# STEREO-BASED STOCHASTIC NOISE COMPENSATION BASED ON TRAJECTORY GMMS

Heiga Zen, Yoshihiko Nankaku, Keiichi Tokuda

Department of Computer Science, Nagoya Institute of Technology, Nagoya, 466-8555 JAPAN {zen,nankaku,tokuda}@sp.nitech.ac.jp

#### ABSTRACT

This paper proposes a novel stereo-based stochastic noise compensation technique based on trajectory GMMs. Although the GMM-based noise compensation techniques such as SPLICE work effective, their performance sometimes degrades due to the inappropriate dynamic characteristics caused by the frame-by-frame mapping. While the use of dynamic feature constraints on the mapping stage can alleviate this problem, it also introduces an inconsistency between training and mapping. The recently proposed trajectory GMM-based feature mapping technique can solve this inconsistency while keeping the benefits of the use of dynamic features, and offers an entire sequence-level transformation rather than the frame-byframe mapping. Results from a noise compensation experiment on the AURORA-2 task show that the proposed trajectory GMM-based noise compensation technique outperforms the conventional ones.

Index Terms— speech recognition, GMM-based mapping, trajectory GMM, dynamic features, noise compensation

#### 1. INTRODUCTION

During these last years, stereo-based continuous stochastic feature mapping techniques based on Gaussian mixture models (GMMs) have been applied to noise compensation [1,2]. This mapping technique consists of three major stages:

- 1. Joint probability density functions (PDF) between noisy (source) and clean (target) features are modeled by a set GMMs using two-channel (so-called "stereo") data;
- 2. Conditional PDFs of clean features for given noisy ones are estimated from the joint PDFs;
- 3. The mapped clean features are determined so as to minimize their mean-square error (MMSE) from the conditional PDFs

Although this technique makes it possible to continuously transform of any sample of the source into that of target, its mapping performance is still insufficient. One of major factors which deteriorate the mapping quality is inappropriate dynamic characteristics caused by the frame-by-frame mapping. To alleviate this problem, GMM-based feature mapping technique with dynamic feature constraints have been proposed in voice conversion area [3]. The use of dynamic features makes it possible to convert a source feature vector sequence to a target one while satisfying the statistics of both static and dynamic features. However, it also introduces an inconsistency between training and mapping [4]. Generally, dynamic features are calculated as regression coefficients from static features of

Heiga Zen is now with the Speech Technology Group at Toshiba Research Europe Ltd. Cambridge Research Laboratory in Cambridge, CB4 0GZ, UK.

their neighboring frames. Thus, the relationship between the static and dynamic features is deterministic. In the conventional mapping techniques with dynamic features, this relationship is ignored on both the joint PDF training and conditional PDF estimation stages but utilized on the mapping stage [3].

Recently, a trajectory model, derived from the HMM by imposing the explicit relationships between static and dynamic features, was proposed [4]. This model, named trajectory HMM, can overcome the frame-wise conditional independence assumption of state output probabilities and piecewise constant statistics of the HMM, without any additional parameters. Based on this idea, Zen et al. also proposed the continuous stochastic feature mapping technique based on trajectory GMMs [5], i.e. joint PDFs between source and target feature vectors are modeled by a set of trajectory GMMs and the conditional PDFs are estimated from them. This technique can solve the inconsistency introduced by the use of dynamic features while keeping its benefits, and offers an entire sequence-level transformation rather than the frame-by-frame mapping. In this paper, we apply this new mapping technique to noise compensation.

This paper is organized as follows. Section 2 describes the mapping technique based on trajectory GMMs. Section 3 shows experimental results in noise compensation on the AURORA-2 task. Concluding remarks and future plans are given in the final section.

