



Statistical Acoustic-to-Articulatory Mapping Unified with Speaker Normalization Based on Voice Conversion

Hidetsugu Uchida, Daisuke Saito, Nobuaki Minematsu, Keikichi Hirose

The University of Tokyo, 7-3-1, Hongo, Bunkyo-ku, Tokyo, 113-8656 Japan
 {uchida,dsk_saito,mine,hirose}@gavo.t.u-tokyo.ac.jp

Abstract

This paper proposes a model of speaker-normalized acoustic-to-articulatory mapping using statistical voice conversion. A mapping function from acoustic parameters to articulatory parameters is usually developed with a single speaker’s parallel data. Hence the constructed mapping model can work appropriately only for this specific speaker, and applying this model to other speakers degrades the performance of acoustic-to-articulatory mapping. In this paper, two models of speaker conversion and acoustic-to-articulatory mapping are implemented using Gaussian Mixture Models (GMM), and by integrating these two models, we propose two methods of speaker-normalized acoustic-to-articulatory mapping. One is concatenating these models sequentially, and the other integrates the two models into a unified model, where acoustic parameters of a speaker can be converted directly to articulatory parameters of another speaker. Experiments show that both methods can improve the mapping accuracy and that the latter method works better than the former method. Especially in the case of velar stop consonants, the mapping accuracy is higher by 0.6 mm.

Index Terms: acoustic-to-articulatory mapping, Gaussian mixture model, voice conversion, speaker normalization

1. Introduction

A new type of pronunciation training system using information of articulatory movements has been studied [1]. The training system provides learners with visual feedback of articulatory movements both of the learners and teachers and it can improve the learners’ pronunciation more effectively than a conventional system that uses only audio features [2].

To implement such systems, it is necessary to establish a method to estimate learners’ articulatory movements from their utterances. Articulatory movements can be measured directly when special instruments are available, e.g. electromagnetic articulography [3] and electropalatography [4][5]. Since they are large-scale instruments, however, it is surely difficult to incorporate them into personalized pronunciation training. These days, methods to estimate articulatory movements from voices have been studied using a parallel corpus of acoustic and articulatory observation. In these studies, several mapping models have been proposed, e.g. Gaussian Mixture Model (GMM) [6], Hidden Markov Model (HMM) [7], and Deep Neural Network (DNN) [8]. A mapping model developed with acoustic-articulatory parallel data of a single speaker has to be a speaker-dependent model. If voices of a different speaker are input to the mapping model, acoustic and articulatory mismatch between the input and the model causes large errors in mapping. The problem can be avoided by constructing a model for each individual with his/her own parallel data. However, this solution is impractical because simultaneous recording of acoustic and articulatory data is expensive. We need a method that can estimate articulatory movements of arbitrary speakers only by using a specific speaker’s acoustic-to-articulatory mapping model.

Statistical mapping methods are widely used in the voice conversion field [9]. Speaker conversion, which can be viewed

as applied voice conversion, is a technique which can convert speaker identity of input voices to a target speaker without changing linguistic content. In this paper, a speaker conversion technique is used to conduct speaker normalization, where speaker identity of input voices is changed to the speaker used for constructing the acoustic-to-articulatory mapping model. For open speakers, mapping accuracy will be improved.

In this study, we investigate two methods which can estimate articulatory movements from arbitrary speakers’ voices. One is concatenating a GMM-based speaker conversion model and a GMM-based acoustic-to-articulatory mapping model sequentially. The other integrates the models into a unified model.

2. Speaker normalization based on voice conversion

2.1. GMM-based parameter mapping

In this paper, both models of speaker conversion and acoustic-to-articulatory mapping are implemented as GMM-based parameter conversion, which is briefly described below.

