an utterance of a target speaker, where  $n$  denotes the dimension of the feature vectors and  $(\cdot)^\top$  notifies the transposition of vector or matrix. Note that these utterances include the same linguistic content. The probability density of the joint vector  $\mathbf{z}_t = [\mathbf{x}_t^\top, \mathbf{y}_t^\top]^\top$  is modeled by GMM as follows:

$$P(\mathbf{z}_t | \lambda^{(z)}) = \sum_{m=1}^M w_m \mathcal{N}(\mathbf{z}_t; \boldsymbol{\mu}_m^{(z)}, \boldsymbol{\Sigma}_m^{(z)}), \quad (1)$$

where  $w_m$ ,  $\boldsymbol{\mu}_m^{(z)}$  and  $\boldsymbol{\Sigma}_m^{(z)}$  denote the weight, the mean vector, and the covariance matrix of the  $m$ -th Gaussian component, respectively.  $\boldsymbol{\mu}_m^{(z)}$ ,  $\boldsymbol{\Sigma}_m^{(z)}$  is indicated by mean vectors and covariance matrices of the source and target speakers separately, which shown as follows:

$$\boldsymbol{\mu}_m^{(z)} = \begin{bmatrix} \boldsymbol{\mu}_m^{(x)} \\ \boldsymbol{\mu}_m^{(y)} \end{bmatrix}, \boldsymbol{\Sigma}_m^{(z)} = \begin{bmatrix} \boldsymbol{\Sigma}_m^{(xx)} & \boldsymbol{\Sigma}_m^{(xy)} \\ \boldsymbol{\Sigma}_m^{(yx)} & \boldsymbol{\Sigma}_m^{(yy)} \end{bmatrix}. \quad (2)$$

$\boldsymbol{\Sigma}_m^{(\cdot)}$  shows covariance matrices of the source and target speakers, or cross-covariance matrices between them, which are usually restricted to diagonal matrices in order to mitigate the influence of overfitting.

A mapping function to convert the source vector  $\mathbf{x}_t$  to the target vector  $\mathbf{y}_t$  is derived based on the conditional probability density of  $\mathbf{y}_t$ . This probability density can be represented by the parameters of the joint density model as follows:

$$P(\mathbf{y}_t | \mathbf{x}_t, \lambda^{(z)}) = \sum_{m=1}^M P(m | \mathbf{x}_t, \lambda^{(z)}) P(\mathbf{y}_t | \mathbf{x}_t, m, \lambda^{(z)}), \quad (3)$$

where

$$P(m | \mathbf{x}_t, \lambda^{(z)}) = \frac{w_m \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_m^{(x)}, \boldsymbol{\Sigma}_m^{(xx)})}{\sum_{m=1}^M w_m \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_m^{(x)}, \boldsymbol{\Sigma}_m^{(xx)})}, \quad (4)$$

$$P(\mathbf{y}_t | \mathbf{x}_t, m, \lambda^{(z)}) = \mathcal{N}(\mathbf{y}_t; \mathbf{E}_{m,t}^{(y)}, \mathbf{D}_{m,t}^{(y)}), \quad (5)$$

$$\mathbf{E}_{m,t}^{(y)} = \boldsymbol{\mu}_m^{(y)} + \boldsymbol{\Sigma}_m^{(yx)} \boldsymbol{\Sigma}_m^{(xx)-1} (\mathbf{x}_t - \boldsymbol{\mu}_m^{(x)}), \quad (6)$$

$$\mathbf{D}_{m,t}^{(y)} = \boldsymbol{\Sigma}_m^{(yy)} - \boldsymbol{\Sigma}_m^{(yx)} \boldsymbol{\Sigma}_m^{(xx)-1} \boldsymbol{\Sigma}_m^{(xy)}. \quad (7)$$