## 2. NOISE COMPENSATION BASED ON TRAJECTORY **GMMS**

# 2.1. Joint Probability Density Function

Noisy and clean MFCC vector sequences, x and y, are written as

$$x = \begin{bmatrix} x_1^\top, \dots, x_T^\top \end{bmatrix}^\top, \quad y = \begin{bmatrix} y_1^\top, \dots, y_T^\top \end{bmatrix}^\top,$$
 (1)

where  $x_t$  and  $y_t$  are the M-dimensional noisy and clean MFCC static feature vectors at the t-th frame, respectively, and T is the total number of frames. In this paper, the joint probability of x and y is modeled by a trajectory GMM [5] as follows:

$$p(z \mid \lambda) = \sum_{\forall q} P(q \mid \lambda) p(z \mid q, \lambda), \qquad (2)$$

$$p(z \mid q, \lambda) = \mathcal{N}\left(z \; ; \; \bar{z}_q, P_q\right),$$
 (3)

$$p(z \mid q, \lambda) = \mathcal{N}\left(z \; ; \; \bar{z}_{q}, P_{q}\right), \tag{3}$$

$$P(q \mid \lambda) = \prod_{t=1}^{T} c_{q_{t}}, \tag{4}$$

$$z = \begin{bmatrix} z_1^\top, \dots, z_T^\top \end{bmatrix}^\top, \quad z_t = \begin{bmatrix} x_t^\top, y_t^\top \end{bmatrix}^\top, \tag{5}$$

where  $\lambda$  denotes the set of model parameters, z is an  $2MT \times 1$  joint MFCC vector sequence,  $z_t$  is the  $2M \times 1$  joint MFCC vector at the tth frame,  $q = \{q_1, \dots, q_T\}$  is a Gaussian component sequence,  $q_t \in$ 

 $\{1, \ldots, N\}$  is the Gaussian component at the t-th frame, N is the total number of Gaussian components in the model set, and  $c_i$  is the mixture prior probability of the *i*-th Gaussian component. In Eq. (3),  $ar{z}_{m{q}}$  and  $m{P}_{m{q}}$  are the  $2MT \times 1$  mean vector and the  $2MT \times 2MT$ covariance matrix for q, respectively. They are given by

$$R_{\boldsymbol{q}}\bar{z}_{\boldsymbol{q}} = r_{\boldsymbol{q}},\tag{6}$$

$$R_{\boldsymbol{q}} = \boldsymbol{W}^{\top} \boldsymbol{\Omega}_{\boldsymbol{q}} \boldsymbol{W} = \boldsymbol{P}_{\boldsymbol{q}}^{-1}, \tag{7}$$

$$r_{\mathbf{q}} = W^{\top} \Omega_{\mathbf{q}} \mu_{\mathbf{q}}, \tag{8}$$

$$\boldsymbol{\mu}_{\boldsymbol{q}} = \left[\boldsymbol{\mu}_{q_1}^\top, \dots, \boldsymbol{\mu}_{q_T}^\top\right]^\top, \tag{9}$$

$$\mu_i = \left[ \mu_i^{(\mathbf{x})\top}, \mu_i^{(\mathbf{y})\top} \right]^{\top}, \tag{10}$$

$$\Omega_{\mathbf{q}} = \operatorname{diag} \left[ \Omega_{q_1}, \dots, \Omega_{q_T} \right],$$
 (11)

$$\Omega_{i} = \begin{bmatrix} \Omega_{i}^{(xx)} & \Omega_{i}^{(xy)} \\ \Omega_{i}^{(yx)} & \Omega_{i}^{(yy)} \end{bmatrix}, \tag{12}$$

where  $\mu_i$  and  $\Omega_i$  are the  $6M \times 1$  mean vector and the  $6M \times 6M$ precision (inverse covariance) matrix associated with the i-th Gaussian component, respectively, and W is a  $6MT \times 2MT$  window matrix which appends dynamic features to z. Unfortunately, estimating this model based on the ML criterion via the exact EM algorithm is computationally infeasible because it is intractable to compute  $p(q \mid z, \lambda)$ . However, approximate training technique based on the Viterbi approximation or Markov chain Monte Carlo (MCMC) has been proposed.

Equations (2)–(4) show that the trajectory GMM can be interpreted as an 2MT-dimensional GMM whose mixture weights are given by products of the mixture prior probabilities. Note that the intra and inter-frame covariance matrix of Eq. (3),  $P_q$ , is generally full. Therefore, the trajectory GMM can model both intra and inter-frame dependencies between noisy and clean MFCC vector sequences without increasing the number of model parameters compared to the GMM with the same model topology.

