

Vowel Creation by Articulatory Control in HMM-based Parametric Speech Synthesis

Zhen-Hua Ling¹, Korin Richmond², Junichi Yamagishi²

¹iFLYTEK Speech Lab, University of Science and Technology of China, P.R.China

²CSTR, University of Edinburgh, United Kingdom

zhling@ustc.edu, korin@cstr.ed.ac.uk, jyamagis@inf.ed.ac.uk

Abstract

This paper presents a method to produce a new vowel by articulatory control in hidden Markov model (HMM) based parametric speech synthesis. A multiple regression HMM (MRHMM) is adopted to model the distribution of acoustic features, with articulatory features used as external auxiliary variables. The dependency between acoustic and articulatory features is modelled by a group of linear transforms that are either estimated context-dependently or determined by the distribution of articulatory features. Vowel identity is removed from the set of context features used to ensure compatibility between the context-dependent model parameters and the articulatory features of a new vowel. At synthesis time, acoustic features are predicted according to the input articulatory features as well as context information. With an appropriate articulatory feature sequence, a new vowel can be generated even when it does not exist in the training set. Experimental results show this method is effective in creating the English vowel / Δ / by articulatory control without using any acoustic samples of this vowel.

Index Terms: Speech synthesis, articulatory features, multiple-regression hidden Markov model

1. Introduction

Hidden Markov model (HMM)-based parametric speech synthesis has become a mainstream speech synthesis method in recent years [1, 2]. This method is able to synthesise highly intelligible and smooth speech sounds [3, 4]. In addition, it makes speech synthesis far more flexible compared to the conventional unit selection and waveform concatenation approach. Several adaptation and interpolation methods have been applied to control model parameters and so diversify the characteristics of the generated speech [5, 6, 7]. However, this flexibility relies upon data-driven machine learning algorithms and it is difficult to integrate phonetic knowledge into the system directly when corresponding training data is not available. In previous work, we have proposed a method to improve the flexibility of HMM-based parametric speech synthesis further by integrating articulatory features [8, 9]. Here, we use “articulatory features” to

refer to the continuous movements of a group of speech articulators, such as the tongue, jaw, lips and velum, recorded by human articulography techniques. In this method, a unified acoustic-articulatory HMM is trained. The dependency between acoustic and articulatory features is modelled by a group of linear transforms which are either trained and tied context-dependently [8] or switched in the articulatory feature space [9]. During synthesis, the characteristics of the synthetic speech can be controlled flexibly by modifying the generated articulatory features according to phonetic rules. Experimental results have shown the effectiveness of this method in controlling the overall character of synthesised speech as well as the quality of a specific vowel [8, 9].

In this paper, we apply this method of articulatory control to the task of vowel creation in HMM-based parametric speech synthesis. In this task, the target vowel to be created does not occur in the training set, but its phonetic characteristics are known beforehand. We aim to produce this target vowel effectively at synthesis time once appropriate articulatory representations are provided. This is potentially useful for applications such as speech synthesis for limited resource languages, cross-language speaker adaptation, and so on. In our previous approach, articulatory features are treated as HMM observation vectors on which the acoustic features depend. In contrast, in this paper we treat the articulatory features as external explanatory variables for the mean vectors of Gaussians. Thus, we can integrate other forms of articulatory prediction model that are simpler to control than with the HMM itself. This model is called a “multiple regression HMM” (MRHMM). Our feature-space transform tying strategy [9] is also applied here, and we compare this with the context-dependent transform tying on the vowel creation task. Furthermore, we remove vowel identity from the set of context features used during context-dependent model training in order to ensure compatibility between the estimated model parameters and the articulatory features of a new target vowel at synthesis time.

This paper is organised as follows. Section 3 describes our proposed method in detail. Section 4 presents the results of our experiments and Section 5 summarises our conclusions.

