

SPEECH PARAMETER GENERATION ALGORITHMS FOR HMM-BASED SPEECH SYNTHESIS

Keiichi Tokuda¹, Takayoshi Yoshimura¹, Takashi Masuko², Takao Kobayashi², Tadashi Kitamura¹,

¹Department of Computer Science, Nagoya Institute of Technology, Nagoya, 466-8555 Japan

²Interdisciplinary Graduate School of Science and Engineering, Tokyo Institute of Technology, Yokohama, 226-8502 Japan

Email: {tokuda,yossie,kitamura}@ics.nitech.ac.jp, {masuko,tkobayas}@ip.titech.ac.jp

ABSTRACT

This paper derives a speech parameter generation algorithm for HMM-based speech synthesis, in which speech parameter sequence is generated from HMMs whose observation vector consists of spectral parameter vector and its dynamic feature vectors. In the algorithm, we assume that the state sequence (state and mixture sequence for the multi-mixture case) or a part of the state sequence is unobservable (i.e., hidden or latent). As a result, the algorithm iterates the forward-backward algorithm and the parameter generation algorithm for the case where state sequence is given. Experimental results show that by using the algorithm, we can reproduce clear formant structure from multi-mixture HMMs as compared with that produced from single-mixture HMMs.

1. INTRODUCTION

The increasing availability of large speech databases makes it possible to construct speech synthesis systems, which are referred to as data-driven or corpus-based approach, by applying statistical learning algorithms. These systems, which can be automatically trained, not only generate natural and high quality synthetic speech but also can reproduce voice characteristics of the original speaker.

For constructing such a system, the use of hidden Markov models (HMMs) has arisen largely. HMMs have successfully been applied to modeling the sequence of speech spectra in speech recognition systems, and the performance of HMM-based speech recognition systems have been improved by techniques which utilize the flexibility of HMMs: context-dependent modeling, dynamic feature parameters, mixtures of Gaussian densities, tying mechanism, speaker and environment adaptation techniques. HMM-based approaches to speech synthesis can be categorized as follows:

1. Transcription and segmentation of speech database [1].
2. Construction of inventory of speech segments [2], [3].
3. Run-time selection of multiple instances of speech segments [4], [5].
4. Speech synthesis from HMMs themselves [6]–[9].

In approaches 1–3, by using a waveform concatenation algorithm, e.g., PSOLA algorithm, a high quality synthetic speech could be produced. However, to obtain various voice characteristics, large amounts of speech data are necessary, and it is difficult to collect, segment, and store these data. On the other hand, in approach 4,

voice characteristics of synthetic speech can be changed by transforming HMM parameters appropriately. From this point of view, we have proposed parameter generation algorithms [10], [11] for HMM-based speech synthesis, and constructed a speech synthesis system [8], [9]. Actually, we have shown that we can change voice characteristics of synthetic speech by applying a speaker adaptation technique [12], [13] or a speaker interpolation technique [14]. The main feature of the system is the use of dynamic feature: by inclusion of dynamic coefficients in the feature vector, the dynamic coefficients of the speech parameter sequence generated in synthesis are constrained to be realistic, as defined by the parameters of the HMMs.

This paper describes algorithms for speech parameter generation from HMMs. In the previously proposed algorithms [10], [11] for speech parameter generation, state sequence is assumed to be given, or determined based on a maximum likelihood criterion. In this paper, we derive a new algorithm in which we assume that the state sequence (state and mixture sequence for the multi-mixture case) or a part of the state sequence is unobservable (i.e., hidden or latent), and show examples of speech spectra generated by the proposed algorithm.

The rest of this paper is organized as follows. Section 2 summarizes the previously proposed algorithms for speech parameter generation, and derives a new algorithm. Experimental results are shown in Section 3. Concluding remarks are presented in the final section.

