



# Modulation Spectrum-Constrained Trajectory Training Algorithm for HMM-Based Speech Synthesis

Shinnosuke Takamichi<sup>1,2</sup>, Tomoki Toda<sup>1</sup>, Alan W. Black<sup>2</sup>, Satoshi Nakamura<sup>1</sup>

<sup>1</sup> Graduate School of Information Science, Nara Institute of Science and Technology (NAIST), Japan

<sup>2</sup> Language Technologies Institute, Carnegie Mellon University (CMU)

shinnosuke-t@is.naist.jp

## Abstract

This paper presents a novel training algorithm for Hidden Markov Model (HMM)-based speech synthesis. One of the biggest issues causing significant quality degradation in synthetic speech is the over-smoothing effect often observed in generated speech parameter trajectories. Recently, we have found that a Modulation Spectrum (MS) of the generated speech parameters is sensitively correlated with the over-smoothing effect, and have proposed the parameter generation algorithm considering the MS. The over-smoothing effect is effectively alleviated by the proposed parameter generation algorithm. On the other hand, it loses the computationally-efficient generation processing of the conventional generation algorithm. In this paper, the MS is integrated into the training stage instead of the parameter generation stage in a similar manner as our previous work on Gaussian Mixture Model (GMM)-based spectral parameter trajectory conversion. The trajectory HMM is trained with a novel objective function consisting of both the conventional trajectory HMM likelihood and a newly implemented MS likelihood. This training framework is further extended to the  $F_0$  component. The experimental results demonstrate that the proposed algorithm yields improvements in synthetic speech quality while preserving a capability of the computationally-efficient generation processing.

**Index Terms:** HMM-based speech synthesis, over-smoothing, global variance, modulation spectrum, trajectory training

## 1. Introduction

Statistical parametric speech synthesis based on Hidden Markov Models (HMMs) [1] is an effective framework for generating diverse types of synthetic speech. Speech parameters, i.e., spectral and excitation features and HMM-state duration are simultaneously modeled with context-dependent HMMs in a unified framework [2]. In synthesis, the speech parameter trajectories are generated by maximizing the likelihood of the HMMs [3]. This approach allows us not only to apply several techniques for flexibly controlling synthetic speech [4, 5, 6] to various speech-based systems [7, 8], but also to build the speech synthesizer without complicated tuning compared to sample-based [9] or deep neural nets-based [10] speech synthesis. The further merit of HMM-based speech synthesis is the computationally-efficient speech parameter generation [3]. This generation algorithm is very helpful to deploy the speech-based systems that need the fast speech synthesis, e.g., speech-to-speech translation system [11].

One of the critical problems in HMM-based speech synthesis is that the parameter trajectories generated from the HMMs are often over-smoothed. This phenomenon causes significant degradation of the perceptual quality and makes

synthetic speech sound muffled [12]. To address this over-smoothing problem, we have found Modulation Spectrum (MS) [13, 14, 15] as a feature well quantify the over-smoothing effect. The MS is defined as the power spectrum of the speech parameter trajectories, and is regarded as an extension of the Global Variance (GV) [16]. The MS of the generated trajectories is often lower than that of natural speech parameter trajectories. [17] integrated a metric on the MS into the parameter generation algorithm to keep the MS close to natural one, and they reported the improvements in synthetic speech quality. However, as [18, 19, 20] reported in the parameter generation considering the GV, the parameter generation algorithm considering the MS also loses the conventional computationally-efficient generation ability because the objective function in synthesis does not solved in a closed form.

As a method to recover such features while adopting the computationally-efficient generation algorithm, [21] have proposed a metric to integrate the GV into the training stage instead of the synthesis stage. They reformulated trajectory HMMs [22] imposing the constraint between the static and dynamic features for spectral and  $F_0$  components. By training the trajectory HMMs with the GV constraint, the computationally-efficient generation algorithm is straightforwardly adopted, but quality benefits by the GV metric is observed in synthetic speech. We can expect that same reformulation with the MS will give us the further gain in synthetic speech quality.

This paper proposes the MS-constrained trajectory training algorithm to HMM-based speech synthesis in the same manner as our previous work [23] on Gaussian Mixture Model (GMM)-based voice conversion [24]. The trajectory HMM is trained with a novel objective function consisting of both the conventional trajectory HMM likelihood and a newly implemented MS likelihood. This training framework is further extended to the  $F_0$  component. The proposed training algorithm is compared to the basic training [2], the conventional trajectory training [22], and GV-constrained trajectory training [21] in term of synthetic speech quality. The result demonstrates the proposed training algorithm achieves the best synthetic speech quality compared to these training algorithms.

