

# MAXIMUM LIKELIHOOD BASED HMM STATE FILTERING APPROACH TO MODEL ADAPTATION FOR LONG REVERBERATION

Chandra Kant Raut Takuya Nishimoto Shigeki Sagayama

Graduate School of Information Science and Technology  
The University of Tokyo  
7-3-1 Hongo, Bunkyo-ku, Tokyo, 113-8656 Japan  
{raut, nishi, sagayama}@hil.t.u-tokyo.ac.jp

## ABSTRACT

In environment with considerably long reverberation time, each frame of speech is affected by reflected energy components from the preceding frames. Therefore to adapt parameters of a state of HMM, it becomes necessary to consider these frames, and compute their contributions to current state. However, these *clean* speech frames preceding to a state of HMM are not known during adaptation of the models. In this paper, we propose to use preceding states as units of preceding speech, and estimate their contributions to current state in maximum likelihood fashion. The experimental results on an isolated word recognition task showed significant improvement in performance of speech recognition system for reverberant speech, compared to other methods.

## 1. INTRODUCTION

Automatic speech recognition (ASR) systems, though usually trained with clean speech, have to operate under real-life environment for any practical purpose. One of the key factors that severely degrade the performance of speech recognition system under such practical environment is convolutional distortion caused by room reverberation or channel characteristics.

Several methods have been developed to deal with such distortion, ranging from front-end methods like inverse filtering, channel normalization and microphone-phone array based techniques to different model based approaches. Channel normalization techniques like Cepstrum Mean Subtraction (CMS) [1] and RASTA [2] have been proved effective to improve the performance of the system. Other variants of normalization technique like SNR-dependent cepstrum normalization and codeword-dependent cepstrum normalization [3] have been also proposed. Model-based approaches [4, 5] have been also applied to compensate the HMMs for channel distortion or reverberant condition. Though these methods have been proved to improve the performance of ASRs, most of them cannot perform well when reverberation time is much longer than analysis window-length and additive noise is also present. However, reverberation time longer than 100 ms is not uncommon [6] in real-life environment, e.g., in office rooms.

In our previous works [7, 8], we proposed a state splitting approach to deal with such long reverberation, which

estimates preceding frames for a state of HMM and finds the compensated parameters by convolving estimated speech frames with channel parameters. In those works, channel characteristic was either assumed to be explicitly given, or estimated from adaptation data using minimum mean square error (MMSE).

This work also deals with speech distorted by long reverberation. In this work, we model the reflected component from preceding speech units in terms of preceding *states*, and estimate their contributions in maximum-likelihood manner from adaptation data. The method differs from our previous work in that it eliminates state-splitting as well as the need to have stereo-data (clean speech and its reverberated counterpart) for estimating channel parameters, making it more practical method for compensating HMMs under real-life environment.

## 2. EFFECT OF LONG REVERBERATION

Reverberation of a room is modeled by passing clean speech signal through a filter with impulse response of the propagation channel as transfer function, such that reverberant speech is given by

$$o[m] = h[m] * s[m] \quad (1)$$

where  $s[m]$  is clean speech,  $h[m]$  is impulse response of the room,  $o[m]$  is reverberant speech,  $m$  is sample number and  $*$  represents convolution in time domain.

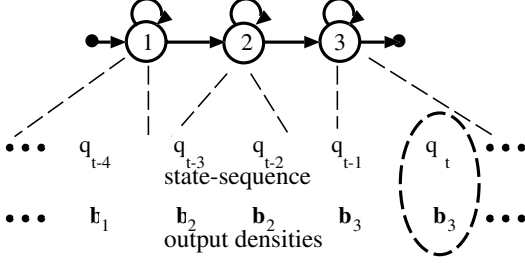
Taking the short-time Fourier transform (STFT) of Eq. 1 gives

$$O(w_k, t) \approx H(w_k, t)S(w_k, t) \quad (2)$$

where  $t$  is frame number and  $w_k$  represents discrete frequency. Parameters  $S(w_k, t)$ ,  $O(w_k, t)$  and  $H(w_k, t)$  are STFTs of clean speech  $s[m]$ , reverberant speech  $o[m]$  and impulse response of channel  $h[m]$ , respectively. However, such a relationship is valid only when length of  $h[m]$  is shorter than window-size used for STFT.