2. SPEECH PARAMETER GENERATION BASED ON MAXIMUM LIKELIHOOD CRITERION

For a given continuous mixture HMM λ , we derive an algorithm for determining speech parameter vector sequence

$$\mathbf{O} = [\mathbf{o}_1^\top, \mathbf{o}_2^\top, \dots, \mathbf{o}_T^\top]^\top \quad (1)$$

in such a way that

$$P(\mathbf{O}|\lambda) = \sum_{\text{all } \mathbf{Q}} P(\mathbf{O}, \mathbf{Q}|\lambda) \quad (2)$$

is maximized with respect to \mathbf{O} , where

$$\mathbf{Q} = \{(q_1, i_1), (q_2, i_2), \dots, (q_T, i_T)\} \quad (3)$$

is the state and mixture sequence, i.e., (q, i) indicates the i -th mixture of state q . We assume that the speech parameter vector \mathbf{o}_t consists of the static feature vector $\mathbf{c}_t = [c_t(1), c_t(2), \dots, c_t(M)]^\top$

This work was partially supported by the Ministry of Education, Science, Sports and Culture Japan, Grant-in-Aid for Scientific Research (B) 2, 1055125, 1998, Encouragement of Young Scientists (A), 10780226, 1998.

(e.g., cepstral coefficients) and dynamic feature vectors $\Delta \mathbf{c}_t, \Delta^2 \mathbf{c}_t$ (e.g., delta and delta-delta cepstral coefficients, respectively), that is, $\mathbf{o}_t = [\mathbf{c}_t^\top, \Delta \mathbf{c}_t^\top, \Delta^2 \mathbf{c}_t^\top]^\top$, where the dynamic feature vectors are calculated by

$$\Delta \mathbf{c}_t = \sum_{\tau=-L_-^{(1)}}^{L_+^{(1)}} w^{(1)}(\tau) \mathbf{c}_{t+\tau} \quad (4)$$

$$\Delta^2 \mathbf{c}_t = \sum_{\tau=-L_-^{(2)}}^{L_+^{(2)}} w^{(2)}(\tau) \mathbf{c}_{t+\tau}. \quad (5)$$

We have derived algorithms [10], [11] for solving the following problems:

Case 1. For given λ and \mathbf{Q} , maximize $P(\mathbf{O}|\mathbf{Q}, \lambda)$ with respect to \mathbf{O} under the conditions (4), (5).

Case 2. For a given λ , maximize $P(\mathbf{O}, \mathbf{Q}|\lambda)$ with respect to \mathbf{Q} and \mathbf{O} under the conditions (4), (5).

In this section, we will review the above algorithms and derive an algorithm for the problem:

Case 3. For a given λ , maximize $P(\mathbf{O}|\lambda)$ with respect to \mathbf{O} under the conditions (4), (5).

2.1. Case 1: Maximizing $P(\mathbf{O}|\mathbf{Q}, \lambda)$ with respect to \mathbf{O}

First, consider maximizing $P(\mathbf{O}|\mathbf{Q}, \lambda)$ with respect to \mathbf{O} for a fixed state and mixture sequence \mathbf{Q} . The logarithm of $P(\mathbf{O}|\mathbf{Q}, \lambda)$ can be written as

$$\log P(\mathbf{O}|\mathbf{Q}, \lambda) = -\frac{1}{2} \mathbf{O}^\top \mathbf{U}^{-1} \mathbf{O} + \mathbf{O}^\top \mathbf{U}^{-1} \mathbf{M} + K \quad (6)$$

where

$$\mathbf{U}^{-1} = \text{diag} [\mathbf{U}_{q_1, i_1}^{-1}, \mathbf{U}_{q_2, i_2}^{-1}, \dots, \mathbf{U}_{q_T, i_T}^{-1}] \quad (7)$$

$$\mathbf{M} = [\boldsymbol{\mu}_{q_1, i_1}^\top, \boldsymbol{\mu}_{q_2, i_2}^\top, \dots, \boldsymbol{\mu}_{q_T, i_T}^\top]^\top \quad (8)$$

$\boldsymbol{\mu}_{q_t, i_t}$ and \mathbf{U}_{q_t, i_t} are the $3M \times 1$ mean vector and the $3M \times 3M$ covariance matrix, respectively, associated with i_t -th mixture of state q_t , and the constant K is independent of \mathbf{O} .

