

# Integration of Acoustic Modeling and Mel-cepstral analysis for HMM-based Speech Synthesis

著者 (英)	Kazuhiro Nakamura, Kei Hashimoto, Yoshihiko Nankaku, Keiichi Tokuda
journal or publication title	2013 IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2013)
page range	7883-7887
year	2013-05
URL	<a href="http://id.nii.ac.jp/1476/00004649/">http://id.nii.ac.jp/1476/00004649/</a>

doi: 10.1109/ICASSP.2013.6639199(<http://dx.doi.org/10.1109/ICASSP.2013.6639199>)

# INTEGRATION OF ACOUSTIC MODELING AND MEL-CEPSTRAL ANALYSIS FOR HMM-BASED SPEECH SYNTHESIS

*Kazuhiro Nakamura, Kei Hashimoto, Yoshihiko Nankaku, Keiichi Tokuda*

Department of Scientific and Engineering Simulation, Nagoya Institute of Technology, Nagoya, Japan

## ABSTRACT

In this paper, a novel approach for integrating acoustic modeling and mel-cepstral analysis is proposed. The aim of HMM-based speech synthesis is to model speech waveforms with a statistical model. However, the conventional techniques divide the modeling process into two steps: the frame by frame feature extraction step and the acoustic modeling step. Although it is reasonably effective, the deterioration of speech quality is caused by the divide of the objective function. In this paper, we propose an approach to modeling them as an integrative model and show the possibility of improving synthesized speech.

*Index Terms*— integrative model, HMM-based speech synthesis, acoustic modeling, mel-cepstral analysis, trajectory HMM

## 1. INTRODUCTION

HMM-based speech synthesis was proposed to enable machines to speak naturally like humans [1]. In this method, spectral and F0 features are extracted and modeled with a statistical technique. Recent large systems are often constructed with combinations of some modules that use statistical models. A famous example is language and acoustic models for speech recognition systems. Recently, the integration of these statistical models is an important research subject. In text-to-speech (TTS) systems, an approach integrating text analysis and speech synthesis modules was proposed [2]. It can optimize linguistic and acoustic models simultaneously. As the essential aim of TTS is to synthesize speech from given texts, the integration of these statistical models is a desirable future of TTS systems. In HMM-based speech synthesis, some heuristic methods have been used to extract spectral features from speech waveforms previously. A statistical method that consists of mel-cepstral analysis was proposed and is widely used [3]. In this method, mel-cepstral coefficients, i.e., frequency transformed cepstral coefficients, are regarded as statistical model parameters and estimated with a parametric method.

In the standard HMM, observation vector sequences are quasi-stationary, and each stationary part can be represented by a state of the HMM. The statistics of each state do not change dynamically, and intra-state time-dependency cannot be represented. Therefore, a technique that augments the dimensionality of an acoustic static feature vector by appending its dynamic feature vectors is widely used. A trajectory model, named a “trajectory HMM” [4], was derived by reformulating the HMM. The standard HMM with static and dynamic features allows inconsistent statistics between the model parameters for static and dynamic features. By imposing the explicit relationship between them, the standard HMM is naturally translated into a trajectory model. The trajectory HMM can overcome the limitations in the standard HMM framework without any additional parameters.

In the mel-cepstral analysis and the acoustic modeling, two statistical models are estimated independently. Since the aim of these

two techniques is to model speech with a statistical model, in this paper, we propose integrating the feature extraction and the acoustic modeling by using a probabilistic representation of extracted features. This approach can model speech waveforms directly without having any intermediate representation, and it can be regarded as a generative model of speech waveforms.

In this model, the spectral extraction process and the spectral modeling process are simultaneously optimized for the given speech waveforms. This framework is similar to the vocal tract transfer function (VTTF) estimation of a speech signal based on a factor analyzed (FA) trajectory HMM [5]. Mel-cepstral coefficients were regarded as factors, and by using the time-varying factor loading matrix, the harmonic components were represented with the FA method. In another approach, the mel-cepstral analysis was integrated into the Gaussian mixture model (GMM) for modeling a quasi-stationary Gaussian process [6]. It can represent mel-cepstral coefficients stochastically with mixture weights of GMM.

The rest of this paper is organized as follows. Section 2 summarizes HMM-based speech synthesis, including the mel-cepstral analysis and the trajectory HMM. In Section 3, the integration algorithm of the mel-cepstral analysis and the acoustic modeling is derived. Experimental results are presented in Section 4. Concluding remarks and future plans are presented in the final section.