## 2.2. Mapping

### 2.2.1. Conditional probability

Equation (3) can be rewritten as follows:

$$p(z \mid q, \lambda) = \mathcal{N}\left(\begin{bmatrix} x \\ y \end{bmatrix}; \begin{bmatrix} \bar{x}_q \\ \bar{y}_q \end{bmatrix}, \begin{bmatrix} P_q^{(xx)} & P_q^{(xy)} \\ P_q^{(yx)} & P_q^{(yy)} \end{bmatrix}\right), \quad (13)$$

where

$$P_q^{(xx)} = C_q^{(xx)^{-1}}, \quad P_q^{(yy)} = C_q^{(yy)^{-1}},$$
 (14)

$$P_{q}^{(xy)} = -R_{q}^{(xx)^{-1}} R_{q}^{(xy)} C_{q}^{(yy)^{-1}} = P_{q}^{(yx)^{\top}}, \quad (15)$$

$$C_q^{(xx)} = R_q^{(xx)} - R_q^{(xy)} R_q^{(yy)^{-1}} R_q^{(yx)}, \tag{16}$$

$$C_q^{(yy)} = R_q^{(yy)} - R_q^{(yx)} R_q^{(xx)^{-1}} R_q^{(xy)}, \tag{17}$$

$$R_q^{(xx)} = W^{(x)}^{\top} \Omega_q^{(xx)} W^{(x)},$$
 (18)

$$R_q^{(yy)} = W^{(y)}^{\top} \Omega_q^{(yy)} W^{(y)}, \tag{19}$$

$$R_q^{(xy)} = W^{(x)}^{\top} \Omega_q^{(xy)} W^{(y)} = R_q^{(yx)}^{\top},$$
 (20)

$$\Omega_{q}^{(xx)} = \operatorname{diag}\left[\Omega_{q_1}^{(xx)}, \dots, \Omega_{q_T}^{(xx)}\right], \tag{21}$$

$$\mathbf{\Omega}_{q}^{(yy)} = \operatorname{diag}\left[\mathbf{\Omega}_{q_{1}}^{(yy)}, \dots, \mathbf{\Omega}_{q_{T}}^{(yy)}\right], \tag{22}$$

$$\Omega_{\mathbf{q}}^{(\mathbf{x}\mathbf{y})} = \operatorname{diag}\left[\Omega_{q_1}^{(\mathbf{x}\mathbf{y})}, \dots, \Omega_{q_T}^{(\mathbf{x}\mathbf{y})}\right] = \Omega_{\mathbf{q}}^{(\mathbf{y}\mathbf{x})},$$
 (23)

and  $W^{(x)} = W^{(y)}$  is the  $3MT \times MT$  window matrix which appends dynamic features to x and y. The mean vectors of Eq. (13),  $\bar{x}_q$  and  $\bar{y}_q$ , are given as follows:

$$\bar{x}_q = P_q^{(xx)} \left( r_q^{(x)} - R_q^{(xy)} R_q^{(yy)^{-1}} r_q^{(y)} \right),$$
 (24)

$$\bar{y}_q = P_q^{(yy)} \left( r_q^{(y)} - R_q^{(yx)} R_q^{(xx)^{-1}} r_q^{(x)} \right),$$
 (25)

where

$$r_q^{(x)} = W^{(x)}^{\top} \left( \Omega_q^{(xx)} \mu_q^{(x)} + \Omega_q^{(xy)} \mu_q^{(y)} \right),$$
 (26)

$$r_q^{(y)} = W^{(y)^{\top}} \left( \Omega_q^{(yy)} \mu_q^{(y)} + \Omega_q^{(yx)} \mu_q^{(x)} \right),$$
 (27)

$$\boldsymbol{\mu}_{\boldsymbol{q}}^{(\boldsymbol{x})} = \left[\boldsymbol{\mu}_{q_1}^{(\boldsymbol{x})^{\top}}, \dots, \boldsymbol{\mu}_{q_T}^{(\boldsymbol{x})^{\top}}\right]^{\top}, \tag{28}$$

$$\boldsymbol{\mu}_{\boldsymbol{q}}^{(\boldsymbol{y})} = \left[\boldsymbol{\mu}_{q_1}^{(\boldsymbol{y})^{\top}}, \dots, \boldsymbol{\mu}_{q_T}^{(\boldsymbol{y})^{\top}}\right]^{\top}.$$
 (29)