The converted feature  $\hat{\mathbf{y}}_t$  when  $\mathbf{x}_t$  is given can be generated by the following equation using maximum likelihood criterion:

$$\hat{\mathbf{y}}_t = \left( \sum_{m=1}^M \gamma_{m,t} \mathbf{D}_{m,t}^{(y)-1} \right)^{-1} \left( \sum_{m=1}^M \gamma_{m,t} \mathbf{D}_{m,t}^{(y)-1} \mathbf{E}_{m,t}^{(y)} \right), \quad (8)$$

$$\gamma_{m,t} = P(m | \mathbf{x}_t, \mathbf{y}_t, \lambda^{(z)}).$$

Dynamic features may also be contained in a GMM-based model [4]. They are usually defined as difference between the previous and next frames, and they are denoted by  $\Delta \mathbf{x}_t$  and  $\Delta \mathbf{y}_t$ .  $\mathbf{X}_t = [\mathbf{x}_t^\top, \Delta \mathbf{x}_t^\top]^\top$ ,  $\mathbf{Y}_t = [\mathbf{y}_t^\top, \Delta \mathbf{y}_t^\top]^\top$  denote  $2D$ -dimensional feature vectors which consist of both static and dynamic features. Time sequences of the source and target feature vectors are vectorized and described as  $\mathbf{x} = [\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_T^\top]^\top$ ,  $\mathbf{y} = [\mathbf{y}_1^\top, \mathbf{y}_2^\top, \dots, \mathbf{y}_T^\top]^\top$ ,  $\mathbf{X} = [\mathbf{X}_1^\top, \mathbf{X}_2^\top, \dots, \mathbf{X}_T^\top]^\top$ , and  $\mathbf{Y} = [\mathbf{Y}_1^\top, \mathbf{Y}_2^\top, \dots, \mathbf{Y}_T^\top]^\top$ . For each frame, a joint vector is defined as  $\mathbf{Z}_t = [\mathbf{X}_t^\top, \mathbf{Y}_t^\top]^\top$ , and the joint probability density is modeled by GMM. When a time sequence  $\mathbf{X}$  is given, the conditional probability density of  $\mathbf{Y}$  can be shown as:

$$P(\mathbf{Y} | \mathbf{X}, \lambda^{(z)}) = \sum_{m=1}^M P(m | \mathbf{X}, \lambda^{(z)}) P(\mathbf{Y} | \mathbf{X}, m, \lambda^{(z)}) \\ = \prod_{t=1}^T \sum_{m=1}^M P(m | \mathbf{X}_t, \lambda^{(z)}) P(\mathbf{Y}_t | \mathbf{X}_t, m, \lambda^{(z)}). \quad (9)$$

The conditional probability density at each frame can also be modeled as GMM. At frame  $t$ , the posterior of the  $m$ -th mixture component and its corresponding conditional probability density of  $\mathbf{Y}_t$  in Equation 9 are described as

$$P(m | \mathbf{X}_t, \lambda^{(z)}) = \frac{w_m \mathcal{N}(\mathbf{X}_t; \boldsymbol{\mu}_m^{(X)}, \boldsymbol{\Sigma}_m^{(XX)})}{\sum_{m=1}^M w_m \mathcal{N}(\mathbf{X}_t; \boldsymbol{\mu}_m^{(X)}, \boldsymbol{\Sigma}_m^{(XX)})}, \quad (10)$$

$$P(\mathbf{Y}_t | \mathbf{X}_t, m, \lambda^{(z)}) = \mathcal{N}(\mathbf{Y}_t; \mathbf{E}_{m,t}^{(Y)}, \mathbf{D}_{m,t}^{(Y)}), \quad (11)$$

$$\mathbf{E}_{m,t}^{(Y)} = \boldsymbol{\mu}_m^{(Y)} + \boldsymbol{\Sigma}_m^{(YX)} \boldsymbol{\Sigma}_m^{(XX)-1} (\mathbf{X}_t - \boldsymbol{\mu}_m^{(X)}) \quad (12)$$

$$\mathbf{D}_{m,t}^{(Y)} = \boldsymbol{\Sigma}_m^{(YY)} - \boldsymbol{\Sigma}_m^{(YX)} \boldsymbol{\Sigma}_m^{(XX)-1} \boldsymbol{\Sigma}_m^{(XY)} \quad (13)$$

Finally, a target sequence  $\hat{\mathbf{y}}$  can be derived as a solution of the following optimization problem:

$$\hat{\mathbf{y}} = \arg \max P(\mathbf{Y} | \mathbf{X}, \lambda^{(z)}) \text{ s.t. } \mathbf{Y} = \mathbf{W} \mathbf{y}. \quad (14)$$

