# A VITERBI ALGORITHM FOR A TRAJECTORY MODEL DERIVED FROM HMM WITH EXPLICIT RELATIONSHIP BETWEEN STATIC AND DYNAMIC FEATURES

*Heiga Zen, Keiichi Tokuda, Tadashi Kitamura*

Department of Computer Science and Engineering, Nagoya Institute of Technology
Gokiso-cho, Showa-ku, Nagoya 466-8555 Japan
`{zen,tokuda,kitamura}@ics.nitech.ac.jp`

## ABSTRACT

This paper introduces a Viterbi algorithm to obtain a sub-optimal state sequence for trajectory-HMM, which is derived from HMM with explicit relationship between static and dynamic features. The trajectory-HMM can alleviate some limitations of HMM, which are i) constant statistics within HMM state and ii) conditional independence of observations given the state sequence, without increasing the number of model parameters. The proposed algorithm was applied to state-boundary optimization for Viterbi training and $N$-best rescoring. In speaker-dependent continuous speech recognition experiment, trajectory-HMM with the proposed algorithm achieved about 14% error reduction over the standard HMM with the conventional Viterbi algorithm.

## 1. INTRODUCTION

Speech recognition technology has achieved significant progress with the introduction of the hidden Markov model (HMM). Its tractability and efficient implementations are achieved by some assumptions: i) constant statistics within an HMM state, ii) conditional independence assumption of observations given the state sequence. Although these assumptions make HMM practically useful, they are not realistic for modeling sequences of speech spectra, especially in spontaneous speech. To overcome the shortcomings of HMM, a lot of alternative models, called segment models, have been proposed, e.g., [1–9]. Although these models can improve the speech recognition performance, they generally require an increase in the number of model parameters and computational complexity. Alternatively, the use of the dynamic features (delta and delta-delta features) [10] also improves the performance of HMM-based speech recognizers. It is a simple mechanism for capturing time dependencies, however, it has been thought of as an ad hoc, not an essential solution.

In [11], a trajectory model, called trajectory-HMM, was derived by reformulating the standard HMM whose state output vector includes static and dynamic feature parameters. 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 addition, trajectory-HMM provides a computational model for coarticulation and the dynamics of human speech. A Viterbi-type training algorithm was also derived in [11]. However, the lack of an algorithm to obtain the most likely state sequence does not permit iterating Viterbi training procedure and designing a Viterbi decoder.

In this paper, a Viterbi algorithm to obtain a sub-optimal state sequence for trajectory-HMM is proposed. This is the first step for developing a Viterbi decoder based on trajectory-HMM framework. The proposed algorithm is applied to state-boundary optimization for Viterbi training and $N$-best rescoring.

The rest of this paper is organized as follows. Section 2 defines the trajectory-HMM. In Section 3 the time-recursive likelihood computation and a Viterbi algorithm for trajectory-HMM are derived. Results of continuous speech recognition experiment are shown in Section 4. Concluding remarks and future plans are presented in the final section.

## 2. REFORMULATING HMM AS TRAJECTORY-HMM

The output probability of a speech parameter vector sequence $\boldsymbol{o} = \left[\boldsymbol{o}_1^\top, \boldsymbol{o}_2^\top, \ldots, \boldsymbol{o}_T^\top\right]^\top$ for the standard HMM is given by

$$P(\boldsymbol{o} \mid \lambda) = \sum_{\text{all } \boldsymbol{q}} P(\boldsymbol{o} \mid \boldsymbol{q}, \lambda)\, P(\boldsymbol{q} \mid \lambda)\,, \qquad (1)$$