## 2. HMM-BASED SPEECH SYNTHESIS

In HMM-based TTS training, spectral envelope, fundamental frequency, and duration are modeled simultaneously by using the corresponding HMMs. Mel-cepstral coefficients  $c$  are widely used as spectral features. They are regarded as statistical model parameters and estimated from a given input signal,  $x$ , in the maximum likelihood (ML) sense:

$$\hat{c} = \operatorname{argmax}_c P(x|c) \quad (1)$$

Extracted mel-cepstral coefficients are used for training HMMs with dynamic (“delta” and “delta-delta”) feature constraints. A speech waveform is finally synthesized from the generated spectral and excitation parameters via the source-filter based production model. Recently, the trajectory HMM was derived by reformulating the HMM. The model parameters  $\Lambda$  are trained in the ML sense by using the static features:

$$\hat{\Lambda} = \operatorname{argmax}_{\Lambda} P(c|\Lambda) \quad (2)$$

### 2.1. MEL-CEPSTRAL ANALYSIS

The synthesis filter  $H(z)$  is represented by mel-cepstral coefficients  $c = [c(0), \dots, c(M)]^T$ <sup>1</sup> defined as frequency-transformed cepstral coefficients:

<sup>1</sup>In section 2.1,  $x$  and  $c$  correspond to not an utterance but a frame. The frame index  $t$  is abbreviated.

$$H(z) = \exp \sum_{m=0}^M c(m) \tilde{z}^{-m} \quad (3)$$

$$\tilde{z}^{-1} = \frac{z^{-1} - \alpha}{1 - \alpha z^{-1}}, \quad |\alpha| < 1 \quad (4)$$

where  $\alpha$  is a frequency warping parameter. If  $\alpha = 0$ , mel-cepstral coefficients are equivalent to cepstral coefficients.

For a given input signal,  $\mathbf{x} = [x(0), \dots, x(N-1)]^\top$ , the mel-cepstral coefficients are determined by minimizing a spectral evaluation function with respect to  $\mathbf{c}$  [7],

$$E(\mathbf{x}, \mathbf{c}) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \{\exp R(\omega) - R(\omega) - 1\} d\omega \quad (5)$$

where

$$R(\omega) = \log I_N(\omega) - \log |H(e^{j\omega})|^2 \quad (6)$$

and  $I_N(\omega)$  is the modified periodogram of weakly stationary process  $x(n)$  with a time window  $w(n)$  of length  $N$ :

$$I_N(\omega) = \frac{\left| \sum_{n=0}^{N-1} w(n) x(n) e^{-j\omega n} \right|^2}{\sum_{n=0}^{N-1} w^2(n)} \quad (7)$$

Mel-cepstral coefficients are determined easily by using an iterative algorithm (e.g., the Newton-Raphson method) because  $E(\mathbf{x}, \mathbf{c})$  is convex with respect to  $\mathbf{c}$ .

When  $x(n)$  is assumed to be a zero-mean Gaussian process, the likelihood can be approximated by

$$P(\mathbf{x}|\mathbf{c}) \simeq \exp \left[ -\frac{N}{2} \left[ \log(2\pi) + \frac{1}{2\pi} \int_{-\pi}^{\pi} \left\{ \log |H(e^{j\omega})|^2 + \frac{I_N(\omega)}{|H(e^{j\omega})|^2} \right\} d\omega \right] \right] \quad (8)$$

and, accordingly, minimization of  $E(\mathbf{x}, \mathbf{c})$  corresponds to the maximization of  $P(\mathbf{x}|\mathbf{c})$ .