Let $\mathbf{x} \in \mathcal{R}^{d_x}$ and $\mathbf{y} \in \mathcal{R}^{d_y}$ be source and target parameter vectors whose dimensions are d_x and d_y , respectively. \mathbf{z} denotes a joint vector consisting of source and target parameters as $\mathbf{z} = [\mathbf{x}^\top, \mathbf{y}^\top]^\top$. The probability density of the joint vector is modeled by using a GMM as follows:

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

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

$\boldsymbol{\lambda}$ is model parameters. $\mathcal{N}(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is a normal distribution with a mean vector $\boldsymbol{\mu}$ and a covariance matrix $\boldsymbol{\Sigma}$. M is the total number of mixture components and α is a weight parameter. The parameter mapping function using the GMM is derived from

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}; \boldsymbol{\lambda}^{(z)}). \quad (3)$$

A GMM-based speaker conversion model is developed with acoustic-acoustic parallel data, where source and target parameters are speech features of source and target speakers, respectively. A GMM-based acoustic-to-articulatory mapping model is developed with acoustic-articulatory parallel data, where a source parameter is a speech feature and a target parameter is an articulatory parameter. In this paper, the speaker used for the acoustic-to-articulatory mapping model is called *model* speaker henceforth. Then, GMM-based speaker conversion is done to change speaker identity of input voices to the model speaker.

Since GMM-based conversion is a frame-to-frame mapping, however, it should be noted that the temporal structure of input utterances is copied onto output utterances¹. Considering this technical fact, the converted speech can be regarded

¹Even if ML-based estimation is employed, the number of frames is the same between the input and the output [6].

as speech generated by the model speaker with a speaking style similar to the input speaker's style. If the speaker-normalized voices are input to the acoustic-to-articulatory mapping model, the output is the model speaker's articulatory movements generated in a speaking style similar to the input speaker's style. In other words, we can regard the output as the model speaker's imitation of input utterances. In the case of pronunciation training, we can assume that an input speaker is a learner and the model speaker is a teacher. Therefore, the estimated articulatory movements are regarded as the teacher's articulation when he/she imitates the learner's utterances. A teacher often imitates a learner's utterances to show clearly what is wrong with the utterances. This imitation is effective to increase learners' awareness of pronunciation. Our proposed method is expected to work as teacher, where the feedback of imitated articulatory movements will be effective for pronunciation training.

2.2. Integrating GMM-based speaker conversion and acoustic-to-articulatory mapping

In this study, we integrate the speaker conversion model and the acoustic-to-articulatory mapping model when two types of parallel data are prepared: one is acoustic-articulatory parallel data of the model speaker, and the other is acoustic-acoustic parallel data between that model speaker and an arbitrary speaker. For integrating these two models, two methods are proposed: a concatenation method and a unification method.

2.3. Concatenation method

This method concatenates two models of speaker conversion and acoustic-to-articulatory mapping sequentially. Let $\mathbf{x}^{(s)}$, $\mathbf{y}^{(s)}$, and $\mathbf{y}^{(a)}$ be a speech vector of the input speaker, that of the model speaker, and an articulatory vector of the model speaker, respectively. The probability density functions of two joint vectors $\mathbf{z}^{(xy)} = [\mathbf{x}^{(s)\top}, \mathbf{y}^{(s)\top}]^\top$ and $\mathbf{z}^{(sa)} = [\mathbf{y}^{(s)\top}, \mathbf{y}^{(a)\top}]^\top$ are characterized respectively as $P(\mathbf{z}^{(xy)}; \boldsymbol{\lambda}^{(xy)}) = \sum_{m=1}^M \alpha_m \mathcal{N}(\mathbf{z}^{(xy)}; \boldsymbol{\mu}_m^{(xy)}, \boldsymbol{\Sigma}_m^{(xy)})$ and $P(\mathbf{z}^{(sa)}; \boldsymbol{\lambda}^{(sa)}) = \sum_{n=1}^N \beta_n \mathcal{N}(\mathbf{z}^{(sa)}; \boldsymbol{\mu}_n^{(sa)}, \boldsymbol{\Sigma}_n^{(sa)})$.

When applying the resulting model to acoustic-to-articulatory mapping, $\mathbf{x}^{(s)}$ is firstly converted to $\hat{\mathbf{y}}^{(s)}$ by Eq. 3, where $\mathbf{x} = \mathbf{x}^{(s)}$ and $\mathbf{y} = \mathbf{y}^{(s)}$. Then, $\hat{\mathbf{y}}^{(s)}$ is converted to articulatory movements by Eq. 3 again, where $\mathbf{x} = \mathbf{y}^{(s)}$ and $\mathbf{y} = \mathbf{y}^{(a)}$.