It is obvious that $P(\mathbf{O}|\mathbf{Q}, \lambda)$ is maximized when $\mathbf{O} = \mathbf{M}$ without the conditions (4), (5), that is, the speech parameter vector sequence becomes a sequence of the mean vectors. Conditions (4), (5) can be arranged in a matrix form:

$$\mathbf{O} = \mathbf{W} \mathbf{C} \quad (9)$$

where

$$\mathbf{C} = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_T]^\top \quad (10)$$

$$\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_T]^\top \quad (11)$$

$$\mathbf{w}_t = [\mathbf{w}_t^{(0)}, \mathbf{w}_t^{(1)}, \mathbf{w}_t^{(2)}] \quad (12)$$

$$\begin{aligned} \mathbf{w}_t^{(n)} = & \left[\underset{1\text{-st}}{\mathbf{0}_{M \times M}}, \dots, \underset{(t-L_-^{(n)})\text{-th}}{\mathbf{0}_{M \times M}}, w^{(n)}(-L_-^{(n)}) \mathbf{I}_{M \times M}, \right. \\ & \dots, w^{(n)}(0) \mathbf{I}_{M \times M}, \dots, w^{(n)}(L_+^{(n)}) \mathbf{I}_{M \times M}, \\ & \left. \underset{T\text{-th}}{\mathbf{0}_{M \times M}}, \dots, \underset{T\text{-th}}{\mathbf{0}_{M \times M}} \right]^\top, \quad n = 0, 1, 2 \quad (13) \end{aligned}$$

$L_-^{(0)} = L_+^{(0)} = 0$, and $w^{(0)}(0) = 1$. Under the condition (9), maximizing $P(\mathbf{O}|\mathbf{Q}, \lambda)$ with respect to \mathbf{O} is equivalent to that with respect to \mathbf{C} . By setting

$$\frac{\partial \log P(\mathbf{W} \mathbf{C}|\mathbf{Q}, \lambda)}{\partial \mathbf{C}} = \mathbf{0}, \quad (14)$$

we obtain a set of equations

$$\mathbf{W}^\top \mathbf{U}^{-1} \mathbf{W} \mathbf{C} = \mathbf{W}^\top \mathbf{U}^{-1} \mathbf{M}^\top. \quad (15)$$

For direct solution of (15), we need $O(T^3 M^3)$ operations¹ because $\mathbf{W}^\top \mathbf{U}^{-1} \mathbf{W}$ is a $TM \times TM$ matrix. By utilizing the special structure of $\mathbf{W}^\top \mathbf{U}^{-1} \mathbf{W}$, (15) can be solved by the Cholesky decomposition or the QR decomposition with $O(TM^3 L^2)$ operations², where $L = \max_{n \in \{1, 2\}, s \in \{-, +\}} L_s^{(n)}$. Equation (15) can also be solved by an algorithm derived in [10], [11], which can operate in a time-recursive manner [16].

2.2. Case 2: Maximizing $P(\mathbf{O}, \mathbf{Q}|\lambda)$ with respect to \mathbf{O} and \mathbf{Q}

This problem can be solved by evaluating $\max_{\mathbf{C}} P(\mathbf{O}, \mathbf{Q}|\lambda) = \max_{\mathbf{C}} P(\mathbf{O}|\mathbf{Q}, \lambda) P(\mathbf{Q}|\lambda)$ for all \mathbf{Q} . However, it is impractical because there are too many combinations of \mathbf{Q} . We have developed a fast algorithm for searching for the optimal or sub-optimal state sequence keeping \mathbf{C} optimal in the sense that $P(\mathbf{O}|\mathbf{Q}, \lambda)$ is maximized with respect to \mathbf{C} [10], [11].