## 2. Basic Framework [2]

### 2.1. Training Algorithm [1]

In HMM-based speech synthesis, A HMM parameter set  $\lambda$  is estimated using the contextual factor sequence  $\mathbf{X}$  of input text and the speech feature sequence  $\mathbf{Y} = [\mathbf{Y}_1^\top, \dots, \mathbf{Y}_t^\top, \dots, \mathbf{Y}_T^\top]^\top$  of  $T$  frames as follows:

$$\lambda = \operatorname{argmax} L_{\text{basic}} = \operatorname{argmax} P(\mathbf{Y} | \mathbf{X}, \lambda). \quad (1)$$

The output probability density function of HMM-state index  $q$  is given as:

$$P(\mathbf{Y}_t | \mathbf{X}, q, \boldsymbol{\lambda}) = \mathcal{N}(\mathbf{Y}_t; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q), \quad (2)$$

where  $\mathbf{Y}_t$  is given by  $3D$ -dimensional joint static and dynamic feature vectors,  $[\mathbf{y}_t^\top, \Delta\mathbf{y}_t^\top, \Delta\Delta\mathbf{y}_t^\top]^\top$ , where  $\mathbf{y}_t = [y_t(1), \dots, y_t(d), \dots, y_t(D)]^\top$  is represented as a  $D$ -dimensional vector at frame  $t$ , and  $d$  is a dimensional index.  $\mathcal{N}(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$  denotes Gaussian distribution of a mean vector  $\boldsymbol{\mu}$  and a covariance matrix  $\boldsymbol{\Sigma}$ . The HMM parameter set  $\boldsymbol{\lambda}$  consists of  $Q$  HMM-states where each HMM-state has the individual mean vector  $\boldsymbol{\mu}_q$  and covariance matrix  $\boldsymbol{\Sigma}_q$ .

## 2.2. Parameter Generation Algorithm [3]

Given the contextual factor sequence  $\mathbf{X}$ , the generated parameter sequence  $\hat{\mathbf{y}}_q = [\hat{\mathbf{y}}_1^\top, \dots, \hat{\mathbf{y}}_t^\top, \dots, \hat{\mathbf{y}}_T^\top]^\top$  is analytically determined by maximizing the output probability of the speech feature vector sequence target  $\mathbf{Y}$  given  $\mathbf{X}$  under a constraint  $\mathbf{Y} = \mathbf{W}\mathbf{y}$  as follows:

$$\hat{\mathbf{y}}_q = \underset{\mathbf{y}}{\operatorname{argmax}} P(\mathbf{W}\mathbf{y} | \mathbf{X}, \hat{\mathbf{q}}, \boldsymbol{\lambda}) \quad (3)$$

$$= \mathbf{R}_q^{-1} \mathbf{r}_q = \left( \mathbf{W}^\top \mathbf{D}_q^{-1} \mathbf{W} \right)^{-1} \mathbf{W}^\top \mathbf{D}_q^{-1} \mathbf{E}_q, \quad (4)$$

where  $\mathbf{W}$  is a  $3DT$ -by- $DT$  weight matrix to calculate the dynamic features [3].  $\hat{\mathbf{q}} = [\hat{q}_1, \dots, \hat{q}_t, \dots, \hat{q}_T]$  is the sub-optimum state sequence determined by maximizing state duration probability distribution function  $P(q | \mathbf{X}, \boldsymbol{\lambda})$ , where  $q_t$  is a sub-optimum HMM-state index at frame  $t$ . The mean vector  $\mathbf{E}_q = [\boldsymbol{\mu}_{q_1}^\top, \dots, \boldsymbol{\mu}_{q_t}^\top, \dots, \boldsymbol{\mu}_{q_T}^\top]^\top$  and the covariance matrix  $\mathbf{D}_q = \operatorname{diag}_{3D} [\boldsymbol{\Sigma}_{q_1}, \dots, \boldsymbol{\Sigma}_{q_t}, \dots, \boldsymbol{\Sigma}_{q_T}]$  are calculated using the corresponding HMM-state, where the notation  $\operatorname{diag}_{3D}$  denotes the construction of a block diagonal matrix that has the  $3D$ -by- $3D$  diagonal elements.

## 3. Trajectory Training [22]

Trajectory HMM is reformulated by imposing the constraint between the static and dynamic features. The objective function for the trajectory training is written as:

$$L_{\text{trj}} = P(\mathbf{y} | \mathbf{X}, \hat{\mathbf{q}}, \boldsymbol{\lambda}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}_q, \mathbf{R}_q^{-1}). \quad (5)$$

The mean vector  $\hat{\mathbf{y}}_q$  is given by Eq. (4) and the inter-frame correlation is effectively modeled by the temporal covariance matrix  $\mathbf{R}_q^{-1}$ . In training, the HMM parameters are updated by maximizing  $L_{\text{trj}}$ .