where $\boldsymbol{q} = \{q_1, q_2, \ldots, q_T\}$ is a state sequence. We assume that the speech parameter vector $\boldsymbol{o}_t$ consists of the static feature vector $\boldsymbol{c}_t = [c_t(1), c_t(2), \ldots, c_t(M)]^\top$ (e.g., cepstral coefficients), and dynamic feature vectors $\Delta \boldsymbol{c}_t, \Delta^2 \boldsymbol{c}_t$ (e.g., delta and delta-delta cepstral coefficients), respectively, that is $\boldsymbol{o}_t = [\boldsymbol{c}_t^\top, \Delta \boldsymbol{c}_t^\top, \Delta^2 \boldsymbol{c}_t^\top]^\top$, where the dynamic feature vectors are calculated by

$$\Delta \boldsymbol{c}_t = \sum_{\tau=-L_-^{(1)}}^{L_+^{(1)}} w^{(1)}(\tau) \boldsymbol{c}_{t+\tau}\,, \qquad \Delta^2 \boldsymbol{c}_t = \sum_{\tau=-L_-^{(2)}}^{L_+^{(2)}} w^{(2)}(\tau) \boldsymbol{c}_{t+\tau}\,. \quad (2)$$

Accordingly, when each state output probability distribution is assumed to be a single Gaussian, $P(\boldsymbol{o} \mid \boldsymbol{q}, \lambda)$ is given by

$$P(\boldsymbol{o} \mid \boldsymbol{q}, \lambda) = \prod_{t=1}^{T} \mathcal{N}(\boldsymbol{o}_t \mid \boldsymbol{\mu}_{q_t}, \boldsymbol{\Sigma}_{q_t}) = \mathcal{N}(\boldsymbol{o} \mid \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q)\,, \qquad (3)$$

where $\boldsymbol{\mu}_{q_t}$ and $\boldsymbol{\Sigma}_{q_t}$ are the $3M \times 1$ mean vector and the $3M \times 3M$ covariance matrix, respectively, associated with $q_t$-th state, and

$$\boldsymbol{\mu}_q = \left[\boldsymbol{\mu}_{q_1}^\top, \boldsymbol{\mu}_{q_2}^\top, \ldots, \boldsymbol{\mu}_{q_T}^\top\right]^\top \qquad (4)$$

$$\boldsymbol{\Sigma}_q = \text{diag}\left[\boldsymbol{\Sigma}_{q_1}, \boldsymbol{\Sigma}_{q_2}, \ldots, \boldsymbol{\Sigma}_{q_T}\right]\,. \qquad (5)$$

Conditions (2) can be arranged in a matrix form:

$$\boldsymbol{o} = \boldsymbol{W} \boldsymbol{c}, \qquad (6)$$

where

$$c = [c_1^\top, c_2^\top, \ldots, c_T^\top]^\top \tag{7}$$

$$W = [w_1, w_2, \ldots, w_T]^\top \tag{8}$$

$$w_t = \left[ w_t^{(0)}, w_t^{(1)}, w_t^{(2)} \right] \tag{9}$$

$$w_t^{(n)} = [\, \underbrace{0_{M\times M}}_{\text{1st}}, \ldots, 0_{M\times M}, \underbrace{w^{(n)}(-L_-^{(n)})I_{M\times M}}_{(t-L_-^{(n)})\text{-th}}, \ldots, \underbrace{w^{(n)}(0)I_{M\times M}}_{t\text{-th}},$$
$$\ldots, \underbrace{w^{(n)}(L_+^{(n)})I_{M\times M}}_{(t+L_+^{(n)})\text{-th}}, 0_{M\times M}, \ldots, \underbrace{0_{M\times M}}_{T\text{-th}}\,]^\top, \quad n = 0, 1, 2 \tag{10}$$