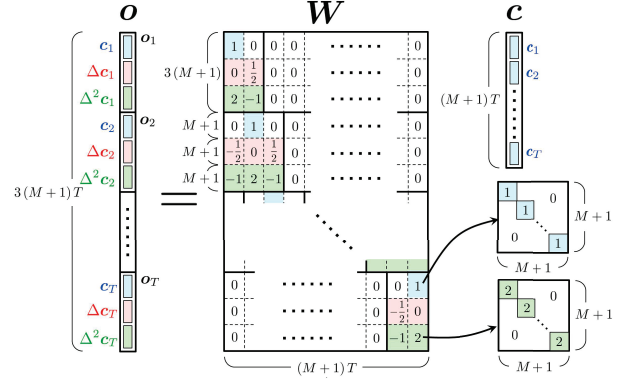
## 2.2. TRAJECTORY HMM

Let a spectral feature vector sequence be  $\mathbf{o} = [\mathbf{o}_1^\top, \dots, \mathbf{o}_T^\top]^\top$ , where  $\mathbf{o}_t = [\mathbf{c}_t^\top, \Delta \mathbf{c}_t^\top, \Delta^2 \mathbf{c}_t^\top]^\top$  includes not only static but also dynamic features. Mel-cepstral coefficients  $c_t$  are a  $M+1$  dimensional vector, and  $T$  is the number of frames. In the conventional model, the probability density of  $\mathbf{o}$  is shown as  $P(\mathbf{o}|\mathbf{q}, \mathbf{\Lambda})$ , where  $\mathbf{q} = (q_1, q_2, \dots, q_T)$  is a state sequence. However, the conventional model is mathematically improper in the sense of statistical modeling. In this model, the static and dynamic features are modeled as independent statistical variables. When it is used as a generative model, it allows inconsistent static and dynamic features. By imposing an explicit relationship between static and dynamic features, which is given by  $\mathbf{o} = \mathbf{W}\mathbf{c}$ , where  $\mathbf{W}$  is a  $3(M+1)T \times (M+1)T$  window matrix as shown in Fig. 1, the conventional HMM is reformed as the trajectory HMM as:

$$P(\mathbf{c}|\mathbf{\Lambda}) = \sum_{\forall \mathbf{q}} P(\mathbf{c}|\mathbf{q}, \mathbf{\Lambda}) P(\mathbf{q}|\mathbf{\Lambda}) \quad (9)$$

$$P(\mathbf{c}|\mathbf{q}, \mathbf{\Lambda}) = \mathcal{N}(\mathbf{c}|\bar{\mathbf{c}}_{\mathbf{q}}, \mathbf{P}_{\mathbf{q}}) = \frac{1}{Z} P(\mathbf{o}|\mathbf{q}, \mathbf{\Lambda}) \quad (10)$$

$$P(\mathbf{q}|\mathbf{\Lambda}) = P(q_1|\mathbf{\Lambda}) \prod_{t=2}^T P(q_t|q_{t-1}, \mathbf{\Lambda}) \quad (11)$$



**Fig. 1.** Example of the relationship between the static feature vector sequence  $\mathbf{c}$  and the speech parameter vector sequence  $\mathbf{o}$  in a matrix form

where  $Z$  is a normalization term. In Eq. (10),  $\bar{\mathbf{c}}_{\mathbf{q}}$  and  $\mathbf{P}_{\mathbf{q}}$  are the  $(M+1)T \times 1$  mean vector and the  $(M+1)T \times (M+1)T$  temporal utterance covariance matrix given by  $\mathbf{q}$ , respectively. They are given by

$$\mathbf{R}_{\mathbf{q}} \bar{\mathbf{c}}_{\mathbf{q}} = \mathbf{r}_{\mathbf{q}} \quad (12)$$

$$\mathbf{R}_{\mathbf{q}} = \mathbf{W}^\top \mathbf{\Sigma}_{\mathbf{q}}^{-1} \mathbf{W} = \mathbf{P}_{\mathbf{q}}^{-1} \quad (13)$$

$$\mathbf{r}_{\mathbf{q}} = \mathbf{W}^\top \mathbf{\Sigma}_{\mathbf{q}}^{-1} \boldsymbol{\mu}_{\mathbf{q}} \quad (14)$$

$$\boldsymbol{\mu}_{\mathbf{q}} = [\boldsymbol{\mu}_{q_1}^\top, \dots, \boldsymbol{\mu}_{q_T}^\top]^\top \quad (15)$$

$$\boldsymbol{\mu}_i = [\boldsymbol{\mu}_i^\top, \Delta \boldsymbol{\mu}_i^\top, \Delta^2 \boldsymbol{\mu}_i^\top]^\top, \quad i = 1, \dots, N \quad (16)$$

$$\mathbf{\Sigma}_{\mathbf{q}} = \text{diag} [\boldsymbol{\Sigma}_{q_1}^\top, \dots, \boldsymbol{\Sigma}_{q_T}^\top]^\top \quad (17)$$