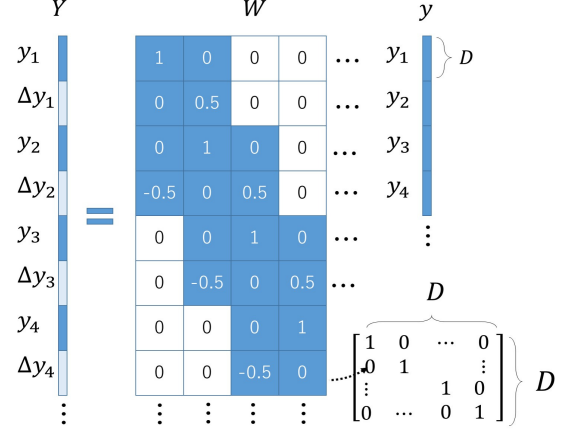


Figure 1: Relation between  $\mathbf{y}_t$  and  $\mathbf{Y}_t$ .  $\Delta \mathbf{y}_t$  is defined as  $0.5(\mathbf{y}_{t+1} - \mathbf{y}_{t-1})$ .

Figure 1 is an example of  $\mathbf{W}$ . In this case, the dynamic feature  $\Delta \mathbf{y}_t$  is defined as  $0.5(\mathbf{y}_{t+1} - \mathbf{y}_{t-1})$ , and the corresponding window matrix  $\mathbf{W}$  is derived. The conversion function is written as the equation below based on maximum likelihood estimation:

$$\hat{\mathbf{y}} = \left( \mathbf{W}^\top \overline{\mathbf{D}^{(Y)-1}} \mathbf{W} \right)^{-1} \mathbf{W}^\top \overline{\mathbf{D}^{(Y)-1}} \mathbf{E}^{(Y)}. \quad (15)$$

$\overline{\mathbf{D}^{(Y)-1}}$ ,  $\overline{\mathbf{D}^{(Y)-1}} \mathbf{E}^{(Y)}$  are parameters in time sequence, which written as:

$$\overline{\mathbf{D}^{(Y)-1}} = \text{diag}[\overline{\mathbf{D}_1^{(Y)-1}}, \overline{\mathbf{D}_2^{(Y)-1}}, \dots, \overline{\mathbf{D}_T^{(Y)-1}}], \quad (16)$$

$$\overline{\mathbf{D}^{(Y)-1}} \mathbf{E}^{(Y)} = [\overline{\mathbf{D}_1^{(Y)-1}} \mathbf{E}^{(Y)}, \overline{\mathbf{D}_2^{(Y)-1}} \mathbf{E}^{(Y)}, \dots, \overline{\mathbf{D}_T^{(Y)-1}} \mathbf{E}^{(Y)}]^\top \quad (17)$$

$$\overline{\mathbf{D}_t^{(Y)-1}} = \sum_{m=1}^M \gamma_{m,t} \mathbf{D}_m^{(Y)-1} \quad (18)$$

$$\overline{\mathbf{D}_t^{(Y)-1}} \mathbf{E}^{(Y)} = \sum_{m=1}^M \gamma_{m,t} \mathbf{D}_m^{(Y)-1} \mathbf{E}_{m,t}^{(Y)} \quad (19)$$

$$\gamma_{m,t} = P(m | \mathbf{X}_t, \mathbf{Y}_t, \lambda^{(z)})$$

Maximum likelihood parameter generation described in Equation 15 achieves a certain improvement of conversion performance since time correlations between adjacent frames are taken into account. However, dynamic features has different properties from those of static features. The joint vector contains both static and dynamic features, which tends to present characteristics of feature space with different properties together. In this concern, this model may not be so reasonable in considering correlations between adjacent frames.

### 3. MV-GMM for VC

In this section, voice conversion based on matrix variate Gaussian model is briefly explained [5]. Here,  $\mathbf{X}$  is a random matrix whose size is  $n \times p$ . When  $\mathbf{X}$  follows a normal distribution of matrix variate, it can be written as follow:

$$\mathbf{X} \sim \mathcal{N}_{mv}(\mathbf{X}; \mathbf{M}, \mathbf{U}, \mathbf{V}), \quad (20)$$

where  $\mathbf{M}$  is a matrix whose size is  $n \times p$ , representing the mean information of the normal distribution.  $\mathbf{U}$  and  $\mathbf{V}$  are matrices whose sizes are  $n \times n$  and  $p \times p$ , which represent the covariance structure of row and column spaces, respectively. A normal distribution of matrix variate corresponds to that of vector variate which is represented as follows [9]:

$$P(\text{vec}(\mathbf{X})|\boldsymbol{\lambda}) = \mathcal{N}(\text{vec}(\mathbf{X}); \text{vec}(\mathbf{M}), \mathbf{V} \otimes \mathbf{U}), \quad (21)$$

where  $\text{vec}()$  is the operator that change a matrix into a vector.  $\mathbf{V} \otimes \mathbf{U}$  means Kronecker product of the two matrices. According to Equation 21, the matrix variate normal distribution differs from the vector variate normal distribution. Compared with the vector variate normal distribution, it has a restricted covariance structure based on Kronecker product. Separating the covariance structure into two matrices  $\mathbf{U}$  and  $\mathbf{V}$  makes it possible to represent properties of the row and column spaces, respectively.

Normal distribution of matrix can be applied to voice conversion, in a similar way to that of normal distribution of vector. Let  $\mathbf{x}_t, \mathbf{y}_t$  be the feature vectors of the source and target speakers.  $\mathbf{x}_t, \mathbf{y}_t$  are combined into one joint matrix  $\mathbf{Z}_t = [\mathbf{x}_t, \mathbf{y}_t] \in \mathcal{R}^{D \times S}$ , where  $D$  means the dimension of feature space, and  $S$  the dimension of speaker space. In the case of one source speaker and one target speaker,  $S = 2$ . The probability density of  $\mathbf{Z}_t$  is represented by a mixture model shown as follows:

$$P(\mathbf{Z}_t|\boldsymbol{\lambda}^{(Z)}) = \sum_{m=1}^M w_m \mathcal{N}_{mv}(\mathbf{Z}_t; \mathbf{M}_m, \mathbf{U}_m, \mathbf{V}_m) \quad (22)$$

According to Equation 22, the probability density is defined as the weighted sum of normal distribution of each component. In each component, the matrix variate normal distribution is indicated by three matrix parameters  $\mathbf{M}_m, \mathbf{U}_m$ , and  $\mathbf{V}_m$ ;  $\mathbf{M}_m$  is the mean matrix,  $\mathbf{U}_m$  the covariance matrix of the feature space, and  $\mathbf{V}_m$  the covariance matrix capturing the correlation between the source and target. These parameters can be estimated by EM algorithm shown as follows:

$$\gamma_{m,t} = \frac{w_m \mathcal{N}_{mv}(\mathbf{Z}_t; \mathbf{M}_m, \mathbf{U}_m, \mathbf{V}_m)}{\sum_{m=1}^M w_m \mathcal{N}_{mv}(\mathbf{Z}_t; \mathbf{M}_m, \mathbf{U}_m, \mathbf{V}_m)} \quad (23)$$

$$\hat{\mathbf{M}}_m = \frac{1}{T_m} \sum_{t=1}^T \gamma_{m,t} \mathbf{Z}_t \quad (24)$$

$$\hat{\mathbf{U}}_m = \frac{1}{ST_m} \sum_{t=1}^T \gamma_{m,t} (\mathbf{Z}_t - \hat{\mathbf{M}}_m) \hat{\mathbf{V}}_m^{-1} (\mathbf{Z}_t - \hat{\mathbf{M}}_m)^\top \quad (25)$$

$$\hat{\mathbf{V}}_m = \frac{1}{DT_m} \sum_{t=1}^T \gamma_{m,t} (\mathbf{Z}_t - \hat{\mathbf{M}}_m)^\top \hat{\mathbf{U}}_m^{-1} (\mathbf{Z}_t - \hat{\mathbf{M}}_m) \quad (26)$$

$$T_m = \sum_{t=1}^T \gamma_{m,t} \quad (27)$$

A conversion function is derived based on the conditional probability density  $P(\mathbf{y}_t|\mathbf{x}_t)$ , referred to Equation 3. The conditional probability density of the  $m$ -th component is shown as follows:

$$P(\mathbf{y}_t|\mathbf{x}_t, m, \boldsymbol{\lambda}^{(Z)}) = \mathcal{N}(\mathbf{y}_t; \mathbf{E}_{m,t}; \mathbf{D}_m), \quad (28)$$