2.4. Unification method

We propose another method that can integrate the two models into a unified model, where input voices are converted directly to articulatory movements of the model speaker.

In the concatenation method, one can find that two probability densities of $P(\mathbf{z}^{(xy)}; \boldsymbol{\lambda}^{(xy)})$ and $P(\mathbf{z}^{(sa)}; \boldsymbol{\lambda}^{(sa)})$ include a common distribution of $\mathbf{y}^{(s)}$ ². Practically speaking, distribution of $\mathbf{y}^{(s)}$ in $P(\mathbf{z}^{(xy)}; \boldsymbol{\lambda}^{(xy)})$ and that in $P(\mathbf{z}^{(sa)}; \boldsymbol{\lambda}^{(sa)})$ may be different when parallel data of $(\mathbf{x}^{(s)}, \mathbf{y}^{(s)})$ and those of $(\mathbf{y}^{(s)}, \mathbf{y}^{(a)})$ are obtained by reading different sentences, for example. By assuming that the two kinds of distributions of $\mathbf{y}^{(s)}$ can be characterized by a single model even in such a case, the total number of parameters that have to be estimated for acoustic-to-articulatory mapping can be reduced [10]. Furthermore, by marginalizing the parameters of $P(\mathbf{y}^{(s)})$, the resulting mapping model becomes a unified model that can convert input voices to articulatory movements of the model speaker directly.

Input voices of speaker x are denoted as $\mathbf{x}^{(s)}$. In the unifi-

² $\mathbf{z}^{(xy)} = [\mathbf{x}^{(s)\top}, \mathbf{y}^{(s)\top}]^\top$ and $\mathbf{z}^{(sa)} = [\mathbf{y}^{(s)\top}, \mathbf{y}^{(a)\top}]^\top$

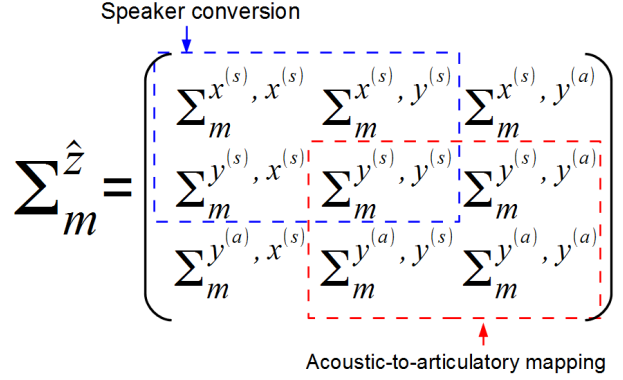


Figure 1: The covariance matrix of a component of GMM for the joint vector $\hat{\mathbf{z}}$

cation method, articulatory movements are obtained by

$$\hat{\mathbf{y}}^{(a)} = \arg \max_{\mathbf{y}^{(a)}} P(\mathbf{y}^{(a)} | \mathbf{x}^{(s)}; \boldsymbol{\lambda}). \quad (4)$$

The above equation can be approximated in the following way:

$$\hat{\mathbf{y}}^{(a)} = \arg \max_{\mathbf{y}^{(a)}} P(\mathbf{y}^{(a)} | \mathbf{x}^{(s)}) \quad (5)$$

$$\approx \arg \max_{\mathbf{y}^{(a)}} \sum_{m=1}^M P(m | \mathbf{x}^{(s)}) \times \int_{\mathbf{y}^{(s)}} P(\mathbf{y}^{(a)} | \mathbf{y}^{(s)}, m) P(\mathbf{y}^{(s)} | \mathbf{x}^{(s)}, m) d\mathbf{y}^{(s)} \quad (6)$$

$$= \arg \max_{\mathbf{y}^{(a)}} \sum_{m=1}^M P(m | \mathbf{x}^{(s)}) \mathcal{N}(\mathbf{y}^{(a)} | \mathbf{E}_m, \mathbf{D}_m), \quad (7)$$