$$\boldsymbol{\Sigma}_i = \text{diag} [\boldsymbol{\Sigma}_i^\top, \Delta \boldsymbol{\Sigma}_i^\top, \Delta^2 \boldsymbol{\Sigma}_i^\top]^\top, \quad i = 1, \dots, N \quad (18)$$

where  $N$  is the total number of state output PDFs, and  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\Sigma}_i$  are the  $3(M+1)T \times 1$  mean vector and the  $3(M+1)T \times 3(M+1)T$  covariance matrix associated with the  $i$ -th state, respectively. The elements of  $\mathbf{W}$  are given as regression window coefficients to calculate delta and delta-delta features as:

$$\Delta^d c_t = \sum_{\tau=-L_-^{(d)}}^{L_+^{(d)}} w^{(d)}(\tau) c_{t+\tau}, \quad d = 1, 2 \quad (19)$$

$$\mathbf{W} = [\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_T]^\top \otimes \mathbf{I}_{(M+1) \times (M+1)} \quad (20)$$

$$\mathbf{W}_t = [\mathbf{w}_t^{(0)}, \mathbf{w}_t^{(1)}, \dots, \mathbf{w}_t^{(D-1)}] \quad (21)$$

$$\mathbf{w}_t^{(d)} = \left[ \underbrace{0, \dots, 0}_{t-L_-^{(d)}-1}, w^{(d)}(-L_-^{(d)}), \dots, w^{(d)}(0), \dots, w^{(d)}(L_+^{(d)}), \underbrace{0, \dots, 0}_{T-(t+L_+^{(d)})} \right]^\top, \quad d = 0, 1, 2 \quad (22)$$

where  $L_-^{(0)} = L_+^{(0)} = 0$ ,  $\mathbf{w}^{(0)} = 1$ , and  $\otimes$  denotes the Kronecker product for matrices.

Note that  $\mathbf{c}$  is modeled by Gaussian distributions whose dimensionality is  $(M+1)T$ , and the covariance matrices  $\mathbf{P}_{\mathbf{q}}$  are generally full. As a result, the trajectory HMM can overcome the deficiencies of the HMM. It is also noted that the parameterization of

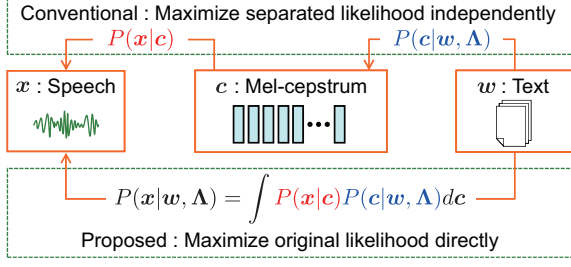


Fig. 2. Basic idea of the proposed approach

the trajectory HMM is completely the same as that of the HMM with the same model topology.

### 3. INTEGRATION ALGORITHM OF ACOUSTIC MODELING AND MEL-CEPSTRAL ANALYSIS

Figure 2 shows the difference between the conventional and the proposed approaches. In the proposed approach, we integrate the statistical mel-cepstral model  $P(x|c)$  and the statistical acoustic model  $P(c|\Lambda)$  as:

$$\begin{aligned} \hat{\Lambda} &= \operatorname{argmax}_{\Lambda} P(x|\Lambda) \\ &= \operatorname{argmax}_{\Lambda} \int P(x, c|\Lambda) dc \\ &= \operatorname{argmax}_{\Lambda} \int P(x|c) P(c|\Lambda) dc \end{aligned} \quad (23)$$

In the proposed method, the mel-cepstral coefficients are represented as a random variable, and the dynamic feature vectors cannot be calculated from static features. Therefore, a modeling technique including the extraction of dynamic features, like trajectory HMM, should be introduced. As far as the conventional method, the standard mel-cepstral analysis can be assumed as extracting 1-best mel-cepstral coefficients.

To train the proposed model, a lower bound of log marginal likelihood  $\mathcal{F}$  is maximized instead of the true likelihood. The lower bound  $\mathcal{F}$  is defined by using Jensen's inequality:

$$\begin{aligned} \mathcal{L}(x|\Lambda) &= \log P(x|\Lambda) \\ &= \log \sum_{\forall q} \int P(x|c) P(c, q|\Lambda) dc \\ &= \log \sum_{\forall q} \int Q(c, q) \frac{P(x|c) P(c, q|\Lambda)}{Q(c, q)} dc \\ &= \log \sum_{\forall q} \int Q(c) Q(q) \frac{P(x|c) P(c, q|\Lambda)}{Q(c) Q(q)} dc \\ &\geq \sum_{\forall q} \int Q(c) Q(q) \log \frac{P(x|c) P(c, q|\Lambda)}{Q(c) Q(q)} dc \\ &= \mathcal{F} \end{aligned} \quad (24)$$