To control temporal structure of speech parameter sequence appropriately, HMMs should incorporate state duration densities. The probability $P(\mathbf{O}, \mathbf{Q}|\lambda)$ can be written as $P(\mathbf{O}, \mathbf{Q}|\lambda) = P(\mathbf{O}, \mathbf{i}|\mathbf{q}, \lambda) P(\mathbf{q}|\lambda)$, where $\mathbf{q} = \{q_1, q_2, \dots, q_T\}$, $\mathbf{i} = \{i_1, i_2, \dots, i_T\}$, and the state duration probability $P(\mathbf{q}|\lambda)$ is given by

$$\log P(\mathbf{q}|\lambda) = \sum_{n=1}^N \log p_{q_n}(d_{q_n}) \quad (16)$$

where the total number of states which have been visited during T frames is N , and $p_{q_n}(d_{q_n})$ is the probability of d_{q_n} consecutive observations in state q_n . If we determine the state sequence \mathbf{q} only by $P(\mathbf{q}|\lambda)$ independently of \mathbf{O} , maximizing $P(\mathbf{O}, \mathbf{Q}|\lambda) = P(\mathbf{O}, \mathbf{i}|\mathbf{q}, \lambda) P(\mathbf{q}|\lambda)$ with respect to \mathbf{O} and \mathbf{Q} is equivalent to maximizing $P(\mathbf{O}, \mathbf{i}|\mathbf{q}, \lambda)$ with respect to \mathbf{O} and \mathbf{i} . Furthermore, if we assume that state output probabilities are single-Gaussian, \mathbf{i} is unique. Therefore, the solution is obtained by solving (15) in the same way as the Case 1.

2.3. Case 3: Maximizing $P(\mathbf{O}|\lambda)$ with respect to \mathbf{O}

We derive an algorithm based on an EM algorithm, which find a critical point of the likelihood function $P(\mathbf{O}|\lambda)$. An auxiliary function of current parameter vector sequence \mathbf{O} and new parameter vector sequence \mathbf{O}' is defined by

$$Q(\mathbf{O}, \mathbf{O}') = \sum_{\text{all } \mathbf{Q}} P(\mathbf{O}, \mathbf{Q}|\lambda) \log P(\mathbf{O}', \mathbf{Q}|\lambda). \quad (17)$$

¹When $\mathbf{U}_{q, i}$ is diagonal, it is reduced to $O(T^3 M)$ since each of the M -dimensions can be calculated independently.

²When $\mathbf{U}_{q, i}$ is diagonal, it is reduced to $O(TML^2)$. Furthermore, when $L_-^{(1)} = -1$, $L_+^{(1)} = 0$, and $w^{(2)}(i) \equiv 0$, it is reduced to $O(TM)$ as described in [15].