where

$$\mathbf{E}_m = \boldsymbol{\mu}_m^{(a)} + \boldsymbol{\Sigma}'_m (\mathbf{x}^{(s)} - \boldsymbol{\mu}_m^{(x)}), \quad (8)$$

$$\mathbf{D}_m = \boldsymbol{\Sigma}_m^{(a,a)} - \boldsymbol{\Sigma}'_m \boldsymbol{\Sigma}_m^{(x,x)} \boldsymbol{\Sigma}'_m{}^\top, \quad (9)$$

$$\boldsymbol{\Sigma}'_m = \boldsymbol{\Sigma}_m^{(a,s)} \boldsymbol{\Sigma}_m^{(y,y)^{-1}} \boldsymbol{\Sigma}_m^{(y,x)}. \quad (10)$$

By using Eq. 7, we can estimate articulatory movements of the model speaker from the input speaker's voices without converting them to the model speaker's voices explicitly.

2.5. Implementation of the unification method

The unification method requires that each GMM component of $\mathbf{y}^{(s)}$ should be the same between the speaker conversion model and the acoustic-to-articulatory mapping model. To satisfy this condition, we first build a GMM of a new joint vector of $\hat{\mathbf{z}} = [\mathbf{x}^{(s)\top}, \mathbf{y}^{(s)\top}, \mathbf{y}^{(a)\top}]^\top$. Fig. 1 illustrates a covariance matrix in the GMM of $\hat{\mathbf{z}}$, where the covariance matrix of $\mathbf{z}^{(xy)}$ and that of $\mathbf{z}^{(sa)}$ coexist, between which the distribution of $\mathbf{y}^{(s)}$ is shared. Namely, the two models of speaker conversion and acoustic-to-articulatory mapping share the distribution of $\mathbf{y}^{(s)}$.

The simplest method to build a GMM of $\hat{\mathbf{z}}$ is to estimate its parameters by using samples of $(\mathbf{x}^{(s)}, \mathbf{y}^{(s)}, \mathbf{y}^{(a)})$. In the context of pronunciation training, however, this method is not practical. The speaker conversion model has to be made by a parallel data of the same language, namely, a learner's native language³. The acoustic-articulatory mapping model should be made by a parallel data of the target language of learning. These conditions require the model speaker to be a bilingual speaker

³Learners are not good speakers of the target language.

of both languages and under these conditions, simultaneous observation of $\mathbf{x}^{(s)}$, $\mathbf{y}^{(s)}$, and $\mathbf{y}^{(a)}$ is practically impossible.

To build a GMM of $\hat{\mathbf{z}}$ by using two non-overlapped parallel databases of $(\mathbf{x}^{(s)}, \mathbf{y}^{(s)})$ and $(\mathbf{y}^{(s)}, \mathbf{y}^{(a)})$, we have to derive a method to build a GMM by using $(\mathbf{x}^{(s)}, \mathbf{y}^{(s)}, \mathbf{y}^{(a)})$, where either of $\mathbf{x}^{(s)}$ or $\mathbf{y}^{(a)}$ is always missing. For derivation, joint vector $\hat{\mathbf{z}}$ is re-introduced as $\mathbf{z}_t^{(xyy')} = [\mathbf{x}_t^{(s)\top}, \mathbf{y}_t^{(s)\top}, \mathbf{y}_t^{(a)\top}]^\top$ when $t = 1, \dots, T_1$ and $\mathbf{z}_t^{(s'sa)} = [\mathbf{x}_t^{(s)\top}, \mathbf{y}_t^{(s)\top}, \mathbf{y}_t^{(a)\top}]^\top$ when $t = T_1 + 1, \dots, T_2$. Here, $\mathbf{y}_t^{(a)}$ is a missing articulatory vector of the model speaker and $\mathbf{x}_t^{(s)}$ is a missing speech vector of the input speaker. The E-step of the general EM algorithm for GMM of $\hat{\mathbf{z}}$ is shown below.