2. Methods

2.1. MRHMM-based parametric speech synthesis

The MRHMM was initially proposed for automatic speech recognition (ASR). The aim is to improve the accuracy of acoustic modelling by introducing auxiliary features that are correlated with the acoustic features [10]. The MRHMM has also been applied to HMM-based parametric speech synthesis, with sentence-level “speech style” vectors as the explana-

This work is partially funded by the National Nature Science Foundation of China (Grant No. 60905010) and the National Natural Science Foundation of China - Royal Society of Edinburgh Joint Project (Grant No. 61111130120). The research leading to these results was partly funded from the European Community’s Seventh Framework Programme (FP7/2007-2013) under grant agreement 256230 (LISTA), and EPSRC grants EP/I027696/1 (Ultrax) and EP/J002526/1. More details about related work are introduced in a paper that has been submitted to IEEE Transactions on Audio, Speech, and Language Processing and which is currently under review.

tory variables [6]. The difference between this model and the standard HMM is that an auxiliary feature sequence is used to supplement the state sequence for determining the distribution of acoustic features. In this paper, that auxiliary feature sequence is comprised of articulatory trajectories. Let $\mathbf{X} = [\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_T^\top]^\top$ and $\mathbf{Y} = [\mathbf{y}_1^\top, \mathbf{y}_2^\top, \dots, \mathbf{y}_T^\top]^\top$ denote the parallel acoustic and articulatory feature sequence of the same length T . For each frame, the feature vector $\mathbf{x}_t \in \mathcal{R}^{3D_X}$ and $\mathbf{y}_t \in \mathcal{R}^{3D_Y}$ consist of static parameters and their velocity and acceleration components, where D_X and D_Y are the dimensionality of the static acoustic and articulatory features respectively. A detailed definition of these dynamic features can be found in [8]. The distribution of \mathbf{X} in the conventional MRHMM [10] can be written as

$$P(\mathbf{X}|\lambda, \mathbf{Y}) = \sum_{\mathbf{q}} \pi_{q_0} \prod_{t=1}^T a_{q_{t-1}q_t} b_{q_t}(\mathbf{x}_t|\mathbf{y}_t), \quad (1)$$

$$b_j(\mathbf{x}_t|\mathbf{y}_t) = \mathcal{N}(\mathbf{x}_t; \mathbf{A}_j \boldsymbol{\xi}_t + \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j), \quad (2)$$

where π_j and a_{ij} represent initial state probability and state transition probability; $b_j(\cdot)$ is the state observation probability density function (PDF) for state j ; $\mathbf{q} = \{q_1, q_2, \dots, q_T\}$ is the state sequence for \mathbf{X} ; $\boldsymbol{\xi}_t = [\mathbf{y}_t^\top, 1]^\top \in \mathcal{R}^{3D_Y+1}$ is the expanded articulatory feature vector; $\mathbf{A}_j \in \mathcal{R}^{3D_X \times (3D_Y+1)}$ is the regression matrix for state j .

To train the MRHMM-based parametric speech synthesis system, the procedures for standard HMM-based synthesiser training [1] (without articulatory features) are first followed. Context-dependent HMMs are trained using rich context information that includes detailed phonetic and prosodic features [1]. To deal with the data-sparsity problem, a decision-tree-based model clustering technique that uses the minimum description length (MDL) criterion [11] is applied to cluster context-dependent HMMs. Then, the estimated mean vector and covariance matrix for each state are used as the initial values of $\boldsymbol{\mu}_j$ and $\boldsymbol{\Sigma}_j$ in the MRHMM. The regression matrix \mathbf{A}_j is initialised as a zero matrix. After introducing articulatory features, these parameters are iteratively updated to maximise $P(\mathbf{X}|\lambda, \mathbf{Y})$ using the EM algorithm¹. The detailed formulae for this are to be found in [10]. Next, a state alignment to the acoustic features is performed using the trained MRHMM in order to train context-dependent state duration probabilities [1].

At synthesis time, acoustic features are generated following the maximum output probability criterion [2]. For the purpose of simplification, only the optimal HMM state sequence is considered. First, this optimal state sequence $\mathbf{q}^* = \{q_1^*, q_1^*, \dots, q_T^*\}$ is determined using the trained duration distributions [1]. Given auxiliary feature sequence \mathbf{Y} , the optimal acoustic feature sequence \mathbf{X}^* is generated by maximising

$$P(\mathbf{X}|\lambda, \mathbf{Y}, \mathbf{q}^*) = \prod_{t=1}^T \mathcal{N}(\mathbf{x}_t; \mathbf{A}_{q_t^*} \boldsymbol{\xi}_t + \boldsymbol{\mu}_{q_t^*}, \boldsymbol{\Sigma}_{q_t^*}). \quad (3)$$

This is the conventional parameter generation problem [2]. The only difference is that the mean vector at each frame is calculated as $\mathbf{A}_{q_t^*} \boldsymbol{\xi}_t + \boldsymbol{\mu}_{q_t^*}$ instead of $\boldsymbol{\mu}_{q_t^*}$.