$L_-^{(0)} = L_+^{(0)} = 0$, and $w^{(0)}(0) = 1$. Thus, $P(o \mid q, \lambda)$ can be viewed as a function of $c$ as follows:

$$P(Wc \mid q, \lambda) = \mathcal{N}(Wc \mid \mu_q, \Sigma_q)$$
$$= K_q \cdot \mathcal{N}(c \mid \bar{c}_q, P_q), \tag{11}$$

where $\bar{c}_q$, $P_q$ and $K_q$ are given by

$$R_q \bar{c}_q = r_q \tag{12}$$

$$R_q = W^\top \Sigma_q^{-1} W = P_q^{-1} \tag{13}$$

$$r_q = W^\top \Sigma_q^{-1} \mu_q \tag{14}$$

$$K_q = \frac{\sqrt{(2\pi)^{MT}|P_q|}}{\sqrt{(2\pi)^{3MT}|\Sigma_q|}} \cdot \exp\left\{ -\frac{1}{2} \left( \mu_q^\top \Sigma_q^{-1} \mu_q - r_q^\top P_q r_q \right) \right\}. \tag{15}$$

From the above expression (11), trajectory-HMM is defined :

$$P(c \mid \lambda) = \sum_{\text{all } q} P(c \mid q, \lambda)\, P(q \mid \lambda), \tag{16}$$

where

$$P(c \mid q, \lambda) = \mathcal{N}(c \mid \bar{c}_q, P_q). \tag{17}$$

In this paper, Eq. (17) is referred to as "trajectory likelihood". Interestingly, the mean $\bar{c}_q$ is exactly the same as the speech parameter trajectory obtained by the speech parameter generation technique (Case 1 in [12]), that is,

$$\bar{c}_q = \arg\max_c P(o \mid q, \lambda) = \arg\max_c P(Wc \mid q, \lambda). \tag{18}$$

By assuming $\bar{c}_q$ is the mean for the spectral parameter vector sequence $c$, corresponding to an utterance, the standard HMM can naturally be translated into a trajectory model. It is also noted that the spectral parameter vector sequence $c$ is modeled by a mixture of Gaussian distributions whose dimensionality is $TM$, and their covariances $P_q$ are generally full. Because of the above characteristics of the mean vector $\bar{c}_q$ and covariance matrix $P_q$, trajectory-HMM can avoid the limitations of HMM.

## 3. VITERBI ALGORITHM FOR TRAJECTORY-HMM

### 3.1. Time recursive likelihood computation

To compute "trajectory likelihood" directly, high dimensional linear algebra (e.g., matrix inversion, determinant, etc.) is required. To avoid it, a time-recursive algorithm to compute "trajectory likelihood" is described in this section.

From Eq. (11), "trajectory likelihood" is given by

$$P(c \mid q, \lambda) = K_q^{-1} \cdot P(o \mid q, \lambda). \tag{19}$$

In Eq. (15), although $|\Sigma_q|$ and $\mu_q^\top \Sigma_q \mu_q$ can be computed time-recursively, it is difficult to compute $|P_q|$ and $r_q^\top P_q r_q$ recursively because of the temporal full-covariance matrix $P_q$. However, by using the special structure of $P_q$, "trajectory likelihood" can be computed in a time-recursive manner.

First, when $\Delta c_t$ and $\Delta^2 c_t$ are computed as regression coefficients from $\{c_{t-L}, \ldots, c_{t+L}\}$, $R_q$ becomes a $(4L + 1)$-diagonal symmetric positive define matrix. Accordingly, $R_q$ can be decomposed by Cholesky decomposition :

$$R_q = U_q^\top U_q, \tag{20}$$

where $U_q$ is an upper $(2L+1)$-band triangular matrix. From Eq. (20), $|P_q|$ can be rewritten as

$$|P_q| = |R_q|^{-1} = |U_q^\top U_q|^{-1} = |U_q|^{-2} = \prod_{t=1}^{T} \left| U_{q_{t+L}}^{(t,t)} \right|^{-2}, \tag{21}$$

where $q_{t+L} = \{q_1, \ldots, q_{t+L}\}$. Since $U_{q_{t+L}}^{(t,t)}$ depends only on the state sequence from time 1 to $t + L$, $|P_q|$ can be computed time-recursively.