$$\mathbf{E}_{m,t} = \boldsymbol{\mu}_m^{(y)} + \frac{v_m^{(yx)}}{v_m^{(xx)}} (\mathbf{x}_t - \boldsymbol{\mu}_m^{(x)}), \quad (29)$$

$$\mathbf{D}_m = \left( v_m^{(yy)} - \frac{v_m^{(yx)} v_m^{(xy)}}{v_m^{(xx)}} \right) \mathbf{U}_m, \quad (30)$$

where  $v_m^{(\cdot)}$  shows elements of  $\mathbf{V}_m$ . Separating the covariance structure into two directions induces an effective estimation as Equations 25 and 26. Hence MVGMM-based VC achieves more concise and proper modeling.

#### 4. MVGMM VC with multiple frame features

In this section, we introduce features from multiple frames into MVGMM-based VC. Let  $\mathbf{x}_t, \mathbf{y}_t$  be feature vectors from utterances of the source and target speakers at frame  $t$ . In order to capture time correlations, features of the source ranging from  $\mathbf{x}_{t-N_x}$  to  $\mathbf{x}_{t+N_x}$  and those of the target ranging from  $\mathbf{y}_{t-N_y}$  to  $\mathbf{y}_{t+N_y}$  are focused on. The joint matrix  $\mathbf{Z}_t$  is defined as  $\mathbf{Z}_t = [\mathbf{x}_{t-N_x}, \dots, \mathbf{x}_t, \dots, \mathbf{x}_{t+N_x}, \mathbf{y}_{t-N_y}, \dots, \mathbf{y}_t, \dots, \mathbf{y}_{t+N_y}] \in \mathcal{R}^{D \times S}$ , where  $D$  denotes the dimension of the feature space, and where  $S$  denotes the total number of dimension considering multiple speakers and frames. The probability density of  $\mathbf{Z}_t$  modeled by MVGMM is described as the same equation of Equation 22. In this case,  $\mathbf{U}_m$ , i.e. the covariance structure for the feature space has the same dimension with that of when multiple frames are not taken into account. Only the dimension of  $\mathbf{V}_m$  increases as Equation 31;

$$\mathbf{V}_m = \begin{bmatrix} \mathbf{V}_m^{(xx)} & \mathbf{V}_m^{(xy)} \\ \mathbf{V}_m^{(yx)} & \mathbf{V}_m^{(yy)} \end{bmatrix}. \quad (31)$$

Two types of the conversion functions are derived according to whether multiple frames of the target are modeled. When a single frame feature of the target speaker is used, the conversion function is derived on each frame. The conditional probability density of  $\mathbf{y}_t$  given  $\mathbf{X}_t = [\mathbf{x}_{t-N_x}, \dots, \mathbf{x}_t, \dots, \mathbf{x}_{t+N_x}]$  is

$$P(\mathbf{y}_t|\mathbf{X}_t, \boldsymbol{\lambda}^{(Z)}) = \sum_{m=1}^M P(m|\mathbf{X}_t, \boldsymbol{\lambda}^{(Z)}) \times P(\mathbf{y}_t|\mathbf{X}_t, m, \boldsymbol{\lambda}^{(Z)}), \quad (32)$$

where

$$P(m|\mathbf{X}_t, \boldsymbol{\lambda}^{(Z)}) = \frac{w_m \mathcal{N}(\text{vec}(\mathbf{X}_t); \text{vec}(\mathbf{M}_m^{(x)}), \mathbf{V}_m^{(xx)} \otimes \mathbf{U}_m)}{\sum_{m=1}^M w_m \mathcal{N}(\text{vec}(\mathbf{X}_t); \text{vec}(\mathbf{M}_m^{(x)}), \mathbf{V}_m^{(xx)} \otimes \mathbf{U}_m)}, \quad (33)$$

$$P(\mathbf{y}_t|\mathbf{X}_t, m, \boldsymbol{\lambda}^{(Z)}) = \mathcal{N}(\mathbf{y}_t; \mathbf{E}_{m,t}^{(y)}, \mathbf{D}_m^{(y)}), \quad (34)$$

$$\mathbf{E}_{m,t}^{(y)} = \boldsymbol{\mu}_m^{(y)} + \mathbf{V}_m^{(yx)} \mathbf{V}_m^{(xx)-1} (\mathbf{X}_t - \boldsymbol{\mu}_m^{(x)}), \quad (35)$$

$$\mathbf{D}_m^{(y)} = \left( \mathbf{V}_m^{(yy)} - \mathbf{V}_m^{(yx)} \mathbf{V}_m^{(xx)-1} \mathbf{V}_m^{(xy)} \right) \mathbf{U}_m. \quad (36)$$

The target  $\hat{\mathbf{y}}_t$  can be generated by a similar way to Equation 8.

