

## ADAPTIVE CONDITIONAL PRONUNCIATION MODELING USING ARTICULATORY FEATURES FOR SPEAKER VERIFICATION

*Ka-Yee Leung<sup>1</sup>, Man-Wai Mak<sup>1</sup>, Manhung Siu<sup>2</sup>, Sun-Yuan Kung<sup>3</sup>*

<sup>1</sup>Center for Multimedia Signal Processing, Dept. of Electronic and Information Engineering,  
The Hong Kong Polytechnic University

<sup>2</sup>Dept. of Electrical and Electronic Engineering, Hong Kong University of Science and Technology

<sup>3</sup>Dept. of Electrical Engineering, Princeton University

### ABSTRACT

This paper proposes an articulatory feature-based conditional pronunciation modeling (AFCPM) technique for speaker verification. The technique models the pronunciation behaviors of speakers by creating a link between the actual phones produced by the speakers and the state of articulations during speech production. Speaker models consisting of conditional probabilities of two articulatory classes are adapted from a set of universal background models (UBMs) using MAP adaptation technique. This adaptation approach aims to prevent over-fitting the speaker models when the amount of speaker data is insufficient for a direct estimation. Experimental results show that the adaptation technique can enhance the discriminating power of speaker models by establishing a tighter coupling between speaker models and the UBMs. Results also show that fusing the scores derived from an AFCPM-based system and a conventional spectral-based system achieves a significantly lower error rate than that of the individual systems. This suggests that AFCPM and spectral features are complementary to each other.

### 1. INTRODUCTION

State-of-the-art text-independent speaker recognition systems typically use Gaussian mixture models (GMMs) [1] to represent the short-term spectral characteristics of speakers. The advantage of spectral-based systems is that promising results are obtainable from a limited amount of training data. However, except for spectral characteristics, these systems ignore other information in speech signals which is useful for human to recognize speakers.

In recent years, researchers have started to investigate the use of high-level features, such as the usage or duration of particular words, prosodic features, etc., for speaker recognition [2]. Their work has demonstrated that these features contain different amount of speaker-dependent information and the best performance was achieved by a system that uses conditional pronunciation modeling (CPM) techniques [3]. Because different speakers have different ways of pronouncing the same phoneme, CPM aims to characterize the pronunciation behaviors of a speaker by computing the correlation between the intended phonemes and the actual phones. The pronunciation behaviors were encoded as discrete probability densities that were used for verifying speakers similar to the

conventional GMMs in spectral-based systems. However, CPM requires multilingual speech data for training the phone models of different languages and long utterances for speaker enrollment and verification.

To avoid the requirement of multilingual training data, Leung et al. [4] proposed using articulatory feature (AF) streams to construct conditional pronunciation models. AFs are abstract classes describing the movements or positions of different articulators during speech production [5]. Compared to phone-based CPM in [3], AF-based CPM provides a more direct coupling between the pronunciation variations and the speech production process. Because the speech production process is a source of speaker variations, AF-based CPM is better than phone-based CPM in terms of speaker modeling. In addition, articulatory properties are the same irrespective of languages, therefore monolingual speech data are sufficient for determining their values. In Leung et al. [4], significantly shorter utterances were used to enroll and verify speakers when compared to those required in Klusáček et al. [3]. This has important computation implication for large-scale deployment.

In Leung et al. [4], the discrete distribution of each speaker model was estimated exclusively from the enrollment data of the corresponding speaker. This may lead to over-trained speaker models unless abundant enrollment data are available. To solve this problem, this paper proposes an adaptation approach in which the discrete distributions of speaker models are adapted from those of universal background models.

### 2. AF-BASED CPM

This section details the notion of articulatory features and explains how AFs can be applied to model the pronunciation characteristics of speakers.