### 3.1. Estimation of Model Parameters

Here, let  $\boldsymbol{\mu} = [\boldsymbol{\mu}_1^\top, \dots, \boldsymbol{\mu}_q^\top, \dots, \boldsymbol{\mu}_Q^\top]^\top$  and  $[\boldsymbol{\Sigma}_1^{-1}, \dots, \boldsymbol{\Sigma}_q^{-1}, \dots, \boldsymbol{\Sigma}_Q^{-1}]^\top$  be the joint parameters of  $\boldsymbol{\mu}_q$  and  $\boldsymbol{\Sigma}_q^{-1}$  over all HMM-states, respectively. The mean vector  $\mathbf{E}_q$  and the precision matrix  $\mathbf{D}_q^{-1}$  are represented as:

$$\mathbf{E}_q = \mathbf{S}_q \boldsymbol{\mu}, \quad (6)$$

$$\mathbf{D}_q^{-1} = \operatorname{diag}_{3D} [\mathbf{S}_q \boldsymbol{\Sigma}_q^{-1}], \quad (7)$$

where  $\mathbf{S}_q = [\mathbf{S}_{q_1}, \dots, \mathbf{S}_{q_T}]^\top \otimes \mathbf{I}_{3D}$  is a  $3DT$ -by- $3DQ$  matrix,  $\mathbf{S}_{q_t}$  is an  $M$ -dimensional vector of which the  $q$ -th component is 1 when  $q = \hat{q}_t$  and otherwise are 0 as shown in Fig. 1, and  $\mathbf{I}_{3D}$  indicates the  $3D$ -by- $3D$  identity matrix.

To optimize these model parameters for the objective func-

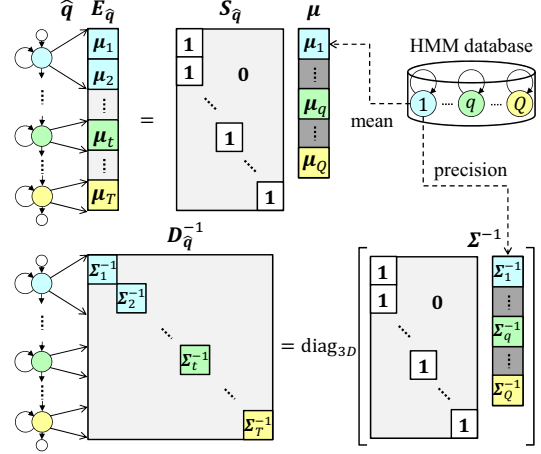


Figure 1: Graphic representation between variables used in this paper.

tion, we employ the steepest descent algorithm. as follows:

$$\boldsymbol{\mu}^{(i+1)} = \boldsymbol{\mu}^{(i)} + \alpha \frac{\partial \log L_{\text{trj}}}{\partial \boldsymbol{\mu}} \Big|_{\boldsymbol{\mu}=\boldsymbol{\mu}^{(i)}}, \quad (8)$$

where  $\alpha$  is a learning rate, and  $i$  is an iteration index.  $\boldsymbol{\Sigma}^{-1}$  are also optimized in the same manner. The gradients are given by:

$$\frac{\partial \log L_{\text{trj}}}{\partial \boldsymbol{\mu}} = \mathbf{S}_q^\top \mathbf{D}_q^{-1} \mathbf{W} (\mathbf{y} - \hat{\mathbf{y}}_q), \quad (9)$$

$$\frac{\partial \log L_{\text{trj}}}{\partial \boldsymbol{\Sigma}^{-1}} = \frac{1}{2} \mathbf{S}_q^\top \operatorname{diag}_{3D}^{-1} \left[ \mathbf{W} \left( \mathbf{R}_q^{-1} + \hat{\mathbf{y}}_q \hat{\mathbf{y}}_q^\top - \mathbf{y} \mathbf{y}^\top \right) - \mathbf{E}_q (\hat{\mathbf{y}}_q - \mathbf{y})^\top \mathbf{W}^\top - \mathbf{W} (\hat{\mathbf{y}}_q - \mathbf{y}) \mathbf{E}_q^\top \right], \quad (10)$$

## 4. Global Variance (GV)-Constrained Trajectory Training [21]

### 4.1. Global Variance (GV) [16]

The GV  $\mathbf{v}(\mathbf{y}) = [v(1), \dots, v(D)]^\top$  is defined as the second order moment of the trajectory  $\mathbf{y}$ , and its  $d$ -th component is given as:

$$v(d) = \frac{1}{T} \sum_{t=1}^T \left( y_t(d) - \frac{1}{T} \sum_{\tau=1}^T y_\tau(d) \right)^2. \quad (11)$$