When multiple frames of the target speaker are taken into account, maximum likelihood sequence estimation [4] is adopted according to Equation 15, where the relation between  $\text{vec}(\mathbf{Y})$  and  $\mathbf{y}$  is shown in Figure 2 [6].

Comparing the proposed method with GMM-based model containing dynamic features, both of them are models considering correlations between adjacent frames. Instead of containing dynamic features which has different properties from static features, the proposed method contains also the static features of adjacent frames which shows similar properties to the current frame. The proposed method is regarded as a more reasonable model because all elements of the joint variate are static features instead of mixing static and dynamic features together.

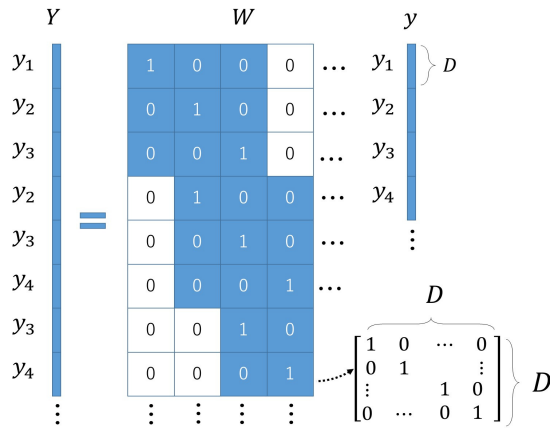


Figure 2: Relation between  $\text{vec}(Y)$  and  $y$

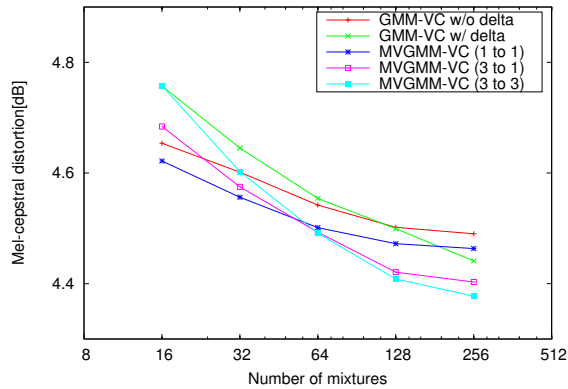


Figure 3: Results of objective evaluations

## 5. Experimental evaluation

### 5.1. Experimental setups

To evaluate the performance of the proposed method, objective and subjective evaluations were carried out. We used speech data of two male speakers (MMY as the source speaker, MHT as the target speaker) from ATR Japanese speech database [7]. 250 pairs of utterances were selected for training, and 50 pairs not included in the training data were selected for testing. As the proposed methods, we considered the method in which only multiple frame features from the source are modeled and the method in which multiple frame features both from the source and target are modeled. GMM-based approaches with/without dynamic features, and MVGMM-based one without multiple frames were compared with the proposed approaches. Diagonal structures were adopted for cross-covariance matrices in the conventional GMM-based methods, and full covariance matrices are used in the MV-GMM methods. 24-dimensional mel-cepstrum derived from STRAIGHT analysis [8] was used as feature vectors.

### 5.2. Objective evaluation

For the objective evaluations, Mel-cepstral distortion was used. From Figure 3, when the number of mixture components is larger than 64, the proposed methods outperformed the conventional ones. The performance was much better when multiple frames features is contained for both the source and tar-

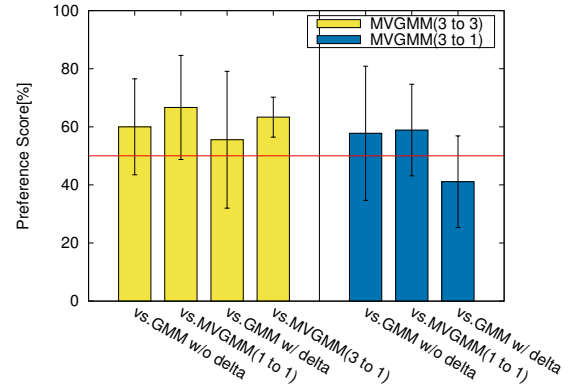


Figure 4: Results of subjective evaluations

get speakers (MVGMM-VC (3 to 3)). Besides, the proposed method (MVGMM-VC (3 to 3)) achieved better performance than the joint vector approach with dynamic features (GMM w/ delta). Note that both the methods take information of neighboring frames into account. This means that the separable structure in the proposed method has a positive effect in voice conversion.