#### 2.1. Articulatory Features

AFs are the representations of some important phonological properties appeared during speech production. More precisely, AFs are abstract classes describing the movements or positions of different articulators during speech production. AFs have been applied to speaker identification [6] and speaker verification [7]. In [6], speaker identification was performed by fusing the scores derived from seven speaker-dependent language models, each of which modeled the classes of a single articulatory property by a discrete conditional distribution. For each utterance, seven articulatory class sequences were obtained from seven HMM-based rec-

This work was supported by The Hong Kong Polytechnic University, Grant No. A-PE44 and Research Grant Council of the Hong Kong SAR (Project No. CUHK 1/02C).

Articulatory properties	Classes	Number of Classes
Manner ( $\mathcal{M}$ )	Silence, Vowel, Stop, Fricative, Nasal, Approximant-Lateral	6
Place ( $\mathcal{P}$ )	Silence, High, Middle, Low, Labial, Dental, Coronal, Palatal, Velar, Glottal	10

**Table 1.** Articulatory properties and the number of classes in each property.

ognizers, each responsible for one articulatory property. The usefulness of AFs in speaker verification was demonstrated in Leung et al. [7], where for each utterance, the probabilities of 26 articulatory classes determined from five multilayer perceptrons (MLPs) were concatenated to form a sequence of articulatory feature vectors. The AF sequence was then fed to a GMM speaker model and a background model to compute a likelihood ratio for decision making.

## 2.2. Articulatory Feature Extraction

The AF extraction approach outlined in [4] was adopted. According to [4], only two articulatory properties, (i.e., the manner and place of articulations listed in Table 1) were used for pronunciation modeling.

The AF-MLPs take  $n$  consecutive frames of Mel-frequency cepstral coefficients (MFCCs)  $X_t$  (with consecutive frame indexes ranging from  $t - \frac{n}{2}$  to  $t + \frac{n}{2}$ ) as inputs at frame  $t$ . For a given  $X_t$ , the outputs of the two AF-MLPs,  $P(\text{Manner} = m|X_t)$  and  $P(\text{Place} = p|X_t)$ , represent the posterior probabilities of different classes in the manner and place of articulation. The manner class label  $l_t^M \in \mathcal{M}$  and the place class label  $l_t^P \in \mathcal{P}$  (the sets of  $\mathcal{M}$  and  $\mathcal{P}$  are listed in Table 1) at frame  $t$  are determined by

$$l_t^M = \arg \max_{m \in \mathcal{M}} P(\text{Manner} = m|X_t) \quad \text{and} \quad (1)$$

$$l_t^P = \arg \max_{p \in \mathcal{P}} P(\text{Place} = p|X_t). \quad (2)$$

The two AF streams—one from the manner MLP and another from the place MLP—for creating the conditional pronunciation models are formed by concatenating  $l_t^M$ 's and  $l_t^P$ 's from  $t = 1, \dots, T$ , where  $T$  is the total number of frames in the utterance.

## 2.3. Speaker Modeling

AF-based CPM (hereafter, referred to as AF-CPM) aims to establish a relationship between the articulatory classes and the actual phonemes obtained from a phoneme-based recognizer. Because different speakers have different ways of pronunciation, their articulatory properties of the same phoneme can be varied.

### 2.3.1. Universal background models

For each phoneme, a set of universal background models (UBMs) is trained from the speech of a large number of speakers to represent the speaker-independent pronunciation characteristics corresponding to that phoneme. Each UBM comprises the joint probabilities of the manner and place classes conditioned on a phoneme. The training procedure begins with aligning two AF streams obtained from the AF-MLPs and a phoneme sequence obtained from

a null-grammar recognizer. For a particular phoneme  $q$ , the joint probabilities of the corresponding UBM are determined by

$$P_{bg}(\text{Manner} = m, \text{Place} = p | \text{Phoneme} = q) = \frac{\#((m, p, q) \text{ in the data of all background speakers})}{\#((*, *, q) \text{ in the data of all background speakers})} \quad (3)$$

where  $m \in \mathcal{M}$ ,  $p \in \mathcal{P}$ ,  $(m, p, q)$  denotes the condition for which  $\text{Manner} = m$ ,  $\text{Place} = p$ , and  $\text{Phoneme} = q$ ,  $*$  represents all possible members in that class, and  $\#(\ )$  represents the total number of frames with phoneme labels and AF labels fulfill the description inside the parentheses. The probabilities of unseen AF combinations are set to zero. For each phoneme, a total of 60 probabilities can be obtained. These probabilities are the products of 6 manner classes and 10 place classes. Therefore, a system with  $N$  phonemes has  $60N$  probabilities in the UBMs.