In the vowel creation task, the articulatory-phonetic characteristics of the target vowel, such as tongue position and similarity with other vowels, are assumed to be known. They are also

¹The acoustic features commonly consist of spectral and F0 parameters extracted from the waveforms of the training sentence. In MRHMM training, \mathbf{X} only contains the spectral feature stream. The relationship between the articulatory features and the F0 features is not considered in this paper.

included in the context features and the question set for training context-dependent models with decision-tree-based clustering. \mathbf{Y} in (3) contains the articulatory representations of the target vowel and we expect to generate the acoustic features of this vowel by solving (3) at synthesis time.

2.2. Feature-space-switched MRHMM

As shown in (2), the regression matrix \mathbf{A}_j in the conventional MRHMM is context-dependent. In previous work [9], we have proposed a feature-space transform tying method to take into account the effect of articulatory features in determining the transform matrices for the unified acoustic-articulatory modelling. This method improved the controllability on synthetic speech when manipulating articulatory features at synthesis time [9]. The same issues apply to the MRHMM-based speech synthesis method introduced in Section 2.1 when attempting the vowel creation task. When solving (3), the articulatory features \mathbf{y}_t of the new vowel may conflict with the transform matrix \mathbf{A}_j which could be estimated using the training samples of other vowels. Therefore, we also apply the feature-space transform tying method in this paper, to create a model that we term a *feature-space-switched MRHMM*. A GMM model $\lambda^{(G)}$ containing M mixture components is trained in advance using only the articulatory stream of the training data to obtain M clusters in the articulatory space. Then, the regression matrices are trained for each mixture component of $\lambda^{(G)}$ instead of for each state of the MRHMM. Mathematically, we rewrite (2) as

$$b_j(\mathbf{x}_t|\mathbf{y}_t) = \sum_{k=1}^M P(\mathbf{x}_t, m_t = k|\mathbf{y}_t, q_t = j, \lambda, \lambda^{(G)}), \quad (4)$$

$$= \sum_{k=1}^M \zeta_k(t) P(\mathbf{x}_t|\mathbf{y}_t, q_t = j, m_t = k, \lambda, \lambda^{(G)}), \quad (5)$$

where m_t denotes the mixture component index of $\lambda^{(G)}$ for the articulatory feature vector at frame t . The HMM state sequence \mathbf{q} and the GMM mixture sequence $\mathbf{m} = \{m_1, m_2, \dots, m_N\}$ are reasonably assumed to be independent of each other, so that

$$\begin{aligned} P(m_t = k|\mathbf{y}_t, q_t = j, \lambda, \lambda^{(G)}) &= P(m_t = k|\mathbf{y}_t, \lambda^{(G)}) \\ &= \zeta_k(t). \end{aligned} \quad (6)$$

For each Gaussian mixture, the dependency between the acoustic features and the articulatory features is described by

$$P(\mathbf{x}_t|\mathbf{y}_t, q_t = j, m_t = k, \lambda, \lambda^{(G)}) = \mathcal{N}(\mathbf{x}_t; \mathbf{A}_k \boldsymbol{\xi}_t + \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j), \quad (7)$$

where $\mathbf{A}_k \in \mathcal{R}^{3D_X \times (3D_Y+1)}$ is the regression matrix for the k -th mixture of $\lambda^{(G)}$. The parameter set $\{\mathbf{A}_k, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j\}$ is estimated using the EM algorithm by maximising $P(\mathbf{X}|\lambda, \mathbf{Y})$. The detailed formulae are similar to those introduced in [9] and so are omitted here.

At synthesis time, the parameter generation criterion in (3) is modified to

$$P(\mathbf{X}|\lambda, \mathbf{Y}, \mathbf{q}^*) = \prod_{t=1}^T \sum_{k=1}^M \zeta_k(t) \mathcal{N}(\mathbf{x}_t; \mathbf{A}_k \boldsymbol{\xi}_t + \boldsymbol{\mu}_{q_t^*}, \boldsymbol{\Sigma}_{q_t^*}), \quad (8)$$

where $\zeta_k(t)$ is calculated based on the given articulatory features \mathbf{Y} . This is a parameter generation problem with mixtures of Gaussians at each frame. It can be solved by either considering only the optimal mixture sequence or by using an EM-based iterative estimation method [2].