The effect of reverberation with reverberation time ( $T_{60}$ ) longer than the analysis window size, on short-time Fourier transform of speech is usually approximated by

$$O(w_k, t) \approx H(w_k, t) * S(w_k, t). \quad (3)$$



**Fig. 1.** HMM Adaptation: To compensate state output  $b_j$  at state  $q_t = j$  for long reverberation, by accounting energy components from previous frames, knowledge of clean observations at time  $t - 1$ ,  $t - 2$  and so on is necessary.

In other words, the effect of long reverberation is no more multiplicative in linear spectral domain, rather it is convolutional.

### 3. ACCOUNTING CONVOLUTIONAL DISTORTION

Eq. 3 shows that the spectral parameters of reverberant speech at frame  $t$  do not depend only upon this frame, but also upon the preceding frames at  $t - 1$ ,  $t - 2$  and so on. Therefore, to adapt the output distribution  $b_j$  at state  $q_t = j$  of given HMM [Fig. 1], frames occurred at time  $t - 1$ ,  $t - 2$  and so on should be considered. However, with such a conventional HMM used in most of speech recognition systems, nothing can be inferred deterministically about the observations, and even the state sequence preceding to a given state cannot be known (as adaptation under consideration is *not* decoding-time adaptation and therefore does not use observations, and in such case even most likely state sequence preceding to a state cannot be estimated).

In [7, 8], we proposed a state splitting approach to estimate preceding frames for a given state of HMM by using state's duration information to estimate state-sequence and using composite mean of the output distributions of those states as preceding observation sequence. Such estimation of preceding frames becomes necessary, when the channel parameters are estimated for the frame-level convolution. Assuming that *average occupation of states does not vary significantly, or if each state is usually occupied only once*, we approximate the convolutional distortion by filtering of *states* (convolution over states) directly, with some optimal channel parameters, without necessitating to work at frame-level and estimate preceding *frames*. Such convolution over states, in linear spectral domain, is represented as

$$\hat{O}(j) = \alpha_0 S(j) + \alpha_1 S(j-1) + \alpha_2 S(j-2) + \dots + \alpha_{N-1} S(j-N+1) \quad (4)$$

where  $S$  and  $\alpha_i$  are clean speech state output distribution and state-filter coefficients, respectively (separate equations in terms of mean and covariance matrix defined later in Eq. 11 and 12);  $k$  represents the dimension of the parameter, and  $j$  represents HMM states of model (please see Fig. 2

for interpretation of  $j$ ). Left contexts of models can be used to account the effect of preceding models; in their absence, only the available preceding states of current model can be used.

With this approach, need for state-splitting and estimation of preceding frames using state's duration information is eliminated; however, optimal  $\alpha_i$  should be estimated for such convolution over states. The next section describes its estimation by maximum likelihood approach.

### 4. MAXIMUM-LIKELIHOOD ESTIMATION OF STATE-FILTER COEFFICIENTS

Reverberant speech model  $\lambda_O$  is composed by using clean speech model  $\lambda_S$  and state-level channel parameters  $\mathbf{A} = \{\alpha_0, \dots, \alpha_i, \dots, \alpha_{N-1}\}$ . Parameter  $\alpha_{ik}$  ( $k$ : dimension of speech parameter) is estimated by maximizing Viterbi-likelihood score  $P(\mathbf{O}, \mathbf{q} | \mathbf{A}, \lambda_S)$  or  $P(\mathbf{O}, \mathbf{q} | \lambda_O)$  of training observation  $\mathbf{O} = \{\mathbf{o}_1, \dots, \mathbf{o}_T\}$  over most likelihood state sequence  $\mathbf{q} = \{q_1, \dots, q_T\}$  given by Viterbi algorithm, as

$$\hat{\alpha}_{ik} = \arg \max_{\alpha_{ik}} P(\mathbf{O} | \alpha_0, \dots, \alpha_{N-1}, \lambda_S). \quad (5)$$