### 2.3.2. Speaker models

Similar to the UBMs, each speaker model consists of the joint probabilities of the manner and place classes. For a particular speaker  $s$ , the joint probabilities corresponding to phoneme  $q$  are given by

$$P_s(\text{Manner} = m, \text{Place} = p | \text{Phoneme} = q) = \frac{\#((m, p, q) \text{ in the data of speaker } s)}{\#((*, *, q) \text{ in the data of speaker } s)}, \quad (4)$$

where only the data from speaker  $s$  are used in the computation. The accuracy of the speaker-dependent joint probabilities is limited by the amount of training data available. For some phoneme (e.g., /th/, /sh/ and /v/), the number of occurrences is too low for an accurate estimation of the joint probabilities. As a result, the pronunciation models of these phonemes are less discriminative.

### 2.3.3. Speaker models by MAP adaptation

To overcome the data sparseness problem, speaker models can be adapted from the UBMs. This approach can also establish a tighter coupling between the speaker models and background models, which can result in a better verification performance [1].

Given the background model corresponding to phoneme  $q$ , the joint probabilities for speaker  $s$  are given by:

$$\hat{P}_s(\text{Manner} = m, \text{Place} = p | \text{Phoneme} = q) = \beta_s^q P_s(\text{Manner} = m, \text{Place} = p | \text{Phoneme} = q) + (1 - \beta_s^q) P_{bg}(\text{Manner} = m, \text{Place} = p | \text{Phoneme} = q), \quad (5)$$

where  $\beta_s^q \in [0, 1]$  is a phoneme-dependent adaptation coefficient controlling the contribution of the speaker model (Eq. 4) and the background model (Eq. 3) on the adapted model. Similar to MAP adaptation of GMM-based systems [1],  $\beta_s^q$  is obtained by

$$\beta_s^q = \frac{\#((*, *, q) \text{ in the data of speaker } s)}{\#((*, *, q) \text{ in the data of speaker } s) + r}, \quad (6)$$

where  $r$  is a fixed relevance factor common to all phonemes and speakers. The purpose of  $r$  is to control the dependence of the adapted model's parameters on speaker's data. The estimation of  $r$  depends on the number of prior occurrences of  $(*, *, q)$  of all  $q$  in the training data. If the number of occurrences of  $(*, *, q)$  is much less than  $r$ , then  $\beta_s^q$  will be very close to 0 and the estimation of the

new model is less dependent on speaker’s data. On the contrary, if the number of occurrences of  $(*, *, q)$  is significantly greater than  $r$ , then  $\beta_s^q$  will be very close to 1 and the adapted model will become more dependent on speaker’s data.

### 2.3.4. Verification

The verification score  $S_{AFCPM}$  of a test utterance is defined as the difference between the speaker score  $S_s$  and background score  $S_b$ :

$$\begin{aligned} S_{AFCPM} &= S_s - S_b & (7) \\ &= \sum_{\substack{t=1, \\ p_s(X_t) \neq 0 \\ p_b(X_t) \neq 0 \\ q_t \neq \text{silence}}}^T (\log p_s(X_t) - \log p_b(X_t)), & (8) \end{aligned}$$

where for each  $t$ ,  $p_s(X_t)$  and  $p_b(X_t)$  are probabilities obtained from a speaker model of the claimed identity  $s$  and a background model, as follows:

$$\begin{aligned} p_s(X_t) &= \left\{ \begin{array}{l} P_s(\text{Manner} = l_t^M, \text{Place} = l_t^P | \text{Phoneme} = q_t) \\ \hat{P}_s(\text{Manner} = l_t^M, \text{Place} = l_t^P | \text{Phoneme} = q_t) \end{array} \right\} \mathfrak{G} \end{aligned}$$

and

$$p_b(X_t) = P_{bg}(\text{Manner} = l_t^M, \text{Place} = l_t^P | \text{Phoneme} = q_t). \quad (10)$$