### 4.2. Objective Function for GV-Constrained Training

$\{\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}\}$  is updated by maximizing the following objective function  $L_{\text{gvtrj}}$  consisting of the trajectory HMM and GV likelihoods:

$$L_{\text{gvtrj}} = P(\mathbf{y} | \mathbf{X}, \hat{\mathbf{q}}, \boldsymbol{\lambda}) P(\mathbf{v}(\mathbf{y}) | \mathbf{X}, \hat{\mathbf{q}}, \boldsymbol{\lambda}, \boldsymbol{\lambda}_v)^{\omega_v T}, \quad (12)$$

$$P(\mathbf{v}(\mathbf{y}) | \mathbf{X}, \hat{\mathbf{q}}, \boldsymbol{\lambda}, \boldsymbol{\lambda}_v) = \mathcal{N}(\mathbf{v}(\mathbf{y}); \mathbf{v}(\hat{\mathbf{y}}_q), \boldsymbol{\Sigma}_v), \quad (13)$$

where  $\omega_v$  is a weight of the GV likelihood,  $\boldsymbol{\Sigma}_v$  is a covariance matrix of the GV, and  $\boldsymbol{\lambda}_v$  is a model parameter set of the GV. This algorithm updates the model parameters to make the GV of the generated parameter sequence close to natural one.

## 5. Modulation Spectrum (MS)-Constrained Trajectory Training

### 5.1. Modulation Spectrum [13]

The MS is defined as the power spectrum of the parameter sequence; i.e., temporal fluctuation of the parameter sequence is

decomposed into individual modulation frequency components and their power values are represented as the MS. In this paper, the MS  $\mathbf{s}(\mathbf{y})$  of the parameter sequence  $\mathbf{y}$  is defined as:

$$\mathbf{s}(\mathbf{y}) = [\mathbf{s}(1)^\top, \dots, \mathbf{s}(d)^\top, \dots, \mathbf{s}(D)^\top]^\top, \quad (14)$$

$$\mathbf{s}(d) = [s_d(0), \dots, s_d(f), \dots, s_d(D'_s - 1)]^\top \quad (15)$$

$$s_d(f) = \hat{R}_{d,f}^2 + \hat{I}_{d,f}^2 \quad (16)$$

$$= \left( \sum_{t=1}^T y_t(d) \cos kt \right)^2 + \left( \sum_{t=1}^T y_t(d) \sin kt \right)^2 \quad (17)$$

where  $2D_s$  is a length of Discrete Fourier Transform (DFT),  $k = -\pi f/D_s$  is a modulation frequency,  $f$  is a modulation frequency index, and  $D'_s$  is the number of MS dimension in each feature dimension, where  $D'_s < D_s$ . We can control the highest modulation frequency considered in this criterion by adjusting the ratio of  $D_s$  to  $D'_s$ . In this paper, the MS is calculated utterance by utterance.

## 5.2. Objective Function for MS-Constrained Trajectory Training

We integrate the MS compensation into the trajectory training. The objective function consists of both the trajectory likelihood and the MS likelihood as follows:

$$L_{\text{mstrj}} = P(\mathbf{y}|\mathbf{X}, \hat{\mathbf{q}}, \lambda) P(\mathbf{s}(\mathbf{y})|\mathbf{X}, \hat{\mathbf{q}}, \lambda, \lambda_s)^{\omega_s T/D'_s}, \quad (18)$$

$$P(\mathbf{s}(\mathbf{y})|\mathbf{X}, \hat{\mathbf{q}}, \lambda, \lambda_s) = \mathcal{N}(\mathbf{s}(\mathbf{y}); \mathbf{s}(\hat{\mathbf{y}}_{\hat{\mathbf{q}}}), \Sigma_s), \quad (19)$$

where  $\lambda_s$  is a model parameter set of the MS, and  $\Sigma_s$  is a  $D'_s D$ -by- $D'_s D$  covariance matrix, and  $\omega_s$  is a weight of the MS likelihood. The trajectory likelihood and the MS likelihood are normalized by the ratio of the number of feature dimensions when  $\omega_s = 1$ .  $\Sigma_s^{-1}$  is represented as  $[\mathbf{p}_s^{(1)}, \dots, \mathbf{p}_s^{(d)}, \dots, \mathbf{p}_s^{(D)}]$ , where  $\mathbf{p}_s^{(d)}$  is  $D'_s D$ -by- $D$  matrix of which columns correspond to  $\mathbf{s}(d)$ . The MS likelihood works as a penalty term to alleviate the reduction of the temporal fluctuation of the generated parameter sequence.  $\Sigma_s$  is in advance estimated using training data.