Maximization of  $P(\mathbf{O}, \mathbf{q} | \mathbf{A}, \lambda_S)$  is done in iterative manner by steepest-descent method, by defining new estimate of  $\alpha_{ik}$  at  $p$ th iteration as

$$\alpha_{ik}(p) = \alpha_{ik}(p-1) + \epsilon \frac{\partial \log (P(\mathbf{O} | \alpha_0, \dots, \alpha_{N-1}, \lambda_S))}{\partial \alpha_{ik}} \quad (6)$$

where  $\epsilon$  is gradient scaling factor.

Computation of likelihood  $P(\mathbf{O}, \mathbf{q} | \mathbf{A}, \lambda_S)$  (and its maximization) is done using cepstrum domain parameters, whereas composition of reverberant speech model  $\lambda_O$  by using  $\lambda_S$ ,  $\mathbf{A}$  and additive noise (when considered) is done in linear-spectral domain, under PMC framework [4]. Therefore, such estimation involves conversion of parameters between these domains at each iteration. The parameters in cepstrum, log and linear spectral domains are specified by sub-script cep, lg and lin, respectively.

Transformation of models from cepstrum-domain to log spectral domain is done as

$$\boldsymbol{\mu}^{S_{lg}} = \mathbf{C}^{-1} \boldsymbol{\mu}^{S_{cep}} \quad (7)$$

$$\boldsymbol{\Sigma}^{S_{lg}} = \mathbf{C}^{-1} \boldsymbol{\Sigma}^{S_{cep}} (\mathbf{C}^{-1})^T \quad (8)$$

where  $\mathbf{C}$  is the discrete Cosine transform (DCT) matrix. These parameters in log spectral domain are further transformed to linear spectral domain, by using

$$\mu_k^{S_{lin}} = \exp \left( \mu_k^{S_{lg}} + \frac{\Sigma_{kk}^{S_{lg}}}{2} \right) \quad (9)$$

$$\Sigma_{kl}^{S_{lin}} = \mu_k^{S_{lin}} \mu_l^{S_{lin}} \left( \exp(\Sigma_{kl}^{S_{lg}}) - 1 \right) \quad (10)$$

where  $k$  and  $l$  are parameter indices.

In linear spectral domain, model for reverberant speech is composed by using clean speech model and estimated  $\alpha_{ik}$  as

$$\begin{aligned} \mu_k^{O_{\text{lin}}}(j) &= \alpha_{0k} \mu_k^{S_{\text{lin}}}(j) \\ &+ \alpha_{1k} \bar{\mu}_k^{S_{\text{lin}}}(j-1) + \alpha_{2k} \bar{\mu}_k^{S_{\text{lin}}}(j-2) \\ &+ \dots + \alpha_{N-1,k} \bar{\mu}_k^{S_{\text{lin}}}(j-N+1) \end{aligned} \quad (11)$$

$$\Sigma_{kl}^{O_{\text{lin}}}(j) = \alpha_{0k} \alpha_{0l} \Sigma_{kl}^{S_{\text{lin}}}(j) \quad (12)$$

From the preceding states, only composite mean (distinguished by overbar) from single component distribution corresponding to Gaussian mixture model of output distributions are used. Such single component composite distribution from M-mixture GMM can be obtained, as

$$\bar{\boldsymbol{\mu}}_{\text{cep}} = \sum_{m=1}^M c_m \boldsymbol{\mu}_{m,\text{cep}} \quad (13)$$

$$\bar{\boldsymbol{\Sigma}}_{\text{cep}} = \sum_{m=1}^M c_m (\boldsymbol{\Sigma}_{m,\text{cep}} + \boldsymbol{\mu}_{m,\text{cep}} \boldsymbol{\mu}_{m,\text{cep}}^T) - \bar{\boldsymbol{\mu}}_{\text{cep}} \bar{\boldsymbol{\mu}}_{\text{cep}}^T \quad (14)$$

where  $m$  represents mixture component, and  $c_m$  is mixture weight. Furthermore, effect of preceding states on covariance matrix is ignored.