To overcome the difficulty of optimization, we assume that  $c$  and  $q$  are independent and that posterior distribution  $Q(c)$  and  $Q(q)$  approximate the true posterior distributions. The optimal posterior distributions can be obtained by maximizing the original objective

function  $\mathcal{F}$  with the variational method as:

$$Q(c) = \frac{1}{Z_c} P(x|c) \exp \sum_{\forall q} Q(q) \log P(c|q, \Lambda) \quad (25)$$

$$Q(q) = \frac{1}{Z_q} P(q|\Lambda) \exp \int Q(c) \log P(c|q, \Lambda) dc \quad (26)$$

where  $Z_c$  and  $Z_q$  are the normalization terms of  $Q(c)$  and  $Q(q)$ , respectively. These optimizations can be effectively performed by iterative calculations as the Expectation and Maximization (EM) algorithm, which increases the value of objective function  $\mathcal{F}$  at each iteration until convergence.

#### 3.1. Posterior Probabilities of Mel-cepstral coefficients

It is difficult to integrate  $Q(c)$  with respect to  $c$ , so we approximate  $Q(c)$  as a Gaussian probability distribution. By using a Laplace approximation,  $Q(c)$  is represented as:

$$Q(c) \simeq \mathcal{N}(c|\tilde{c}, A^{-1}) \quad (27)$$

$$\tilde{c} = \operatorname{argmax}_c Q(c) \quad (28)$$

$$A = \frac{N}{2} H|_{c=\tilde{c}} + \sum_{\forall q} Q(q) P_q^{-1} \quad (29)$$

where

$$H = -\frac{2}{N} \frac{\partial^2}{\partial c \partial c^T} \log P(x|c) \quad (30)$$

$$= \operatorname{diag} \left( \left[ H_1^T, H_2^T, \dots, H_T^T \right]^T \right) \quad (31)$$

$H_t$  is the Hessian matrix of the spectral evaluation function at time  $t$ :

$$H_t = \frac{\partial^2}{\partial c_t \partial c_t^T} E(x_t, c_t) = -\frac{2}{N} \frac{\partial^2}{\partial c_t \partial c_t^T} \log P(x_t|c_t) \quad (32)$$

In the standard mel-cepstral analysis, mel-cepstral coefficients for each frame can be estimated frame by frame independently. However, in the proposed method, mel-cepstral coefficients for all utterances should be estimated simultaneously. Thus, a large computation cost is required for this Newton-Raphson method.

#### 3.2. Posterior Probabilities of State Sequences

The expectation with respect to  $c$  in Eq. (26) is given by

$$\begin{aligned} &\int Q(c) \log P(c|q, \Lambda) dc \\ &= \log \mathcal{N}(\tilde{c}|\tilde{c}_q, P_q) - \frac{1}{2} \operatorname{tr}(\mathbf{R}_q A^{-1}) \end{aligned} \quad (33)$$

where the matrix  $\mathbf{R}_q$  is a positive definite  $(4L(M+1)+1)$ -diagonal band symmetric matrix. Thus,  $\mathbf{R}_q$  can be decomposed into its Cholesky factorization:

$$\mathbf{R}_q = \mathbf{U}_q^T \mathbf{U}_q \quad (34)$$

where  $\mathbf{U}_q$  is an upper-triangular  $(2L(M+1)+1)$ -diagonal matrix. Elements of  $\mathbf{U}_q$  are calculated in a recursive manner and depend only on substates from time 1 to  $t+2L$ .

$$\begin{aligned} (\mathbf{R}_q A^{-1})^{(t,t)} &= (\mathbf{U}_q^T \mathbf{U}_q A^{-1})^{(t,t)} = (\mathbf{U}_q A^{-1} \mathbf{U}_q^T)^{(t,t)} \\ &= \sum_{i=t}^{t+2L} \sum_{j=t}^{t+2L} \mathbf{U}_q^{(t,i)} (A^{-1})^{(i,j)} \mathbf{U}_q^{(t,j)} \end{aligned} \quad (35)$$

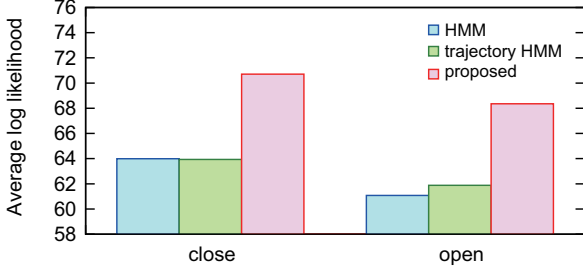


Fig. 3. Log likelihood per frame for closed and opened data sets

Thus, the delayed decision Viterbi algorithm [4] can be applied to the proposed method.