## 5.3. Estimation of Model Parameters

The model parameters are estimated in the same way as for GV-constrained trajectory training. Let  $L_{\text{ms}}$  be the MS likelihood  $\mathcal{N}(\mathbf{s}(\mathbf{y}); \mathbf{s}(\hat{\mathbf{y}}_{\hat{\mathbf{q}}}), \Sigma_s)$ . The logarithm function of  $L_{\text{mstrj}}$  is given by:

$$\log L_{\text{mstrj}} = \log L_{\text{trj}} + \omega_s \frac{T}{D'_s} \log L_{\text{ms}}, \quad (20)$$

and the gradients of  $\log L_{\text{ms}}$  are given as:

$$\frac{\partial \log L_{\text{ms}}}{\partial \boldsymbol{\mu}} = \mathbf{S}_{\hat{\mathbf{q}}}^\top \mathbf{D}_{\hat{\mathbf{q}}}^{-1} \mathbf{W} \mathbf{R}_{\hat{\mathbf{q}}}^{-1} \mathbf{s}_{\hat{\mathbf{q}}}, \quad (21)$$

$$\frac{\partial \log L_{\text{ms}}}{\partial \Sigma^{-1}} = \mathbf{S}_{\hat{\mathbf{q}}}^\top \text{diag}_{3D}^{-1} [\mathbf{W} \mathbf{R}_{\hat{\mathbf{q}}}^{-1} \mathbf{s}_{\hat{\mathbf{q}}} (\mathbf{E}_{\hat{\mathbf{q}}} - \mathbf{W} \hat{\mathbf{y}}_{\hat{\mathbf{q}}})] \quad (22)$$

where

$$\mathbf{s}_{\hat{\mathbf{q}}} = [\mathbf{s}'_1{}^\top, \dots, \mathbf{s}'_t{}^\top, \dots, \mathbf{s}'_T{}^\top]^\top, \quad (23)$$

$$\mathbf{s}'_t = [s_t(1), \dots, s_t(d), \dots, s_t(D)]^\top, \quad (24)$$

$$s_t(d) = 2\mathbf{f}_t(d) \mathbf{p}_s^{(d)\top} (\mathbf{s}(\mathbf{y}) - \mathbf{s}(\hat{\mathbf{y}})), \quad (25)$$

$$\mathbf{f}_t(d) = [f_{t,d}(0), \dots, f_{t,d}(f), \dots, f_{t,d}(D'_s - 1)]^\top \quad (26)$$

$$f_{t,d}(f) = \hat{R}_{d,f} \cos kt + \hat{I}_{d,f} \sin kt, \quad (27)$$

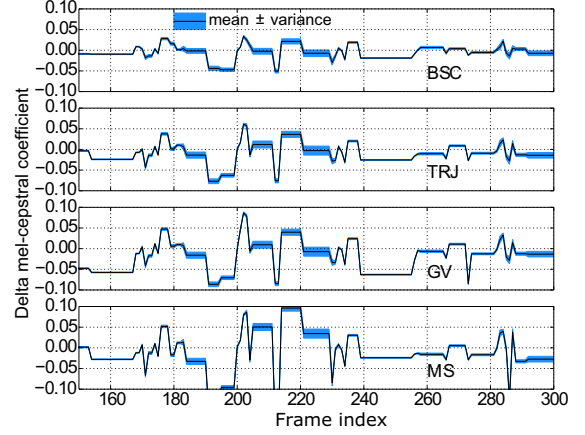


Figure 2: An example of the output probabilities of HMMs trained by several training algorithms.

where  $\hat{R}_{d,f}$  and  $\hat{I}_{d,f}$  are calculated using the  $d$ -th dimensional components of  $\hat{\mathbf{y}}_{\hat{\mathbf{q}}}$ .

## 5.4. Discussion

As reported in [23, 21], it is unnecessary to consider the MS in parameter generation because the HMM parameters are optimized to make the MS of the generated parameter sequence close to the natural one. Consequently, the basic computationally-efficient parameter generation algorithm is employed. This also enables to avoid the large footprint discussed in the parameter generation algorithm considering the MS [17].

Multi-Space probability Distribution (MSD)-HMM [25] is unsuitable for the implementation of the proposed algorithm for  $F_0$  contour because the MS modeling of the non-continuous sequence is inaccurate [17]. To solve this problem, this paper adopted continuous  $F_0$  modeling [26]. Moreover, 0-mean MS modeling is also adopted [17], which means that the MS is calculated from the  $F_0$  contour that the utterance-level  $F_0$  is subtracted.

Fig. 2 draws the output probabilities at each frame. We can see that the variance of trajectory training (“TRJ”) is slightly larger than that of basic training (“BSC”), and the mean of GV-constrained (“GV”) or MS-constrained (“MS”) trajectory training is significantly changed compared to “TRJ.” It is observed that the mean of “MS” tends to transit greatly from the neighbor HMM-state<sup>1</sup>.