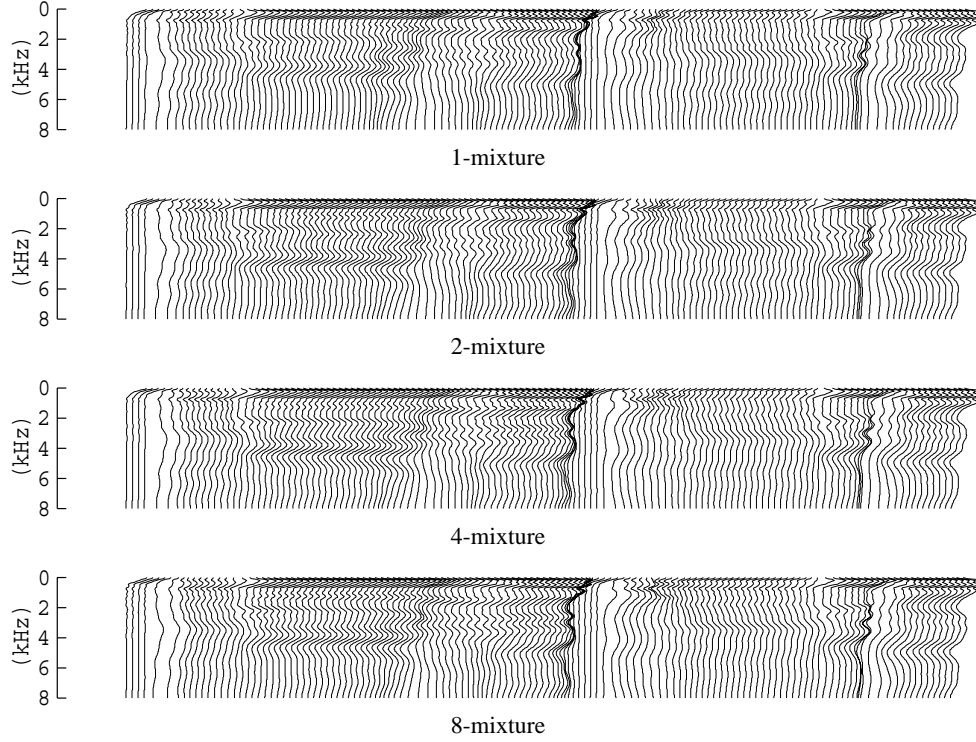


Figure 1: Generated spectra for a sentence fragment “kiNzokuhiroo.”

It can be shown that by substituting \mathbf{O}' which maximizes $Q(\mathbf{O}, \mathbf{O}')$ for \mathbf{O} , the likelihood increases unless \mathbf{O} is a critical point of the likelihood. Equation (17) can be written as

$$Q(\mathbf{O}, \mathbf{O}') = P(\mathbf{O}|\lambda) \left\{ -\frac{1}{2} \mathbf{O}'^\top \overline{\mathbf{U}^{-1}} \mathbf{O}' + \mathbf{O}'^\top \overline{\mathbf{U}^{-1}} \overline{\mathbf{M}} + \overline{\mathbf{K}} \right\} \quad (18)$$

where

$$\overline{\mathbf{U}^{-1}} = \text{diag} \left[\overline{\mathbf{U}_1^{-1}}, \overline{\mathbf{U}_2^{-1}}, \dots, \overline{\mathbf{U}_T^{-1}} \right] \quad (19)$$

$$\overline{\mathbf{U}_t^{-1}} = \sum_{q,i} \gamma_t(q,i) \mathbf{U}_{q,i}^{-1} \quad (20)$$

$$\overline{\mathbf{U}^{-1}} \overline{\mathbf{M}} = \left[\overline{\mathbf{U}_1^{-1}} \overline{\boldsymbol{\mu}_1}^\top, \overline{\mathbf{U}_2^{-1}} \overline{\boldsymbol{\mu}_2}^\top, \dots, \overline{\mathbf{U}_T^{-1}} \overline{\boldsymbol{\mu}_T}^\top \right]^\top \quad (21)$$

$$\overline{\mathbf{U}_t^{-1}} \overline{\boldsymbol{\mu}_t} = \sum_{q,i} \gamma_t(q,i) \mathbf{U}_{q,i}^{-1} \boldsymbol{\mu}_{q,i} \quad (22)$$

and the constant $\overline{\mathbf{K}}$ is independent of \mathbf{O}' . The occupancy probability $\gamma_t(q,i)$ defined by

$$\gamma_t(q,i) = P(q_t = (q,i) | \mathbf{O}, \lambda) \quad (23)$$

can be calculated with the forward-backward inductive procedure. Under the condition $\mathbf{O}' = \mathbf{W} \mathbf{C}'$, \mathbf{C}' which maximizes $Q(\mathbf{O}, \mathbf{O}')$ is given by the following set of equations:

$$\mathbf{W}^\top \overline{\mathbf{U}^{-1}} \mathbf{W} \mathbf{C}' = \mathbf{W}^\top \overline{\mathbf{U}^{-1}} \overline{\mathbf{M}}. \quad (24)$$

The above set of equations has the same form as (15). Accordingly, it can be solved by the algorithm for solving (15).