$$\gamma_{m,t}^{(xyy')} = \frac{\pi_m \mathcal{N}(\mathbf{z}_t^{(xyy')} | \boldsymbol{\mu}_m^{(\hat{\mathbf{z}})}, \boldsymbol{\Sigma}_m^{(\hat{\mathbf{z}})})}{\sum_{j=1}^M \pi_j \mathcal{N}(\mathbf{z}_t^{(xyy')} | \boldsymbol{\mu}_j^{(\hat{\mathbf{z}})}, \boldsymbol{\Sigma}_j^{(\hat{\mathbf{z}})})} \quad (11)$$

$$\gamma_{m,t}^{(s'sa)} = \frac{\pi_m \mathcal{N}(\mathbf{z}_t^{(s'sa)} | \boldsymbol{\mu}_m^{(\hat{\mathbf{z}})}, \boldsymbol{\Sigma}_m^{(\hat{\mathbf{z}})})}{\sum_{j=1}^M \pi_j \mathcal{N}(\mathbf{z}_t^{(s'sa)} | \boldsymbol{\mu}_j^{(\hat{\mathbf{z}})}, \boldsymbol{\Sigma}_j^{(\hat{\mathbf{z}})})} \quad (12)$$

Since $\mathbf{y}_t^{(a)}$ and $\mathbf{x}_t^{(s)}$ are missing, they have to be replaced by some appropriate terms. Here, their expected values are used.

$$\mathbf{y}'_t^{(a)} \leftarrow \sum_{m=1}^M P(m | \boldsymbol{\mu}_m^{(xy)}, \boldsymbol{\Sigma}_m^{(xy)}) E_{m,t}^{(y^{(a)} | x^{(s)}, y^{(s)})}. \quad (13)$$

$$\mathbf{x}'_t^{(s)} \leftarrow \sum_{m=1}^M P(m | \boldsymbol{\mu}_m^{(sa)}, \boldsymbol{\Sigma}_m^{(sa)}) E_{m,t}^{(x^{(s)} | y^{(s)}, y^{(a)})}. \quad (14)$$

$\boldsymbol{\mu}_m^{(xy)}$ and $\boldsymbol{\Sigma}_m^{(xy)}$ are estimated from samples of joint vector $(\mathbf{x}^{(s)}, \mathbf{y}^{(s)})$. Similarly, $\boldsymbol{\mu}_m^{(sa)}$ and $\boldsymbol{\Sigma}_m^{(sa)}$ are estimated from samples of joint vector $(\mathbf{y}^{(s)}, \mathbf{y}^{(a)})$. $E_{m,t}^{(y^{(a)} | x^{(s)}, y^{(s)})}$ and $E_{m,t}^{(x^{(s)} | y^{(s)}, y^{(a)})}$ are obtained as follows:

$$E_{m,t}^{(y^{(a)} | x^{(s)}, y^{(s)})} = \boldsymbol{\mu}_m^{(a)} + \boldsymbol{\Sigma}_m^{(a,xy)} \boldsymbol{\Sigma}_m^{(xy,xy)^{-1}} (\mathbf{z}_t^{(xy)} - \boldsymbol{\mu}_m^{(xy)}) \quad (15)$$

$$E_{m,t}^{(x^{(s)} | y^{(s)}, y^{(a)})} = \boldsymbol{\mu}_m^{(s)} + \boldsymbol{\Sigma}_m^{(s,sa)} \boldsymbol{\Sigma}_m^{(sa,sa)^{-1}} (\mathbf{z}_t^{(sa)} - \boldsymbol{\mu}_m^{(sa)}) \quad (16)$$

$\boldsymbol{\mu}_m$, $\boldsymbol{\Sigma}_m$, and π_m are updated in the M-step below.