# 6. Experimental Evaluation

## 6.1. Experimental Condition

We trained a context-dependent phoneme Hidden Semi-Markov Model (HSMM) [27] for an English male speaker “RMS” from the CMU ARCTIC database [28]. We used 593 sentences from subset A for training and 100 sentences from subset B for evaluation. Speech signals were sampled at 16 kHz. The shift length was set to 5 ms. The 0th-through-24th mel-cepstral coefficients were extracted as a spectral parameter and log-scaled  $F_0$  and 5 band-aperiodicity [29, 30] were extracted as excitation pa-

<sup>1</sup>Note that the frames that have same statistics correspond to the same HMM-state.

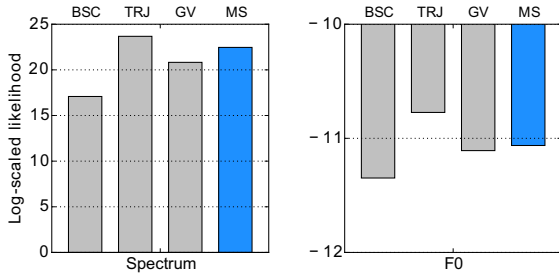


Figure 3: Trajectory likelihoods for the natural spectral parameter sequence and continuous  $F_0$  contour.

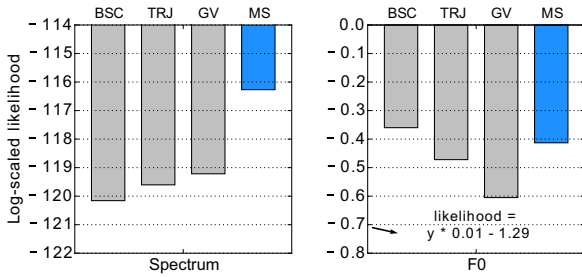


Figure 4: MS likelihoods for the natural spectral parameter sequence and continuous  $F_0$  contour.

rameters. The STRAIGHT analysis-synthesis system [31] was employed for parameter extraction and waveform generation. The feature vector consisted of spectral and excitation parameters and their delta and delta-delta features. 5-state left-to-right HMMs were used. The DFT length to calculate the MS was set to 2048. Diagonal covariance matrices were used in the HMM, the GV probability density function, and the MS probability density function. The likelihood weight  $\omega_v$  and  $\omega_s$  were set to 0.5 [21] and 1.0, respectively.  $D'_s$  for spectrum and  $F_0$  components were set to  $D_s/2 (= 50 \text{ Hz})$  [23] and  $D_s/10 (= 10 \text{ Hz})$  [13], respectively.

We compared the following training algorithms:

- BSC:** basic training ( $L_{\text{basic}}$ ) [2]
- TRJ:** trajectory training ( $L_{\text{trj}}$ ) [22]
- GV:** GV-constrained trajectory training ( $L_{\text{gvtrj}}$ ) [21]
- MS:** proposed MS-constrained trajectory training ( $L_{\text{mstrj}}$ )

These training algorithms were applied to spectral and  $F_0$  components, and the “BSC” training algorithm was applied to aperiodic component. First, “BSC” training is performed to estimate the HMM parameter set  $\lambda$ , then, “TRJ” training algorithm updates the parameters. “GV” and “MS” training algorithm updates after “TRJ” training. Note that the voiced/unvoiced regions of  $F_0$  contour never changes in all training algorithms.

The trajectory likelihood and the MS likelihood for the natural parameter trajectories of the evaluation data were firstly calculated to analyze the effect of the proposed algorithm. Then, the speech quality of the synthetic speech are evaluated in the perceptual evaluation.

## 6.2. Objective Evaluation

Fig. 3 and Fig. 4 illustrate the trajectory likelihood  $L_{\text{trj}}$  and the MS likelihood  $L_{\text{ms}}$  for the natural parameter trajectories of the

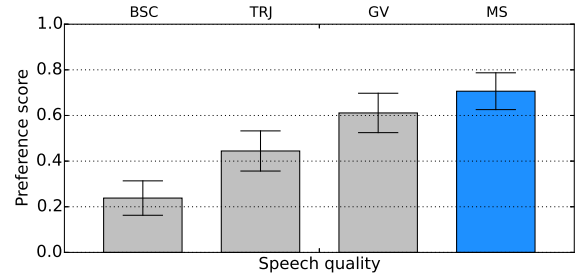


Figure 5: Preference scores on synthetic speech quality with 95% confidence interval.

evaluation data, respectively. Note that the trajectory likelihood is normalized by the total number of frames  $T$ , and the MS likelihood is similarly normalized by the number of dimension  $D'_s$ . We can see that the proposed MS-constrained trajectory training (“MS”) dramatically improves the MS likelihood compared to other training algorithms for the spectral component. On the other hand, the MS likelihood of “MS” is lower than that of the basic training algorithm (“BSC”) for  $F_0$  component. To discuss this result, we calculated the proposed objective function divided by  $T$ , and the values are  $-12.64$  and  $-12.35$  for “BSC” and “MS,” respectively. Therefore, the proposed objective function is certainly improved by the proposed algorithm. Therefore, we can confirm the implementation of the proposed algorithm for both spectral and  $F_0$  components.