The whole procedure is summarized as follows:

Step 0. Choose an initial parameter vector sequence \mathbf{C} .

Step 1. Calculate $\gamma_t(q,i)$ with the forward-backward algorithm.

Step 2. Calculate $\overline{\mathbf{U}^{-1}}$ and $\overline{\mathbf{U}^{-1}} \overline{\mathbf{M}}$ by (19)–(22), and solve (24).

Step 3. Set $\mathbf{C} = \mathbf{C}'$. If a certain convergence condition is satisfied, stop; otherwise, goto Step 1.

From the same reason as Case 2, HMMs should incorporate state duration densities. If we determine the state sequence \mathbf{q} only by $P(\mathbf{q}|\lambda)$ independently of \mathbf{O} in a manner similar to the previous section, only the mixture sequence \mathbf{i} is assumed to be unobservable³. Further, we can also assume that \mathbf{Q} is unobservable but phoneme or syllable durations are given.

3. EXAMPLE

We used phonetically balanced 450 sentences from ATR Japanese speech database for training. Speech signal were sampled at 16 kHz and windowed by a 25.6-ms Blackman window with a 5-ms shift, and then mel-cepstral coefficients were obtained by a mel-cepstral analysis technique [18]. Feature vector consists of 25 mel-cepstral coefficients including the zeroth coefficient, their delta and delta-delta coefficients. We used 5-state left-to-right HMMs. Decision-tree based state clustering [19] was applied to the context-dependent phoneme model set, and the resultant HMM set has approximately 900 states. We assume that the state and mixture sequence is unobservable but phoneme durations are given from phoneme duration densities in a manner similar to [20]. It has found that a few iterations are sufficient for convergence of the

³For this problem, an algorithm based on a direct search has also been proposed in [17].

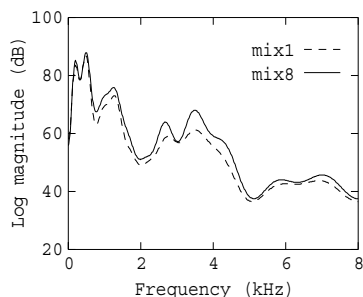


Figure 2: Spectra obtained from 1-mixture HMMs and 8-mixture HMMs.

proposed algorithm. Fig. 1 shows generated spectra for a Japanese sentence fragment “kiNzokuhiroo” taken from a sentence which is not included in the training data. Fig. 2 compares two spectra obtained from single-mixture HMMs and 8-mixture HMMs, respectively, for the same temporal position of the sentence fragment.

It is seen from Fig. 1 and Fig. 2 that with increasing mixtures, the formant structure of the generated spectra get clearer. From informal listening of the synthetic speech, it has been observed that the quality of the synthetic speech is considerably improved by increasing mixtures.

When we use single-mixture HMMs, the formant structure of spectrum corresponding to each mean vector $\mu_{q,i}$ might be vague since $\mu_{q,i}$ is the average of different speech spectra. One can increase the number of decision tree leaf clusters. However, it might result in perceivable discontinuities in synthetic speech since overly large tree will be overspecialized to training data and generalized poorly. We expect that the proposed algorithm can avoid this situation in a simple manner.

4. CONCLUSION

This paper has derived an algorithm for speech parameter generation from HMM whose observation vector consists of spectral parameter vector and its dynamic feature vectors. In the algorithm, we assume that the state sequence (state and mixture sequence for the multi-mixture case) or a part of the state sequence is unobservable (i.e., hidden or latent). The algorithm is based on an EM algorithm. As a result, it iterates the forward-backward algorithm and the parameter generation algorithm for the case where state sequence is given. Experimental results have shown that by using the algorithm, we can reproduce clear formant structure from multi-mixture HMMs as compared with that produced from single-mixture HMMs. Results of subjective evaluation tests of the synthetic speech will be presented in the near future.