$$\boldsymbol{\mu}_m = \frac{1}{\gamma_m} \left(\sum_{t=1}^{T_1} \gamma_{m,t}^{(xyy')} \mathbf{z}_t^{(xyy')} + \sum_{t=T_1+1}^{T_2} \gamma_{m,t}^{(s'sa)} \mathbf{z}_t^{(s'sa)} \right) \quad (17)$$

$$\boldsymbol{\Sigma}_m = \frac{1}{\gamma_m} \left(\sum_{t=1}^{T_1} \gamma_{m,t}^{(xyy')} \left\{ (\mathbf{z}_t^{(xyy')} - \boldsymbol{\mu}_m) (\mathbf{z}_t^{(xyy')} - \boldsymbol{\mu}_m)^\top - \hat{\mathbf{D}}_m^{(xyy')} \right\} + \sum_{t=T_1+1}^{T_2} \gamma_{m,t}^{(s'sa)} \left\{ (\mathbf{z}_t^{(s'sa)} - \boldsymbol{\mu}_m) (\mathbf{z}_t^{(s'sa)} - \boldsymbol{\mu}_m)^\top - \hat{\mathbf{D}}_m^{(s'sa)} \right\} \right) \quad (18)$$

$$\pi_m = \frac{\gamma_m}{\sum_{m=1}^M \gamma_m} \quad (19)$$

where $\gamma_m = \sum_{t=1}^{T_1} \gamma_{m,t}^{(xyy')} + \sum_{t=T_1+1}^{T_2} \gamma_{m,t}^{(s'sa)}$. In the M-step, Eq.15 and Eq.16 are substituted again for missing $\mathbf{y}_t^{(a)}$ and $\mathbf{x}_t^{(s)}$. $\hat{\mathbf{D}}_m^{(xyy')}$ and $\hat{\mathbf{D}}_m^{(s'sa)}$ in Eq.18 are obtained as

$$\hat{\mathbf{D}}_m^{(xyy')} = \begin{bmatrix} \mathbf{0}^{(d_1, d_1)} & \mathbf{0}^{(d_1, d_2)} \\ \mathbf{0}^{(d_2, d_1)} & \mathbf{D}_m^{(y^{(a)} | x^{(s)}, y^{(s)})} \end{bmatrix} \quad (20)$$

$$\hat{\mathbf{D}}_m^{(s'sa)} = \begin{bmatrix} \mathbf{D}_m^{(x^{(s)} | y^{(s)}, y^{(a)})} & \mathbf{0}^{(d_3, d_4)} \\ \mathbf{0}^{(d_4, d_3)} & \mathbf{0}^{(d_4, d_4)} \end{bmatrix}, \quad (21)$$

where

$$\mathbf{D}_m^{(y^{(a)} | x^{(s)}, y^{(s)})} = \boldsymbol{\Sigma}_m^{(a,a)} - \boldsymbol{\Sigma}_m^{(a,xy)} \boldsymbol{\Sigma}_m^{(xy,xy)^{-1}} \boldsymbol{\Sigma}_m^{(xy,a)}, \quad (22)$$

$$\mathbf{D}_m^{(x^{(s)} | y^{(s)}, y^{(a)})} = \boldsymbol{\Sigma}_m^{(s,s)} - \boldsymbol{\Sigma}_m^{(s,sa)} \boldsymbol{\Sigma}_m^{(sa,sa)^{-1}} \boldsymbol{\Sigma}_m^{(sa,s)}. \quad (23)$$

$\mathbf{0}^{(m,l)}$ is an $m \times l$ zero matrix, and d_1, d_2, d_3 , and d_4 are dimensions of $\mathbf{z}^{(xy)}$, $\mathbf{y}^{(a)}$, $\mathbf{z}^{(x)}$, and $\mathbf{z}^{(sa)}$, respectively.

3. Experimental evaluations

3.1. The database used for experiments

To evaluate the performance of our proposal of speaker-normalized acoustic-to-articulatory mapping model, we conducted experimental evaluations using MOCHA database [11]. This database includes acoustic-articulatory parallel data of one male speaker and one female speaker. In the experiments, either of the two was used as model speaker and the other was an input speaker. Articulatory data are measured by an electromagnetic articulography, where its sensors are placed at lower incisor (LI), upper and lower lips (UL and LL), tongue tip (TT), tongue body (TB), tongue dorsum (TD) and velum (V) in the med-sagittal plane. The articulatory data of each point are two-dimensional data of horizontal (x-axis) and vertical (y-axis) directions. The parallel data were divided into 5 parts, where 4 parts were used for training and the other was used for testing.