In Eqs. 9 and 10,  $q_t$  is the phoneme at frame  $t$ . Because no speaker information is carried in the silence frames, they can be removed to improve the accuracy of the verification score. Moreover, only the “seen” AF combinations (i.e.,  $p_s(X_t) \neq 0$  and  $p_b(X_t) \neq 0$ ) appeared in both speaker and background models are considered during verification.

## 3. EXPERIMENTS AND RESULTS

Speaker verification was evaluated on the SPIDRE corpus [8], a subset of the Switchboard corpus. Genuine verification trials involved one handset-match conversation and two handset-mismatch conversations from each of the 44 target speakers (speaker sp1007 was discarded due to corrupted data); impostor attempts involved 200 conversations from 160 nontarget speakers. The same set of nontarget speakers’ conversations was applied to all target speaker models in the impostor attempts. Each of the testing utterances, which contains 5 minutes of speech (including silence), was split into short segments, with each segment ranging from 1 to 15 seconds according to the speaker turns labeled in the transcriptions [9]. All silence frames were removed by a voice activity detector.

The training conversation of all target speakers were used to train the phoneme models. The phoneme set consisted of 46 context-independent phonemes [9], including one silence and four noise, each of which was modeled by a three-state left-to-right HMM with 16 diagonal-covariance Gaussian mixtures per state. The HTK [10] was used to train the HMMs. Acoustic vectors of 39 dimensions—each comprising of 12 MFCCs, the normalized energy, and their first- and second-order derivatives—were used for training the phoneme models and for recognition.

The software Quicknet [11] was used to train two AF-MLPs, each of which was composed of 234 input nodes (nine frames of

26-dimensional MFCCs: 12 MFCCs, log energy, and the corresponding delta coefficients), 50 hidden nodes, and either 6 or 10 output nodes. To improve the robustness of AFs against handset variations, a total of 3,794 utterances randomly selected from all of the 10 handsets in the HTIMIT [12] corpus were used to train the AF-MLPs.

The aligned AF streams and phoneme sequences of all target speakers were used to train a set of UBMs ( $\Lambda_b^{AFCPM}$ ) representing the probabilities of 60 manner and place class combinations conditioned on 41 phonemes (excluding the silence and noise) in the phone set. The way to obtain the phoneme alignments of the training utterances was consistent with that of the verification utterances, which will be discussed in detail in Section 3.2.

Two approaches were adopted to obtain an AFCPM-based speaker model  $\Lambda_s^{AFCPM}$ . For the first approach, the probabilities in  $\Lambda_s^{AFCPM}$  were computed based on the AF streams and phoneme sequences of a given speaker  $s$  according to Eq. 4. This approach was referred to as AFCPM. In the second approach, the speaker probabilities were adapted from those of  $\Lambda_b^{AFCPM}$  using the training data from speaker  $s$  according to Eqs. 6 and 6 with  $r$  set to 18. Hereafter, this adaptation approach is referred to as A-AFCPM.

### 3.1. Spectral-based system and Score Fusion

The AFCPM and the conventional spectral features (MFCCs) characterize speakers at two different levels; the former represents the pronunciation behaviors of individual speakers, whereas the latter look at their vocal tract’s characteristics. Therefore, fusing the scores of AFCPM- and MFCC-based systems is expected to enhance speaker verification performance.