5. REFERENCES

- [1] A. Ljolje, J. Hirschberg and J. P. H. van Santen, “Automatic speech segmentation for concatenative inventory selection,” in *Progress in Speech Synthesis*, ed. J. P. H. van Santen, R. W. Sproat, J. P. Olive and J. Hirschberg, Springer-Verlag, New York, 1997.
- [2] R. E. Donovan and P. C. Woodland, “Automatic speech synthesiser parameter estimation using HMMs,” in *Proc. ICASSP*, 1995, pp.640–643.
- [3] X. Huang, A. Acero, H. Hon, Y. Ju, J. Liu, S. Meredith and M. Plumpe, “Recent improvements on Microsoft’s trainable text-to-speech system -Whistler,” in *Proc. ICASSP*, 1997, pp.959–962.
- [4] H. Hon, A. Acero, X. Huang, J. Liu and M. Plumpe, “Automatic generation of synthesis units for trainable text-to-speech synthesis,” in *Proc. ICASSP*, 1998, pp.293–306.
- [5] R. E. Donovan and E. M. Eide, “The IBM Trainable Speech Synthesis System,” in *Proc. ICSLP*, 1998, pp.1703–1706.
- [6] A. Falaschi, M. Giustiniani and M. Verola, “A hidden Markov model approach to speech synthesis,” in *Proc. EUROSPEECH*, 1989, pp.187–190.
- [7] M. Giustiniani and P. Pierucci, “Phonetic ergodic HMM for speech synthesis,” in *Proc. EUROSPEECH*, 1991, pp.349–352.
- [8] T. Masuko, K. Tokuda, T. Kobayashi and S. Imai, “Speech synthesis from HMMs using dynamic features,” in *Proc. ICASSP*, 1996, pp.389–392.
- [9] 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*, 1999, pp.2347–2350.
- [10] K. Tokuda, T. Kobayashi and S. Imai, “Speech parameter generation from HMM using dynamic features,” in *Proc. ICASSP*, 1995, pp.660–663.
- [11] K. Tokuda, T. Masuko, T. Yamada, T. Kobayashi and S. Imai, “An algorithm for speech parameter generation from continuous mixture HMMs with dynamic features,” in *Proc. EUROSPEECH*, 1995, pp.757–760.
- [12] T. Masuko, K. Tokuda, T. Kobayashi and S. Imai, “Voice characteristics conversion for HMM-based speech synthesis system,” in *Proc. ICASSP*, 1997, pp.1611–1614.
- [13] M. Tamura, T. Masuko, K. Tokuda and T. Kobayashi, “Speaker adaptation for HMM-based speech synthesis system using MLLR,” in *Proc. ESCA/COCOSDA Third International Workshop on Speech Synthesis*, 1998, pp.273–276.
- [14] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi and T. Kitamura, “Speaker Interpolation in HMM-Based Speech Synthesis System,” in *Proc. EUROSPEECH*, 1997, pp.2523–2526.
- [15] A. Acero, “Formant analysis and synthesis using hidden Markov models,” in *Proc. EUROSPEECH*, 1999, pp.1047–1050.
- [16] K. Koishida, K. Tokuda, T. Masuko and T. Kobayashi, “Vector quantization of speech spectral parameters using statistics of dynamic features,” in *Proc. ICSP*, 1997, pp.247–252.
- [17] W. Tachiwa and S. Furui, “A study of speech synthesis using HMMs,” in *Proc. Spring Meeting of Acoustical Society of Japan*, Mar. 1999, pp.239–240 (in Japanese).
- [18] T. Fukada, K. Tokuda, T. Kobayashi and S. Imai, “An adaptive algorithm for mel-cepstral analysis of speech,” in *Proc. ICASSP*, 1992, pp.137–140.
- [19] J. J. Odell, “The use of context in large vocabulary speech recognition,” PhD thesis, Cambridge University, 1995.
- [20] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi and T. Kitamura, “Duration modeling for HMM-based speech synthesis,” in *Proc. ICSLP*, 1998, pp.29–32.