3.2. Conditions

For the experiments, three kinds of features were extracted, which are related to acoustic and articulatory observation at time t : 1) 0-24th MFCCs and their deltas, 2) the first 75 principal components obtained from PCA of 11 frames of MFCCs centering at time t [6], and 3) 28-dimensional data consisting of 14-dimensional articulatory parameters (2-dimensional data obtained from the 7 sensors) and their deltas. It should be noted that, by following [6], acoustic-to-articulatory conversion was implemented from 2) to 3), not from 1) to 3).

In the concatenation method, firstly, conversion of MFCCs of an input speaker to MFCCs of the model speaker was carried out. After applying PCA, acoustic-to-articulatory mapping was done. On the other hand in the unification method, joint vectors of PCA-based features of an input speaker, those of the model speaker, and articulatory parameters of the model speaker were formed to estimate their GMM, where some values are missing as explained in section 2.5. Then, the PCA-based features of the input speaker was converted directly to the articulatory parameters of the model speaker. For each conversion process, ML-based estimation using dynamic features was adopted. For time alignment between two sequences, we employed DTW.

For comparison, we built another acoustic-to-articulatory mapping model using MOCHA's acoustic-articulatory parallel data of the model speaker. This model is called as reference model. We conducted mapping experiments in four conditions:

1. reference: the reference model is used, where both the input and output speakers are the model speaker.
2. baseline: the reference model is used again but the input speaker is different from the model (output) speaker.
3. concatenation: the speaker-normalized model based on concatenation is used, where the input speaker is different from the model (output) speaker.
4. unification: the speaker-normalized model based on unification is used, where the input speaker is different from the model (output) speaker.

Five-fold cross-validation tests were employed. For that, MOCHA database of two speakers was divided into several

5. References

- [1] P. Badin, A. Youssef, G. Bailly, F. Elisei, and T. Hueber, "Visual articulatory feedback for phonetic correction in second language learning," In *L2SW, Workshop on "Second Language Studies: Acquisition, Learning, Education and Technology."* pp. 1-10, 2010.
- [2] A. Suemitsu, J. Dang, T. Ito, M. Tiede, "A study on effect of real-time articulatory feedback presentation in American English pronunciation learning," In *Proc. Acoustic society of Japan Autumn Meeting.* pp. 427-428, 2013.
- [3] T. Kaburagi, K. Wakamiya, and M. Honda, "Three-dimensional electromagnetic articulography," *Journal of the Acoustical Society of America*, vol. 118, pp. 428-443, 2005
- [4] W. Hardcastle, W. Jones and C. Knight, "New developments in electropalatography: A state-of-the-art report," *Clinical Linguistics & Phonetics*, vol. 3, No. 1, pp. 1-38, 1989.
- [5] A. Wrench, F. Gibbon, A. M. McNeill, and S. Wood, "An EPG therapy protocol for remediation and assessment of articulation disorders," In *Proc. ICSLP*, pp. 965-968, 2002.
- [6] T. Toda, W. A. Black, and K. Tokuda, "Statistical mapping between articulatory movements and acoustic spectrum using a Gaussian mixture model," *Speech Commun.*, vol. 50, pp. 215 -227, 2008.
- [7] S. Hiroya, and M. Honda, "Estimation of articulatory movements from speech acoustics using an HMM-based speech production model," *IEEE Trans. Speech and Audio Processing*, vol. 12, pp. 175-185, 2004.
- [8] B. Uria, S. Renals, and K. Richmond, "A deep neural network for acoustic-articulatory speech inversion," In *Proc. NIPS 2011 Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [9] Y. Stylianou, O. Cappé, and E. Moulines, "Continuous probabilistic transform for voice conversion," *IEEE Trans. Speech Audio Process.*, vol. 6, No. 2, pp. 131-142, 1998.
- [10] Y. Ohtani, T. Toda, H. Saruwatari, and K. Shikano, "Many-to-many eigenvoice conversion with reference voice," In *Proc. INTERSPEECH*, pp. 1623-1626, 2009
- [11] MOCHA-TIMIT - Centre for Speech Technology Research, 5 June 2014, "<http://www.cstr.ed.ac.uk/research/projects/artic/mocha.html>"