Secondly, from Eqs. (12), (13), and (20), $r_q^\top P_q r_q$ can be rewritten by

$$r_q^\top P_q r_q = r_q^\top P_q^\top R_q P_q r_q = \bar{c}_q^\top U_q^\top U_q \bar{c}_q \tag{22}$$

$$= g_q^\top g_q \qquad \left( g = U_q \bar{c}_q = U_q^{-1} r_q \right) \tag{23}$$

$$= \sum_{t=1}^{T} \left[ g_{q_{t+L}}^{(t)} \right]^\top \cdot g_{q_{t+L}}^{(t)} \tag{24}$$

where $g_q$ is a vector computed from $U_q$ and $r_q$ by forward substitution. Since $g_{q_{t+L}}^{(t)}$ depends only on the state sequence from time 1 to $t + L$, $r_q^\top P_q r_q$ can be also computed time-recursively.

As a result, "trajectory likelihood" (Eq. (17)) can be computed time-recursively as follows :

$$P(c \mid q, \lambda) = \prod_{t=1}^{T} \frac{1}{K_{q_{t+L}}^{(t)}} \cdot P(o_t \mid q_t, \lambda), \tag{25}$$
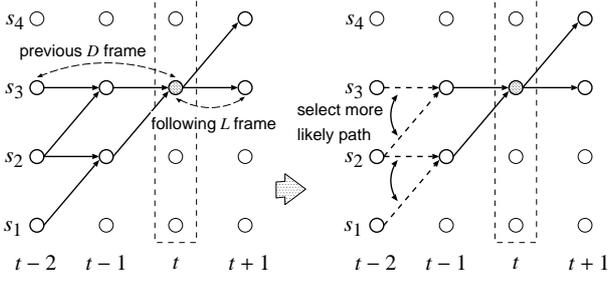
where

$$K_{q_{t+L}}^{(t)} = \frac{\sqrt{(2\pi)^M} \left| U_{q_{t+L}}^{(t,t)} \right|^{-1}}{\sqrt{(2\pi)^{3M} |\Sigma_{q_t}|}} \exp\left\{ -\frac{1}{2} \left( \mu_{q_t}^\top \Sigma_{q_t}^{-1} \mu_{q_t} - \left[ g_{q_{t+L}}^{(t)} \right]^\top g_{q_{t+L}}^{(t)} \right) \right\}. \tag{26}$$

### 3.2. Viterbi algorithm for trajectory-HMMs

In the trajectory-HMM framework, observations and the given state sequence depend on each other because the temporal covariance matrix $P_q$ is full. Accordingly, it is difficult to obtain the most likely state sequence for trajectory-HMM by dynamic programming. In this section, we describe an algorithm to obtain a "sub-optimal" state sequence based on dynamic programming. In this algorithm, the state $q_{t-D}$ is determined at time $t$, according to "trajectory likelihood" for $q_{t+L}$. It can be viewed as a Viterbi algorithm with $D$-frame delayed decision.

The proposed algorithm illustrated in Fig. 1 involves the followings:

**Fig. 1**. An overview for proposed Viterbi algorithm ($D = 2, L = 1$).

**0)** Initialize : $t = 1$ ; $\forall \boldsymbol{q}_{1-D}^{1+L}$

$$\delta_1\left(\boldsymbol{q}_{1-D}^{1+L}\right) = \pi_{q_1} b\left(\boldsymbol{q}_{1+L}\right)$$
$$\psi_1\left(\boldsymbol{q}_{1-D}^{1+L}\right) = 0$$

**1)** Iterate : $t = 2, \ldots, T$ ; $\forall \boldsymbol{q}_{t-D}^{t+L}$

$$\delta_t\left(\boldsymbol{q}_{t-D}^{t+L}\right) = \max_{q_{t-D-1}}\left[\delta_{t-1}\left(\boldsymbol{q}_{t-D-1}^{t+L-1}\right) a_{q_{t-1}, q_t}\right] b\left(\boldsymbol{q}_{t+L}\right)$$
$$\psi_t\left(\boldsymbol{q}_{t-D}^{t+L}\right) = \arg\max_{q_{t-D-1}}\left[\delta_{t-1}\left(\boldsymbol{q}_{t-D-1}^{t+L-1}\right) a_{q_{t-1}, q_t}\right]$$