Once the models are adapted in linear spectral domain, they are transformed back to log spectral domain by using

$$\mu_k^{O_{\text{lg}}} = \log(\mu_k^{O_{\text{lin}}}) - \frac{1}{2} \log\left(\frac{\Sigma_{kk}^{O_{\text{lin}}}}{\mu_k^{O_{\text{lin}}2}} + 1\right) \quad (15)$$

$$\Sigma_{kl}^{O_{\text{lg}}} = \log\left(\frac{\Sigma_{kl}^{O_{\text{lin}}}}{\mu_k^{O_{\text{lin}}} \mu_l^{O_{\text{lin}}}} + 1\right), \quad (16)$$

and to cepstrum domain by using

$$\boldsymbol{\mu}^{O_{\text{cep}}} = \mathbf{C} \boldsymbol{\mu}^{O_{\text{lg}}} \quad (17)$$

$$\boldsymbol{\Sigma}^{O_{\text{cep}}} = \mathbf{C} \boldsymbol{\Sigma}^{O_{\text{lg}}} \mathbf{C}^T. \quad (18)$$

Such formulations for transformation of parameters from one-domain to another and composition of models are used while estimating  $\alpha_{ik}$  as well. We use similar approach as in [5] to maximize likelihood and estimate filter coefficients. As under large mixture GMMs, estimation of  $\alpha_{ik}$  becomes complex, they can be first reduced to single-component distribution using Eqs.13 and 14 and used while estimating  $\alpha_{ik}$ . The new estimate for  $\alpha_{ik}$ , under single-mixture case, is given by

$$\begin{aligned} \alpha_{ik}(p) &= \alpha_{ik}(p-1) \\ &+ \epsilon \frac{\partial}{\partial \alpha_{ik}} \sum_{\forall t} \left\{ -\frac{1}{2} \log((2\pi)^D | \boldsymbol{\Sigma}_t^{O_{\text{cep}}} |) \right. \\ &\left. - \frac{1}{2} (\mathbf{O}_t - \boldsymbol{\mu}_t^{O_{\text{cep}}})^T \boldsymbol{\Sigma}_t^{O_{\text{cep}-1}} (\mathbf{O}_t - \boldsymbol{\mu}_t^{O_{\text{cep}}}) \right\} \end{aligned} \quad (19)$$

where  $\{(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \dots, (\boldsymbol{\mu}_T, \boldsymbol{\Sigma}_T)\}$  corresponds to output distributions of most likely state-sequence decoded by Viterbi

algorithm. Ignoring the change in covariance w.r.t.  $\alpha_{ik}$  leads to

$$\begin{aligned} \alpha_{ik}(p) &= \alpha_{ik}(p-1) \\ &+ \epsilon \sum_{\forall t} \left( \frac{1}{2} \frac{\partial \boldsymbol{\mu}_t^{O_{\text{cep}T}}}{\partial \alpha_{ik}} \boldsymbol{\Sigma}_t^{O_{\text{cep}-1}} (\mathbf{O}_t - \boldsymbol{\mu}_t^{O_{\text{cep}}}) \right. \\ &+ \left. \frac{1}{2} (\mathbf{O}_t - \boldsymbol{\mu}_t^{O_{\text{cep}}})^T \boldsymbol{\Sigma}_t^{O_{\text{cep}-1}} \frac{\partial \boldsymbol{\mu}_t^{O_{\text{cep}}}}{\partial \alpha_{ik}} \right) \quad (20) \\ &= \alpha_{ik}(p-1) \\ &+ \epsilon \sum_{\forall t} \left( \frac{1}{2} (\mathbf{C} \frac{\partial \boldsymbol{\mu}_t^{O_{\text{lg}}}}{\partial \alpha_{ik}})^T \boldsymbol{\Sigma}_t^{O_{\text{cep}-1}} (\mathbf{O}_t - \boldsymbol{\mu}_t^{O_{\text{cep}}}) \right. \\ &+ \left. \frac{1}{2} (\mathbf{O}_t - \boldsymbol{\mu}_t^{O_{\text{cep}}})^T \boldsymbol{\Sigma}_t^{O_{\text{cep}-1}} \mathbf{C} \frac{\partial \boldsymbol{\mu}_t^{O_{\text{lg}}}}{\partial \alpha_{ik}} \right). \end{aligned} \quad (21)$$