For the MFCC system, 24-dimensional MFCC vectors were used as features. Each feature vector  $\mathbf{x}_t$  comprises 12 MFCCs and the corresponding delta coefficients computed every 14ms using a Hamming window of 28ms. A 128-center universal background GMM  $\Lambda_b^{MFCC}$  was trained using all training conversations of all target speakers. For a speaker  $s$  in the target speaker set, a speaker GMM  $\Lambda_s^{MFCC}$  was adapted from  $\Lambda_b^{MFCC}$  using MAP adaptation [1].

Scores from the AFCPM and MFCC systems were fused according to the frame-weighted fusion proposed in [4]. The fusion weights were determined by  $K$ -fold cross validations. More specifically, the test data of the target and nontarget speakers were divided into  $K$  disjoint subsets, and the fusion weight was selected such that the average error obtained from the  $K$ -fold evaluations was minimized. It was suggested in [4], that the probabilities from the manner MLP are more reliable than those from the place MLP. Therefore, probabilities from the manner MLP ( $P(\text{Manner} = l_t^M | X_t)$ ) were adopted as  $w(t)$ .

### 3.2. Results and Discussions

Table 2 shows two sets of experimental results: recognized alignment (*Rec.*) and forced alignment (*F.A.*). In the former, the phoneme sequences were obtained from a null-grammar recognizer; in the latter the phoneme sequences were obtained by forced aligning the utterances with the transcribed word sequences and lexicon obtained from [9]. The forced alignments aim to minimize the effect of incorrect phoneme alignments on verification performance by assuming that a nearly perfect phoneme recognizer is available, thereby providing an upper bound performance of the AFCPM system. Table 2 also shows the performance of the

	Features	EER (%)		
		Matched	Mis-matched	All
	MFCC	8.55	18.18	15.84
Rec.	AFCPM	19.52	27.69	25.83
	A-AFCPM	18.07	26.69	24.04
	MFCC+AFCPM (error red. %)	8.50 (0.58)	16.61 (8.63)	14.44 (8.83)
	MFCC+A-AFCPM (error red. %)	8.25 (3.5)	16.04 (11.77)	13.76 (13.13)
F.A.	AFCPM	17.92	24.98	22.69
	A-AFCPM	16.60	23.98	21.72
	MFCC+AFCPM (error red. %)	8.08 (5.49)	15.24 (16.17)	13.26 (16.28)
	MFCC+A-AFCPM (error red. %)	7.74 (9.47)	14.68 (19.25)	12.95 (18.24)

**Table 2.** EERs and relative error reduction (in %) obtained from the MFCC system, the AFCPM system, and the fusion of the two systems. *Rec.* represents the recognized alignments and *F.A.* represents the forced alignments. *A-AFCPM* denotes the adaptive AFCPM system whose speaker models are adapted from the UBMs. *MFCC+AFCPM* (*MFCC + A-AFCPM*) denotes the fusion of frame-weighted MFCC scores and AFCPM (adaptive AFCPM) scores suggested in [4]. *Matched* (*Mismatched*) refers to the cases where the handset used by a claimant in a verification session is identical to (different from) the one used by the target speaker during the enrollment session. The test data from nontarget speakers under *Matched* and *Mismatched* are identical. *All* represents the overall EERs obtained from gathering all test data from the target speakers using both matched and mismatched handsets. Note that the MFCC system does not depend on any phoneme alignments.

MFCC system and the fused systems. The fusion weights were determined from a four-fold cross validation on all the testing data of the target and nontarget speakers. Note that the MFCC system does not require any phoneme alignments. The results of the MFCC system is the baseline for comparison.

When recognized alignments were used, an overall EER of 24.04% was obtained from AFCPM with adaptation (labeled as A-AFCPM). This represents a relative improvement of 7.0% when compared to the AFCPM system without adaptation (labeled as AFCPM). This suggests that better speaker models can be obtained by adapting the UBMs. Through the adaptation, speaker models can become tightly coupled to the UBMs. This helps prevent over-fitting the speaker models and improve their discriminative power. When forced alignments were used, the A-AFCPM system achieves an overall EER of 21.72%. The reduction from 24.04% (using recognized alignments) to 21.72% (using forced alignments) suggests that the accuracy of phoneme alignments is critical to the verification performance of the AFCPM system.