Table 1: Summary of the various systems used in the experiments.

Label	Model Structure		
	HMM	Context Features	Regression Matrix
<i>STD-F</i>	<i>standard</i>	<i>full</i>	<i>N/A</i>
<i>STD-T</i>	<i>standard</i>	<i>tailored</i>	<i>N/A</i>
<i>MR-FC</i>	<i>MRHMM</i>	<i>full</i>	<i>context-dependent</i>
<i>MR-TC</i>	<i>MRHMM</i>	<i>tailored</i>	<i>context-dependent</i>
<i>MR-TF</i>	<i>MRHMM</i>	<i>tailored</i>	<i>feature-space-switched</i>

2.3. Context feature tailoring

In the MRHMM-based speech synthesis method, the model parameters A_j , μ_j , and Σ_j in (2) are trained context-dependently. To reconstruct speech signals more accurately, the context features describing each phone commonly consist of detailed segmental and suprasegmental information, such as current and surrounding phone identifiers, prosodic boundaries, stress and accent positions, part of speech, and so on. In the vowel creation task, no training samples of the target vowel are available however. At synthesis time, the model parameters A_j , μ_j , Σ_j determined by the context features of the target vowel are actually estimated using the samples of other phones. Therefore, these model parameters may be incompatible with the input articulatory features of the target vowel which are phone-dependent and are unseen at training time. In the feature-space-switched MRHMM method described in Section 2.2, the regression matrices are not context-dependent but assigned according to the posterior probability of each frame in the articulatory feature space. However, the μ_j and Σ_j in (7) are still context-dependent. In order to ensure the compatibility between the context-dependent model parameters and the input articulatory features in the vowel creation task, the vowel identity feature is removed from the set of context features during MRHMM training. This ‘‘context feature tailoring’’ method is expected to improve the generalisation property of the trained model parameters for unseen vowels, and we compare the use of this method with the standard use of full context features in our experiments.

3. Experiments

3.1. Experimental conditions

We used the same multi-channel articulatory database as we used in our previous work [8, 9] for the experiments in this paper. It contains acoustic waveforms recorded concurrently with EMA data using a Carstens AG500 electromagnetic articulograph [12]. Around 1300 phonetically balanced sentences were read by a male British English speaker. The waveforms were in 16kHz PCM format with 16 bit precision. Six EMA sensors were placed at the *tongue dorsum* (T3), *tongue body* (T2), *tongue tip* (T1), *lower lip* (LL), *upper lip* (UL), and *lower incisor* (LI) of the speaker. Each sensor recorded spatial location in 3 dimensions at a 200Hz sample rate: coordinates on the x- (front to back), y- (bottom to top) and z-(left to right) axes (relative to viewing the speaker’s face from the front). Only the x- and y-coordinates of the six sensors were used in our experiments because the movements in the z-axis were relatively small. There were a total of 12 static articulatory features at each frame. The static acoustic features were composed of F0 and 40-order frequency-warped line spectral pairs (LSPs) [4] plus an extra gain dimension, which were derived using

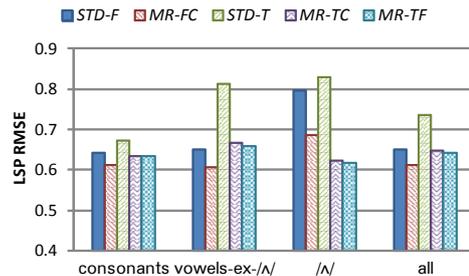


Figure 1: LSP RMSEs for the five systems listed in Table 1. ‘‘vowels-ex-/Λ/’’ indicates all vowels excluding the target vowel /Λ/.

STRAIGHT[13] analysis with a frame shift of 5ms.