The term  $\partial \boldsymbol{\mu}_t^{O_{\text{lg}}}/\partial \alpha_{ik}$  (each  $k$ th component represented as  $\partial \mu_k^{O_{\text{lg}}}(j)/\partial \alpha_{ik}$ , where  $j$  is the aligned state to frame  $t$  of training observation) can be obtained by taking derivative of Eq. 15 as

$$\begin{aligned} \frac{\partial \mu_k^{O_{\text{lg}}}(j)}{\partial \alpha_{ik}} &= \frac{\mu_k^{S_{\text{lin}}}(j-i)}{\mu_k^{O_{\text{lin}}}(j)} \\ &+ \frac{\Sigma_{kk}^{O_{\text{lin}}}(j) \mu^{S_{\text{lin}}}(j-i)}{\mu_k^{O_{\text{lin}}}(j) \Sigma_{kk}^{O_{\text{lin}}}(j) + (\mu_k^{O_{\text{lin}}}(j))^3}. \end{aligned} \quad (22)$$

While also considering additive noise, mean and covariance matrix terms for it can be included in Eqs. 11 and 12, and can be estimated together, or if already estimated (e.g. using signal during non-speech activity), they can be used during estimation of  $\alpha_{ik}$ . When likelihood score converges for the training observation, the estimated values of  $\alpha_{ik}$  are used to transform all the clean models to models for reverberant speech. The procedure is also depicted in Fig. 2.

## 5. EVALUATION

The proposed method was evaluated on a speaker-dependent isolated word recognition task. Clean speech HMMs were trained with 2620 words of the same speaker taken from ATR speech database A-Set. Clean speech HMMs comprised of 425 context-dependent biphone models with left-context, each with three emitting states single mixture Gaussian model. The speech signal was single channel with sampling frequency of 16 kHz. The speech signal was analyzed with Hamming window of 25 ms window-size and frame shift of 10 ms into 13-dimensional MFCC feature vectors including 0th-order coefficient, using 24 mel filter-banks. The test set consisted of 655 words of the same speaker taken exclusively from the ATR speech database A-set, and HTK 3.1 was used as decoder.

For evaluation, reverberant speech was simulated by a linear convolution of clean speech and impulse responses of the environment (viz. E1B and OFC) taken from RWCP Sound Scene Database in Real Environment. The performance of the speech recognition system decreased with clean