The system based on the frame-weighted fusion of the MFCC system and the AFCPM (A-AFCPM) system is labeled as *MFCC + AFCPM* (*MFCC + A-AFCPM*) in Table 2. The overall EERs on *MFCC + AFCPM* were reduced to 14.44% (an 8.83% error reduction) and 13.26% (an 16.28% error reduction) when recognized alignments and forced alignments were adopted. More significant error reductions were obtained from *MFCC + A-AFCPM*, the overall EER were reduced to 13.76% (an 13.13% error reduction) and 12.95% (an 18.24% error reduction) with the recognized

alignments and forced alignments. The results suggested that a better representation of speaker's pronunciation characteristics is achieved by obtaining the speaker models from adaptation. The results also show that A-AFCPM provides complementary speaker information, which is more useful than those from AFCPM, to the fusion system.

#### 4. CONCLUSIONS

This paper has presented an AFCPM speaker verification system in which the conditional pronunciation probabilities of the speaker models are adapted from those of the universal background models. The system distinguishes speakers by capturing their pronunciation characteristics via the conditional pronunciation modeling of two articulatory property streams. Experimental results have demonstrated the effectiveness of the AFCPM system in telephone-based speaker verification. A better verification performance of the AFCPM system was achieved when the speaker models were adapted from the background models because this increases speaker discrimination by establishing a tighter coupling between the speaker models and background models.

A lower error rate was achieved by the frame-weighted fusion of conventional MFCC and the adapted AFCPM scores. This suggests that the adapted speaker models incorporate more speaker information complementary to the spectral-feature comparing to those without adaptation.

#### 5. REFERENCES

- [1] D. A. Reynolds, T. F. Quatieri, and R. B. Dunn, "Speaker verification using adapted Gaussian mixture models," *Digital Signal Processing*, vol. 10, pp. 19–41, 2000.
- [2] D. Reynolds, et. al., "The superSID project: exploiting high-level information for high-accuracy speaker recognition," in *Proc. of ICASSP'03*, Hong Kong, April 2003, vol. 4, pp. 784–787.
- [3] D. Klusáček, J. Navrátil, D. A. Reynolds, and J. P. Campbell, "Conditional pronunciation modeling in speaker detection," in *Proc. of ICASSP'03*, 2003, vol. 4, pp. 804–807.
- [4] K.Y. Leung, M.W. Mak, and S.Y. Kung, "Articulatory feature-based conditional pronunciation modeling for speaker verification," in *Proc. of ICSLP'04*, 2004.
- [5] K. Kirchhoff, *Robust Speech Recognition Using Articulatory Information*, PhD thesis, University of Bielefeld, 1999.
- [6] <http://www.clsp.jhu.edu/ws2002/groups/supersid/>.
- [7] K.Y. Leung, M.W. Mak, and S.Y. Kung, "Applying articulatory features to telephone-based speaker verification," in *Proc. of ICASSP'04*, Montreal, May 2004, vol. 1, pp. 85–88.
- [8] J. P. Campbell and D. A. Reynolds, "Corpora for the evaluation of speaker recognition systems," in *Proc. of ICASSP'99*, 1999, vol. 2, pp. 829–832.
- [9] <http://www.isip.msstate.edu/projects/switchboard/>.
- [10] S. Young, D. Kershaw, J. Odell, D. Ollason, V. Valtchev, and P. Woodland, "The HTK book for HTK 3.0," Tech. Rep., Microsoft Corporation, 2000.
- [11] P. Farber, "Quicknet on multispart: fast parallel neural network training," Tech. Rep. TR-97-047, ICSI, 1997.
- [12] D. A. Reynolds, "HTIMIT and LLHDB: speech corpora for the study of handset transducer effects," in *Proc. ICASSP'97*, 1997, vol. 2, pp. 1535–1538.