### 3.3. Update Model Parameters

Model parameters  $\mathbf{m}$  and  $\phi$  are defined by concatenating the mean vectors and covariance matrices of all unique Gaussian components in the model set as:

$$\mathbf{m} = [\boldsymbol{\mu}_1^\top, \boldsymbol{\mu}_2^\top, \dots, \boldsymbol{\mu}_N^\top]^\top \quad (36)$$

$$\phi = [\boldsymbol{\Sigma}_1^\top, \boldsymbol{\Sigma}_2^\top, \dots, \boldsymbol{\Sigma}_N^\top]^\top \quad (37)$$

where  $\boldsymbol{\mu}_n$  and  $\boldsymbol{\Sigma}_n$  are the mean vector and covariance matrix of the  $n$ -th unique Gaussian component in the model set, and  $N$  is the total number of Gaussian components in the model set, respectively.

By setting the partial derivative of  $\mathcal{F}$  with respect to  $\mathbf{m}$  to 0, a set of linear equations for determining  $\mathbf{m}$  maximizing  $\mathcal{F}$  are obtained as:

$$\sum_{\forall q} Q(q) \mathbf{S}_q^\top \mathbf{W} \mathbf{P}_q \mathbf{W}^\top \mathbf{S}_q \boldsymbol{\Phi}^{-1} \mathbf{m} = \sum_{\forall q} Q(q) \mathbf{S}_q^\top \mathbf{W} \tilde{\mathbf{c}} \quad (38)$$

where

$$\boldsymbol{\mu}_q = \mathbf{S}_q \mathbf{m} \quad (39)$$

$$\boldsymbol{\Phi}^{-1} = \text{diag}(\phi) \quad (40)$$

$$\boldsymbol{\Sigma}_q^{-1} = \text{diag}(\mathbf{S}_q \phi) \quad (41)$$

$$\mathbf{S}_q \boldsymbol{\Phi}^{-1} = \boldsymbol{\Sigma}_q^{-1} \mathbf{S}_q \quad (42)$$

In the above equations,  $\mathbf{S}_q$  is a  $3(M+1)T \times 3(M+1)T$  matrix whose elements are 0 or 1 determined by the Gaussian component sequence  $q$ .

For maximizing  $\mathcal{F}$  with respect to  $\phi$ , a gradient method is applied by using its partial derivative

$$\begin{aligned} \frac{\partial \mathcal{F}}{\partial \phi} = \sum_{\forall q} Q(q) \left[ \frac{1}{2} \mathbf{S}_q^\top \text{diag}^{-1} \left\{ \mathbf{W} \mathbf{P}_q \mathbf{W}^\top - \mathbf{W} \mathbf{A}^{-1} \mathbf{W}^\top \right. \right. \\ \left. \left. - \mathbf{W} \tilde{\mathbf{c}} \tilde{\mathbf{c}}^\top \mathbf{W}^\top + 2\boldsymbol{\mu}_q \tilde{\mathbf{c}}^\top \mathbf{W}^\top \right. \right. \\ \left. \left. + \mathbf{W} \tilde{\mathbf{c}}_q \tilde{\mathbf{c}}_q^\top \mathbf{W}^\top - 2\boldsymbol{\mu}_q \tilde{\mathbf{c}}_q^\top \mathbf{W}^\top \right\} \right] \quad (43) \end{aligned}$$

because Eq. (43) is not a quadratic function of  $\phi$ .