### 5.3. Subjective evaluation

A subjective evaluation was processed in XAB method, where each listener was provided with utterances converted from two different kinds of conversion methods and he/she was asked to select the more similar one to the original speech of the target speaker. Ten Japanese native speakers aging from 20 to 35, each of whom were provided with 70 pairs of sentences in 7 groups of comparison, including comparison between the proposed and conventional methods, and among the proposed ones. Mixture of 256 was used for all conversion methods. For waveform generation,  $\log F_0$  values were linearly converted based on mean and variance values, and then STRAIGHT vocoder was adopted.

Figure 4 shows the result of the subjective evaluation. The error bars in the figure represent 95% confidence intervals. From Figure 4, the proposed method with multiple frame features of both the source and target (MVGMM-VC (3 to 3)) is more likely to be selected comparing to all other methods, while the proposed method with multiple frames feature only on the source side (MVGMM-VC (3 to 1)) performs better those not considering relationship between frames, however, it is not as preferable when compared with the joint vector approach with dynamic features (GMM w/ delta). This result suggests that modeling the multiple frames of both the source and target has a considerable influence on perception.

## 6. Conclusion

This paper presents a novel voice conversion method based on matrix variate Gaussian mixture model (MV-GMM) using features of multiple frames. The proposed method can effectively model features of utterances from the source and target speakers. Both objective and subjective evaluations show that the proposed method containing multiple frame features of both the source and target achieved a certain improvement of conversion performance, compared with the conventional methods.

## 7. Acknowledgment

This work was supported by MEXT KAKENHI Grant Number JP26118002 and JSPS KAKENHI Grant Number JP25730105.

## 8. References

- [1] A. Kain and M. Macon, "Spectral voice conversion for text-to-speech synthesis," IEEE International Conference on, 1998, pp. 285-288
- [2] S. Desai, E. V. Raghavendra, B. Yegnanarayana, A. W. Black and K. Prahallad, "Voice conversion using artificial neural networks," ICASSP, pp. 3893-3896 (2009).
- [3] Ryoichi Takashima, Tetsuya Takiguchi, and Yasuo Ariki, "Exemplar-based voice conversion in noisy environment," in SLT, pp. 313-317, (2012).
- [4] Tomoki Toda, Alan W Black, and Keiichi Tokuda, "Voice Conversion Based on Maximum-Likelihood Estimation of Spectral Parameter Trajectory," IEEE TRANSACTIONS 15 (8) pp. 2222-2235 (2007).
- [5] Daisuke Saito, Hidenobu Doi, Nobuaki Minematsu, and Keikichi Hirose, "Application of Matrix Variate Gaussian Mixture Model to Statistical Voice Conversion," INTER-SPEECH 2014
- [6] Ling-Hui Chen, Zhen-Hua Ling, and Li-Rong Dai, "Voice Conversion Using Generative Trained Deep Neural Networks with Multiple Frame Spectral Envelopes," Inter-speech, pp. 2313-2317, (2014).
- [7] Akira Kurematsu, Kazuya Takeda, Yoshinori Sagisaka, Shigeru Katagiri, Hisao Kuwabara, and Kiyohiro Shikano "ATR Japanese speech database as a tool of speech recognition and synthesis," Speech Communication, vol.9, pp.357-363, (1990).
- [8] H. Kawahara, I. Masuda-Katsuse, and A.de Cheveigne, "Re-structuring speech representations using a pitch-adaptive time- frequency smoothing and an instantaneous-frequency-based F0 extraction: Possible role of a repetitive structure in sounds," Speech Communication, vol.27, pp.187-207, 1999.
- [9] A. K. Gupta, D. K. Nagar, Matrix Variate Distributions, 2000