In our experiments, the scenario of vowel creation was simulated by selecting a target vowel from the British English phone set and removing all sentences containing this target vowel from the training set. Vowel /Λ/ was selected as the target vowel in the experiment described here. 809 sentences in the database which contain no instances of this vowel were selected for a training set. Five acoustic models were trained in total. Descriptions of these models are shown in Table 1. A five-state, left-to-right HMM structure with no skips was adopted and diagonal covariance matrices were used for all five systems. The *STD-F* and *STD-T* systems were trained following the conventional HMM-based parametric speech synthesis approach [1]. The difference between these two systems was that vowel identity was removed from the context features in the *STD-T* system. The number of regression matrices in the *MR-FC*, *MR-TC*, and *MR-TF* systems was set to 65, 65, and 64 respectively.

3.2. Objective evaluation

We randomly selected 50 sentences from the remaining 454 sentences which contain instances of vowel /Λ/ to form a test set. These sentences were synthesised using the five systems listed in Table 1. The LSPs were generated using state durations derived from state alignment against the natural speech performed using each system. Natural articulatory recordings were used as the input for the *MR-FC*, *MR-TC*, and *MR-TF* systems at synthesis time. LSP root mean squared error (RMSE) for different types of phone was calculated and these are shown in Fig. 1.²

From this figure, we see that the *STD-F* system had much higher LSP RMSE for /Λ/ than for the other vowels and consonants, because no samples for /Λ/ were available during training. Once the MRHMM-based speech synthesis method was applied and the natural articulatory features were available at synthesis time, the *MR-FC* system could achieve much lower LSP RMSEs than the *STD-F* system, especially for the new vowel /Λ/. Comparing the context feature tailoring method with using standard full context features, the *STD-T* system was worse than the *STD-F* system, while the *MR-TC* system had better prediction accuracy than the *MR-FC* system for the vowel /Λ/. In fact, the LSP RMSE on the /Λ/ vowel for system *MR-TC* is not higher than that of the *STD-F* system on the other vowels. The feature-space-switched MRHMM modelling can further improve the prediction accuracy slightly. However, the difference between the performance of the *MR-TC* and *MR-*

²Some examples of the synthetic speech can be found at <http://staff.ustc.edu.cn/~zhling/VowCreIS2012/demo.html>.

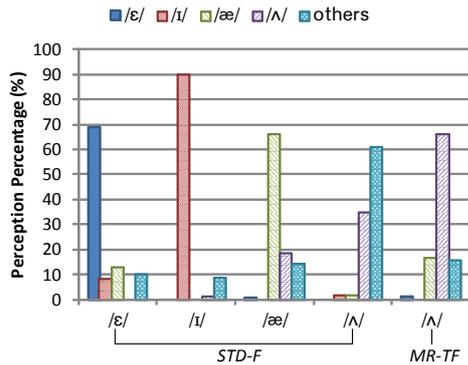


Figure 2: Vowel identity perception results for synthesising different vowels using the *STD-F* system and creating vowel /ʌ/ by articulatory control using the *MR-TF* system.

TF systems is not significant in this task because the context feature tailoring method can also deal with the possible conflict between the context-dependent regression matrix and the articulatory features input during synthesis.

3.3. Subjective evaluation

We also carried out a vowel identity perception test to further evaluate the effectiveness of creating the target /ʌ/ vowel. Five monosyllabic words (“but”, “hum”, “puck”, “tun”, “dud”) containing the /ʌ/ vowel were selected and embedded within a carrier sentence “Now we’ll say ... again”. These sentences were synthesised using the *STD-F* system and the *MR-TF* system respectively. Because natural articulatory recordings of these sentences were not available, the articulatory features generated from the HMM-based articulatory prediction model [14]³ were adopted for the *MR-TF* system. For the purpose of comparison, we substituted the vowel /ʌ/ in the five monosyllabic words with /ɛ/, /ɪ/ and /æ/, and then synthesised the respective test sentences using the *STD-F* system. Thus, we created twenty-five stimuli for the vowel identity perception test. Thirty-two native English listeners were asked to listen to these stimuli and to write down the key word in the carrier sentence they heard. Then, we calculated the percentages for how the vowels were perceived. These results are shown in Fig. 2. We see that only 35% of the synthesised vowels /ʌ/ were perceived correctly using the *STD-F* system, due to the lack of acoustic training samples for this vowel. This percentage is above chance level because the state models used here to synthesise vowels /ʌ/ may be estimated using the acoustic features of the vowels which have similar pronunciation to /ʌ/ due to the decision-tree-based model clustering. Using the *MR-TF* system and the generated articulatory features, this percentage increased to 66.25%, which is close to the perception accuracy of synthesising vowel /ɛ/ (68.75%) and /æ/ (66.25%) using the *STD-F* system.