## 6.3. Subjective Evaluation

A preference AB test was conducted by 6 listeners. We presented every pair of generated speech of 4 training algorithms in a random order, and we forced listeners to select speech sample that sounds better quality.

Fig. 5 illustrates the result. We can see that the proposed algorithm “MS” achieves the best score. This result demonstrates that the proposed MS-constrained trajectory training for spectral and  $F_0$  components makes it possible to improve the synthetic speech quality while preserving the computationally-efficient generation ability.

## 7. Conclusion

This paper proposed a novel training algorithm for HMM-based speech synthesis in order to produce the high-quality speech while preserving the computationally-efficient generation algorithm. The Modulation Spectrum (MS) have been integrated into the trajectory HMM training for both spectral and  $F_0$  components. The experimental results yielded the quality improvement in synthetic speech. As a future work, we combine both the proposed algorithm and the rich context modeling [32].

**Acknowledgements:** Part of this work was supported by JSPS KAKENHI Grant Number 26280060 and Grant-in-Aid for JSPS Fellows Grant Number 26 · 10354, and part of this work was executed under “JSPS Strategic Young Researcher Overseas Visits Program for Accelerating Brain Circulation.”

## 8. References

- [1] K. Tokuda, Y. Nankaku, T. Toda, H. Zen, J. Yamagishi, and K. Oura, "Speech synthesis based on hidden Markov models," *Proceedings of the IEEE*, vol. 101, no. 5, pp. 1234–1252, 2013.
- [2] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, and T. Kitamura, "Simultaneous modeling of spectrum, pitch and duration in HMM-based speech synthesis," in *Proc. EUROSPEECH*, Budapest, Hungary, Apr. 1999, pp. 2347–2350.
- [3] K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, and T. Kitamura, "Speech parameter generation algorithms for HMM-based speech synthesis," in *Proc. ICASSP*, Istanbul, Turkey, June 2000, pp. 1315–1318.
- [4] T. Yoshimura, T. Masuko, K. Tokuda, T. Kobayashi, and T. Kitamura, "Speaker interpolation for HMM-based speech synthesis system," *J. Acoust. Soc. Jpn. (E)*, vol. 21, no. 4, pp. 199–206, 2000.
- [5] J. Yamagishi and T. Kobayashi, "Average-voice-based speech synthesis using HSMM-based speaker adaptation and adaptive training," *IEICE Trans., Inf. and 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. and Syst.*, vol. E90-D, no. 9, pp. 1406–1413, 2007.
- [7] K. Shirota, K. Nakamura, K. Hashimoto, K. Oura, Y. Nankaku, and K. Tokuda, "Integration of speaker and pitch adaptive training for HMM-based singing voice synthesis," in *Proc. ICASSP*, Florence, Italy, May 2014, pp. 2578–2582.
- [8] J. Yamagishi, C. Veaux, S. King, and S. Renals, "Speech synthesis technologies for individuals with vocal disabilities: Voice banking and reconstruction," *Acoust. Sci. technol.*, vol. 33, pp. 1–5, 2012.
- [9] Y. Sagisaka, "Speech synthesis by rule using an optimal selection of non-uniform synthesis units," in *Proc. ICASSP*, New York, U.S.A., Apr. 1988, pp. 679–682.
- [10] H. Zen and A. Senior, "Deep mixture density networks for acoustic modeling in statistical parametric speech synthesis," in *Proc. ICASSP*, Florence, Italy, May 2014, pp. 3872–3876.
- [11] T. Kano, S. Takamichi, S. Sakti, T. T. G. Neubig, and S. Nakamura, "Generalizing continuous-space translation of paralinguistic information," in *Proc. INTERSPEECH*, Lyon, France, Aug 2013, pp. 2614–2618.
- [12] S. King and V. Karaiskos, "The blizzard challenge 2011," in *Proc. Blizzard Challenge workshop*, Turin, Italy, Sept. 2011.
- [13] S. Takamichi, T. Toda, G. Neubig, S. Sakti, and S. Nakamura, "A postfilter to modify modulation spectrum in HMM-based speech synthesis," in *Proc. ICASSP*, Florence, Italy, May 2014, pp. 290–294.