**2)** Finish :

$$\hat{P} = \max_{\boldsymbol{q}_{T-D}^{T+L}}\left[\delta_T\left(\boldsymbol{q}_{T-D}^{T+L}\right)\right]$$
$$\hat{\boldsymbol{q}}_{T-D}^{T+L} = \arg\max_{\boldsymbol{q}_{T-D}^{T+L}}\left[\delta_T\left(\boldsymbol{q}_{T-D}^{T+L}\right)\right]$$
$$= \{\hat{q}_{T-D}, \ldots, \hat{q}_{T+L}\}$$

**3)** Back track : $t = T, \ldots, D + 2$

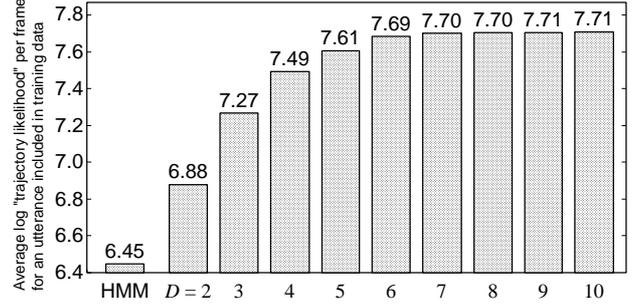$$\hat{q}_{t-D-1} = \psi_t\left(\hat{\boldsymbol{q}}_{t-D}^{t+L}\right)$$
$$\hat{\boldsymbol{q}}_{t-D-1}^{t+L-1} = \{\hat{q}_{t-D-1}, \ldots, \hat{q}_{t+L-1}\}$$

where $\pi_{q_1}$ is an initial state probability of state $q_1$, $a_{q_{t-1}, q_t}$ is a state transition probability from $q_{t-1}$ to $q_t$, $\boldsymbol{q}_{t-D}^{t+L} = \{q_{t-D}, \ldots, q_{t+L}\}$, and $b\left(\boldsymbol{q}_{t+L}\right)$ is given by

$$b(\boldsymbol{q}_{t+L}) = K_{q_{t+L}}^{(t)} \cdot P\left(\boldsymbol{o}_t \mid q_t, \lambda\right) . \qquad (27)$$

When $D$ is equal to $T$, the most likely state sequence can be obtained by the above algorithm. Although the proposed algorithm with larger $D$ can obtain more likely state sequence, it requires huge amount of computational cost. Hence, $D$ have to be set to a proper value balancing its performance and computational complexity.

It is generally considered that the coarticulation affect neighboring frames within 100–200 ms. This indicates that optimal state sequence is probably obtained by the proposed algorithm when $D$ is set to larger than 10 for 10-ms frame shift.



**Fig. 2**. Average log "trajectory likelihood" per frame for the standard HMM segmentation and state sequences obtained by the proposed algorithm.

## 4. EXPERIMENT

### 4.1. Experimental conditions

We used phonetically balanced 503 sentences uttered by a male speaker MHT from the ATR Japanese speech database b-set. The 450 sentences were used for training monophone standard HMMs and trajectory-HMMs. The remaining 53 sentences were used for testing. Speech signals were sampled at 16 kHz and windowed by a 25.6-ms Blackman window with a 10-ms shift, and then mel-cepstral coefficients were obtained by a mel-cepstral analysis technique . Feature vector consists of 19 mel-cepstral coefficients including the zeroth coefficient, their delta and delta-delta coefficients. We used 3-state left-to-right with no skip HMM structure for modeling 36 Japanese phonemes. Each state output probability distribution is represented by a single Gaussian with diagonal covariance matrix.