³In our experiments, this model was trained using the full database with full context features. In practical scenarios, however, we would not expect the articulatory-acoustic features of target vowel to be available in the training set. Therefore, an articulatory prediction method that is capable of integrating phonetic knowledge to generate articulatory representations for a new vowel is necessary. This will be one focus of our future work.

4. Conclusions

An MRHMM-based speech synthesis method has been presented in this paper for creating a new vowel without acoustic training instances, but using articulatory representations at synthesis time instead. In this method, articulatory features are combined with context information to determine the distribution of acoustic features at each frame. A method for feature-space regression matrix switching and a strategy of context feature tailoring have been introduced to ensure compatibility between context-dependent model parameters and unseen articulatory features of new vowel. Our experiment on producing the English vowel /ʌ/ has shown the effectiveness of our method.

5. References

- [1] K. Tokuda, H. Zen, and A. W. Black, “HMM-based approach to multilingual speech synthesis,” in *Text to speech synthesis: New paradigms and advances*, S. Narayanan and A. Alwan, Eds. Prentice Hall, 2004.
- [2] K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, and T. Kitamura, “Speech parameter generation algorithms for HMM-based speech synthesis,” in *ICASSP*, vol. 3, 2000, pp. 1315–1318.
- [3] H. Zen, T. Toda, M. Nakamura, and K. Tokuda, “Details of Nitech HMM-based speech synthesis system for the Blizzard Challenge 2005,” *IEICE Trans. Inf. & Syst.*, vol. E90-D, no. 1, pp. 325–333, 2007.
- [4] Z.-H. Ling, Y.-J. Wu, Y.-P. Wang, L. Qin, and R.-H. Wang, “USTC system for Blizzard Challenge 2006: an improved HMM-based speech synthesis method,” in *Blizzard Challenge Workshop*, 2006.
- [5] J. Yamagishi and T. Kobayashi, “Average-voice-based speech synthesis using HSMM-based speaker adaptation and adaptive training,” *IEICE Trans. on Inf. & Syst.*, vol. E90-D, no. 2, pp. 533–543, 2007.
- [6] T. Nose, J. Yamagishi, T. Masuko, and T. Kobayashi, “A style control technique for HMM-based expressive speech synthesis,” *IEICE Trans. Inf. & Syst.*, vol. E90-D, no. 9, pp. 1406–1413, 2007.
- [7] M. Tachibana, J. Yamagishi, T. Masuko, and T. Kobayashi, “Speech synthesis with various emotional expressions and speaking styles by style interpolation and morphing,” *IEICE Trans. Inf. & Syst.*, vol. E88-D, no. 11, pp. 2484–2491, 2005.
- [8] Z.-H. Ling, K. Richmond, J. Yamagishi, and R.-H. Wang, “Integrating articulatory features into HMM-based parametric speech synthesis,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 6, pp. 1171–1185, Aug. 2009.
- [9] Z.-H. Ling, K. Richmond, and J. Yamagishi, “Feature-space transform tying in unified acoustic-articulatory modelling for articulatory control of HMM-based speech synthesis,” in *Interspeech*, 2011, pp. 117–120.
- [10] K. Fujinaga, M. Nakai, H. Shimodaira, and S. Sagayama, “Multiple-regression hidden Markov model,” in *ICASSP*, 2001, pp. 513–516.
- [11] K. Shinoda and T. Watanabe, “MDL-based context-dependent subword modeling for speech recognition,” *J. Acoust. Soc. Japan (E)*, vol. 21, no. 2, pp. 79–86, 2000.
- [12] K. Richmond, P. Hoole, and S. King, “Announcing the electromagnetic articulography (day 1) subset of the mngu0 articulatory corpus,” in *Interspeech*, 2011, pp. 1505–1508.
- [13] H. Kawahara, I. Masuda-Katsuse, and A. de Cheveigne, “Restructuring speech representations using 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.
- [14] Z.-H. Ling, K. Richmond, and J. Yamagishi, “An analysis of HMM-based prediction of articulatory movements,” *Speech Communication*, vol. 52, no. 10, pp. 834–846, 2010.