- [14] S. Takamichi, T. Toda, A. W. Black, and S. Nakamura, "Modulation spectrum-based post-filter for GMM-based voice conversion," in *Proc. APSIPA ASC*, Siem Reap, Cambodia, Dec. 2014.
- [15] —, "Modified modulation spectrum-based post-filter for HMM-based speech synthesis," in *Proc. GlobalSIP*, Atlanta, United States, Dec. 2014, pp. 710–714.
- [16] T. Toda and K. Tokuda, "A speech parameter generation algorithm considering global variance for HMM-based speech synthesis," *IEICE Trans.*, vol. E90-D, no. 5, pp. 816–824, 2007.
- [17] S. Takamichi, T. Toda, A. W. Black, and S. Nakamura, "Parameter generation algorithm considering modulation spectrum for HMM-based speech synthesis," in *Proc. ICASSP*, Brisbane, Australia, Apr. 2015.
- [18] T. Nose and A. Ito, "Analysis of spectral enhancement using global variance in HMM-based speech synthesis," in *Proc. INTERSPEECH*, MAX Atria, Singapore, May 2014, pp. 2917–2921.
- [19] H. Hwang, Y. Tsao, H. Wang, Y. Wang, and S. Chen, "Incorporating global variance in the training phase of GMM-based voice conversion," in *Proc. APSIPA*, Kaohsiung, Taiwan, Oct. 2013, pp. 1–6.
- [20] M. Shannon and W. Byrne, "Fast, low-artifact speech synthesis considering global variance," in *Proc. ICASSP*, Vancouver, Canada, May. 2013, pp. 7869–7873.
- [21] T. Toda and S. Young, "Trajectory training considering global variance for HMM-based speech synthesis," in *Proc. ICASSP*, Taipei, Taiwan, Aug. 2009, pp. 4025–4028.
- [22] H. Zen, K. Tokuda, and T. Kitamura, "Reformulating the HMM as a trajectory model by imposing explicit relationships between static and dynamic feature vector sequences," *Computer Speech and Language*, vol. 21, no. 1, pp. 153–173, Jan. 2007.
- [23] S. Takamichi, T. Toda, A. W. Black, and S. Nakamura, "Modulation spectrum-constrained trajectory training algorithm for GMM-based voice conversion," in *Proc. ICASSP*, Brisbane, Australia, Apr. 2015.
- [24] T. Toda, A. W. Black, and K. Tokuda, "Voice conversion based on maximum likelihood estimation of spectral parameter trajectory," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 15, no. 8, pp. 2222–2235, 2007.
- [25] K. Tokuda, T. Masuko, B. Miyazaki, and T. Kobayashi, "Multi-space probability distribution HMM," *IEICE Trans., Inf. and Syst.*, vol. E85-D, no. 3, pp. 455–464, 2002.
- [26] K. Yu and S. Young, "Continuous F0 modeling for HMM based statistical parametric speech synthesis," *IEEE Trans. Audio, Speech and Language*, vol. 19, no. 5, pp. 1071–1079, 2011.
- [27] H. Zen, K. Tokuda, T. K. T. Masuko, and T. Kitamura, "Hidden semi-Markov model based speech synthesis system," *IEICE Trans., Inf. and Syst.*, E90-D, no. 5, pp. 825–834, 2007.
- [28] J. Kominek and A. W. Black, "The CMU ARCTIC speech databases for speech synthesis research," in *Tech. Rep. CMU-LTI-03-177*, Language Technologies Institute, Carnegie Mellon University, Pittsburgh, U.S.A., 2003.
- [29] H. Kawahara, J. Estill, and O. Fujimura, "Aperiodicity extraction and control using mixed mode excitation and group delay manipulation for a high quality speech analysis, modification and synthesis system STRAIGHT," in *MAVEBA 2001*, Firentze, Italy, Sept. 2001, pp. 1–6.
- [30] Y. Ohtani, T. Toda, H. Saruwatari, and K. Shikano, "Maximum likelihood voice conversion based on GMM with STRAIGHT mixed excitation," in *Proc. INTERSPEECH*, Pittsburgh, U.S.A., Sep. 2006, pp. 2266–2269.
- [31] H. Kawahara, I. Masuda-Katsuse, and A. D. Cheveigne, "Restructuring 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 Commun.*, vol. 27, no. 3–4, pp. 187–207, 1999.
- [32] S. Takamichi, T. Toda, Y. Shiga, S. Sakti, G. Neubig, and S. Nakamura, "Parameter generation methods with rich context models for high-quality and flexible text-to-speech synthesis," *IEEE Journal of Selected Topics in Signal Processing*, vol. 8, no. 2, pp. 239–250, 2014.