### 4.2. Experimental results

First, we evaluated the performance of the proposed Viterbi algorithm in the likelihood of the obtained state sequence. Figure 2 shows the average log "trajectory likelihood" per frame of the obtained state sequence for an utterance included in the training data. In this figure, "HMM" means that the standard HMM segmentation was used for state sequence $\boldsymbol{q}$, and "$D = 2, \ldots, 10$" mean that the state sequences were obtained by the proposed algorithm. In this experiment, standard HMMs were used for acoustic models. It shows that the proposed algorithm could obtain more likely state sequence than the standard HMM segmentation. Furthermore, as $D$ increased, "trajectory likelihood" was converged. When $D$ was greater than 5, further likelihood improvement was not achieved. It suggests that the proposed algorithm with $D = 5$ can obtain approximately optimal state sequence.

Secondly, the proposed algorithm was applied to state-boundary optimization in Viterbi training procedure. After obtaining the state sequences by the proposed algorithm, model parameters were updated by formulas derived in [11] according to the obtained state sequences. This procedure was iterated several times. Table 1 shows the average log "trajectory likelihood" per frame for whole training data. In Table 1 and 2, "Iteration 0" means that initial models were used. In this experiment, the average log likelihood of the standard HMM segmentation for whole training data was 6.873. These results show that the proposed algorithm could obtain

**Table 1**. Average log "trajectory likelihood" per frame of the trained trajectory-HMMs for whole training data

|       | Iteration | | | | |
|-------|-------|-------|-------|-------|-------|
|       | 0     | 1     | 2     | 3     | 4     |
| $D = 2$ | 8.112 | 14.76 | 15.01 | 15.08 | 15.13 |
| $D = 3$ | 8.501 | 14.93 | 15.17 | 15.28 | 15.32 |
| $D = 4$ | 8.642 | 15.01 | 15.26 | 15.36 | 15.40 |
| $D = 5$ | 8.699 | 15.04 | 15.30 | 15.39 | 15.42 |

**Table 2**. Phoneme error rate (%) for iteratively trained trajectory-HMMs (the number of $N$-best candidates was 20).

|       | Iteration | | | | | |
|-------|------|------|------|------|------|------|
|       | 0    | 1    | 2    | 3    | 4    | 5    |
| $D = 2$ | 20.1 | 18.8 | 18.5 | 18.8 | 18.8 | 18.4 |
| $D = 3$ | 19.5 | 18.6 | 18.7 | 18.6 | 18.8 | 18.8 |
| $D = 4$ | 19.7 | **18.1** | 18.2 | 18.7 | 18.8 | 18.5 |
| $D = 5$ | 19.5 | **18.1** | 18.5 | 18.5 | 18.6 | 18.6 |

**Table 3**. Phoneme error rates (%) for the several number of $N$-best candidates.

|       | Number of candidates | | | |
|-------|------|------|------|------|
|       | 20   | 30   | 100  | 200  |
| $D = 2$ | 18.8 | 18.7 | 18.2 | 18.6 |
| $D = 3$ | 18.6 | 18.4 | 17.9 | 17.8 |
| $D = 4$ | 18.1 | 17.8 | 17.6 | 17.5 |
| $D = 5$ | 18.1 | 17.8 | 17.5 | **17.2** |

more likely state sequence than the standard HMM segmentation. In addition, iterative training improved the model likelihood for training data.

Thirdly, trained trajectory-HMMs were evaluated in speaker-dependent continuous phoneme recognition experiments. First, an $N$-best list was generated for each test utterance by using the HTK Viterbi decoder with the phonotactic constraints of phoneme sequences in Japanese. For $N$-best list generation, initial models were used. Then, each candidate was re-segmented by the proposed algorithm and rescored with the "trajectory likelihood". The proposed algorithm with the same $D$ were used for both training and re-segmentation. The phoneme error rate for the baseline system based on the standard HMM was 20.1%.