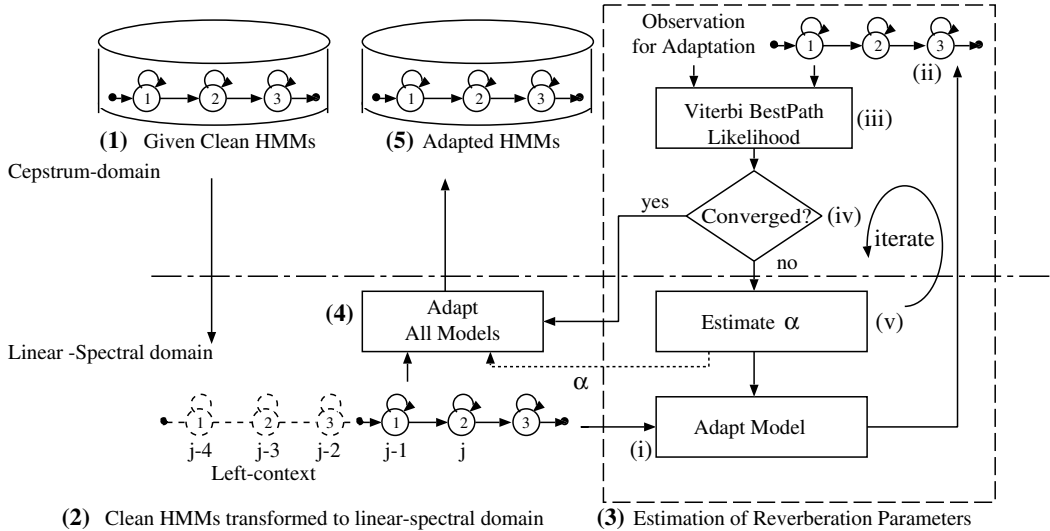


Fig. 2. Adaptation of model parameters

Table 1. Experimental Results (Word Recognition Rate %)

Data	$T_{60}$	Clean	CMS	SF(N=4)
Clean	—	97.9	—	—
E1B	310 ms	67.6	77.3	83.2
OFC	780 ms	44.8	47.5	72.5

model as listed under “Clean” in Table 1 for the reverberant speech. The recognition performance of reverberant speech was evaluated with CMS as well. For this purpose, CMS was performed on the same training set data, and the model was retrained with it. CMS was applied to test set also, and performance was evaluated with the retrained model.

To evaluate the proposed method, ten words of reverberant speech was used as adaptation data to estimate  $\alpha_{ik}$ , with filter-order of  $N = 4$ , and states of left-context were considered for the convolution.

Experimental results as listed under Table 1 show better performance of proposed state-filtering (SF) approach compared to clean model and CMS, which demonstrates its effectiveness for recognizing reverberant speech. With longer reverberation time, the improvement with state-filtering approach is more pronounced compared to other methods.

## 6. CONCLUSION

In this paper, we proposed a method for model adaptation based on state-filtering, for reverberant speech with considerably long reverberation time. The filter coefficients are estimated by using maximum-likelihood approach, using small amount of training data. The experimental result shows the effectiveness of the method for improving performance of the system under reverberant condition.

Future work includes evaluation of the method on large vocabulary continuous speech recognition task and with Gaussian mixture models. Its application for adapting models for moving-speaker case will be also investigated.

## 7. REFERENCES

- [1] B. S. Atal, “Effectiveness of linear prediction characteristics of the speech wave for automatic speaker identification and verification,” *J. Acoust. Soc. Am.*, vol. 55, pp. 1304–1312, 1974.
- [2] B. E. Kingsbury and N. Morgan, “Recognizing reverberant speech with RASTA-PLP,” in *Proc. ICASSP*, Munich, Germany, 1997, pp. 1259–1262.
- [3] A. Acero and R. M. Stern, “Environmental robustness in automatic speech recognition,” *Proc. ICASSP*, pp. 849–852, 1990.
- [4] M. J. F. Gales and S. J. Young, “Robust speech recognition in additive and convolutional noise using parallel model combination,” *Computer Speech and Language*, vol. 9, pp. 289–307, 1995.
- [5] Y. Minami and S. Furui, “A maximum likelihood procedure for a universal adaptation method based on HMM composition,” in *Proc. ICASSP*, 1995, pp. 129–132.
- [6] X. Huang, A. Acero, and H.-W. Hon, *Spoken Language Processing: A Guide to Theory, Algorithm, and System Development*, Prentice Hall PTR, New Jersey, 1st edition, 2001.
- [7] C. K. Raut, T. Nishimoto, and S. Sagayama, “Model adaptation by state splitting of HMM for long reverberation,” in *Proc. Interspeech*, Sep. 2005.
- [8] C. K. Raut, T. Nishimoto, and S. Sagayama, “Model convolution by state splitting of HMM for robust speech recognition in presence of convolutional noise,” in *Proc. The Acoustical Society of Japan*, Mar. 2005, vol. 3-5-5, pp. 85–86.