## 4. EXPERIMENTS

### 4.1. Experimental Conditions

To evaluate the effectiveness of the proposed method, objective comparison tests of likelihood and subjective comparison test of Mean opinion Score (MOS) were conducted. For training, 50 sentences of the phonetically balanced 503 sentences from the ATR

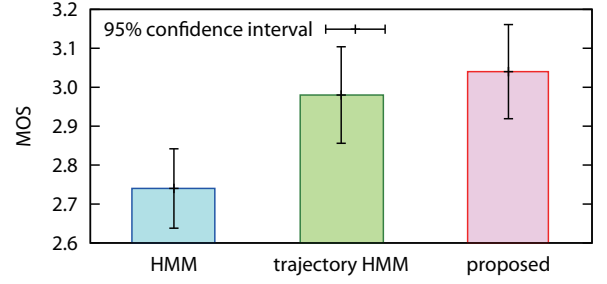


Fig. 4. Mean opinion scores for synthesized speech obtained by standard HMM, trajectory HMM and proposed model.

Japanese speech database (Set B) uttered by a male speaker M001 in Nitech, were used. Fifty other sentences were used for evaluating. The speech data was recorded at 48 kHz and windowed at a frame rate of 5-ms by using a 25-ms Hamming window. The windowed waveforms were used as the input data in the proposed method, and 35 mel-cepstral coefficients, which include the zero coefficient estimated with the standard mel-cepstral analysis technique were used in the conventional method. The dimension of the hidden mel-cepstral coefficients of the proposed method was set to the same as that of the conventional method. The frequency warping parameter  $\alpha$  was set to 0.55. Each state output probability distribution was modeled by a single Gaussian distribution with a diagonal covariance matrix. In the subjective test, 10 subjects were asked to rate the naturalness of the synthesized speech voices on a MOS with a scale from 1 (poor) to 5 (good). Fifteen randomly selected sentences were presented to each subject. The experiments were carried out in a sound-proof room.

### 4.2. Experimental Results

After estimating the standard HMMs, the proposed models were re-estimated by using the standard HMMs as their initial models in accordance with the training procedure described in Section 3. Figure 3 shows the average log likelihood per frame for a training data set (close) and test data set (open) to compare the proposed models with the standard HMMs. The proposed models outperformed the baseline standard HMMs for both data sets. This result for the test data set especially indicates the possibility of improving synthesized speech. Figure 4 shows the subjective listening results. In Figure 4, the MOS of the proposed technique was higher than that of the standard HMMs and equivalent to or higher than that of the trajectory HMMs. These results mean that the proposed technique can represent speech waveforms better by modeling them directly, even though baseline and proposed models have the same number of parameters.

## 5. CONCLUSION

In this paper, we defined a novel kind of acoustic model for modeling speech waveforms directly by integrating the mel-cepstral analysis and the acoustic modeling. In experiments, the objective and subjective evaluation scores of proposed models were equivalent to or higher than the baseline models. These results indicate the possibility of improving the quality of synthesized speech. Experiments on larger data sets will be future work.

## 6. ACKNOWLEDGMENTS

The research leading to these results was partly funded by the Core Research for Evolutionary Science and Technology (CREST) from the Japan Science and Technology Agency (JST).

## 7. REFERENCES

- [1] K. T. Masuko, T. Tokuda, K. Kobayashi, and S. Imai, "Speech synthesis from HMMs using dynamic features," *Proceedings of ICASSP*, pp. 389–392, 1996.
- [2] K. Oura, Y. Nankaku, T. Toda, and K. Tokuda, "Simultaneous acoustic, prosodic, and phrasing model training for TTS conversion systems," *Proceedings of ISCSLP2008*, pp. 1–4, 2008.
- [3] T. Fukada, K. Tokuda, T. Kobayashi, and S. Imai, "An adaptive algorithm for mel-cepstral analysis of speech," *Proceedings of ICASSP*, vol. 1, pp. 137–140, 1992.
- [4] K. Tokuda, H. Zen, and T. Kitamura, "Trajectory modeling based on HMMs with the explicit relationship between static and dynamic features," *Proceedings of Eurospeech*, pp. 865–868, 2003.
- [5] T. Toda and K. Tokuda, "Statistical approach to vocal tract transfer function estimation based on factor analyzed trajectory HMM," *Proceedings of ICASSP*, pp. 3925–3928, 2008.
- [6] T. Takahashi, K. Tokuda, T. Kobayashi, and T. Kitamura, "Mixture density models based on mel-cepstral representation of gaussian process," *Proceedings of IEICE*, pp. 1971–1978, 2003.
- [7] S. Imai and C. Furuichi, "Unbiased estimator of log spectrum and its application to speech signal processing," *Proceedings of EURASIP*, pp. 203–206, 1988.