Table 2 shows the phoneme error rates for iteratively trained trajectory-HMMs. Recognition results were shown in phoneme error rates. When trajectory-HMMs trained by single Viterbi training iteration with the proposed algorithm setting $D$ to 4 or 5 were used, about 10% relative error reduction was achieved over the standard HMM. On the other hand, recognition performance did not improve by iterative training. This suggests that iterative training causes over-fitting to training data.

Table 3 shows the phoneme error rates for the several numbers of $N$-best candidates. In this experiment, trajectory-HMMs trained by Viterbi training with single iteration were used. As the number of $N$-best candidates increased, the recognition performance improved significantly. When 200-best candidates were used, about 14% relative error reduction over the standard HMM was achieved. It indicates that if we search wider space by implementing a Viterbi decoder based on the proposed algorithm, it is expected that larger improvement can be achieved.

## 5. CONCLUSION

This paper describes a Viterbi algorithm to obtain a sub-optimal state sequence for trajectory-HMM. It can be considered as a Viterbi algorithm with $D$-frame delayed decision.

It was applied to state-boundary optimization for Viterbi training and $N$-best rescoring. In speaker-dependent continuous speech recognition experiment using rescoring scheme, about 14% relative error reduction over the standard HMM with the conventional Viterbi algorithm was achieved.

Future work includes the implementation of Viterbi decoder based on proposed algorithm and derivation of a Baum-Welch-type training algorithm.

## 7. REFERENCES

[1] M.J.F. Gales and S.J. Young, "Segmental hidden Markov models," *Proc of Eurospeech'93*, pp.1579–1582, 1993.

[2] H. Gish and K. Ng, "Parametric trajectory models for speech recognition," *Proc of ICSLP'96*, vol.I, pp.466–469, 1996.

[3] W.J. Holmes and M.J. Russell, "Speech recognition using a linear dynamic segmental HMM," *Proc of ICASSP'95*, pp.1611–1614, 1995.

[4] M. Ostendorf, V. Digalakis, and O.A. Kimball, "From HMMs to segment models," *IEEE Trans. on Speech & Audio Process.*, vol.4, no.5, pp.360–378, 1996.

[5] L. Deng, M. Aksmanovic, X. Sun, and J. Wu, "Speech recognition using hidden Markov models with polynomial regression functions as nonstationary states," *IEEE Trans. on Speech & Audio Process.*, vol.2, no.4, pp.507–520, 1994.

[6] Y. Gong and J.P. Haton, "Stochastic trajectory modeling for speech recognition," *Proc. of ICASSP'94*, vol.I, pp57–60, 1994.

[7] K.K. Paliwal, "Use of temporal correlation between successive frames in hidden Markov model based Speech recognizer," *Proc. of ICASSP'93*, pp.215-218, 1993.

[8] M. Ostendorf and S. Roukos, "A stochastic segment model for phoneme-based continuous speech recognition," *IEEE Trans. Acoust., Speech & Signal Process.*, vol.37, no.12, pp.1857–1869, 1989.

[9] C.J. Wellekens, "Explicit correlation in hidden Markov model for speech recognition," *Proc. of ICASSP'87*, pp.383–386, 1987.

[10] S. Furui, "Speaker independent isolated word recognition using dynamic features of speech spectrum," *IEEE Trans. Acoust., Speech, & Signal Process.*, vol.34, pp.52–59, 1986.

[11] K. Tokuda, H. Zen, and T. Kitamura, "Trajectory modeling based on HMMs with the explicit relationship between static and dynamic features," *Proc of Eurospeech2003*, pp.865–868, 2003.

[12] K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, and T. Kitamura, "Speech parameter generation algorithms for HMM-based speech synthesis," *Proc. of ICASSP2000*, vol.3, pp.1315–1318, 2000.