As a result, the conditional PDF of y for given x and  $\lambda$  can be expressed as

$$p(y \mid x, \lambda) = \sum_{\forall q} \gamma_{q} \cdot p(y \mid x, q, \lambda), \qquad (30)$$

$$p(y \mid x, q, \lambda) = \mathcal{N}\left(y \; ; \; \tilde{y}_{q}, \tilde{P}_{q}^{(yy)}\right),$$
 (31)

$$\tilde{y}_{\boldsymbol{q}} = \bar{y}_{\boldsymbol{q}} + P_{\boldsymbol{q}}^{(yx)} C_{\boldsymbol{q}}^{(xx)} \left( x - \bar{x}_{\boldsymbol{q}} \right), \tag{32}$$

$$\tilde{y}_{q} = \tilde{y}_{q} + P_{q}^{(yx)} C_{q}^{(xx)} (x - \tilde{x}_{q}), \qquad (32)$$

$$\tilde{P}_{q}^{(yy)} = P_{q}^{(yy)} - P_{q}^{(yx)} C_{q}^{(xx)} P_{q}^{(xy)}, \qquad (33)$$

where  $\gamma_q$  is a posterior probability of q.

## 2.2.2. MMSE-based mapping

The MMSE estimates of the clean static MFCC vector sequence  $\hat{y}$ is determined as follows:

$$\hat{y} = E[y \mid x] = \int p(y \mid x, \lambda) y dy$$
 (34)

$$= \int \sum_{\forall q} \gamma_q \cdot p(y \mid x, q, \lambda) y dy$$
 (35)

$$= \sum_{\forall q} \gamma_q \cdot \tilde{y}_q, \tag{36}$$

where  $E[\cdot]$  means expectation. The estimated MFCC static feature vector sequence  $\hat{y}$  is defined as the weighted sum of the conditional mean vectors  $\tilde{y}_{q}$ , where the posterior probabilities of q are used as weights.

Unfortunately, computing  $\gamma_q$  is computationally intractable. Hence, we also need to use the Viterbi or MCMC approximation like the ML estimation of trajectory GMMs. In the all experiments described in this paper, we approximated the marginalization over all possible q by a single  $\hat{q}$  and did not iterate the EM algorithm.

#### 3. EXPERIMENTS

## 3.1. Experimental conditions

Speech recognition experiments were conducted under various noise conditions using the Aurora-2 database and task [6]. The Aurora-2 database is a subset of the TI-DIGITS, which contains a set of connected digit utterances spoken in English; while the task consists of the recognition of the connected digit utterances interfered with various noise sources at different signal-to-noise ratios (SNRs), in which the test-sets A (seen noise types) and B (unseen noise types) are artificially contaminated with eight different types of real world noises in a wide range of SNRs (-5dB, 0dB, 5dB, 10dB, 15dB, 20dB, and Clean) and the test-set C additionally includes the channel distortion. The acoustic model for each digit was a left-to-right continuous density HMM with 16 states, and each state had a mixture of three Gaussian distributions. Two additional silence models were defined. One had three states with a mixture of six Gaussian distributions per state for modeling the silence at the beginning and at the end of each utterance. The other one had one state with a mixture of six Gaussian distributions for modeling the inter-word short pause. A 39-dimensional feature vector was composed at each frame, including 12 MFCCs and the logarithm of the energy extracted by the reference Aurora front-end version 2.0, as well as their corresponding dynamic coefficients. The training and recognition tests used HTK, which followed the setup originally defined for the ETSI evaluations. All the experimental results reported below are based on clean-condition training, i.e., the acoustic models were trained only with the clean (uncontaminated) training utterances. The clean acoustic model training data consists of 8,440 utterances. The multistyle acoustic model training data was used to train mapping functions for noise compensation. It consists of the same utterances synthetically mixed with four different noise types (Subway, Car, Babble, and Exhibition) at varying amplitudes (5dB, 10dB, 15dB, 20dB, and Clean), for a total of 17 unique noise conditions. Using the pairs of clean and noisy speech data, joint static feature vectors consisted of pairs of clean and noisy MFCCs were composed. They were augmented with their dynamic features (78 dimensions in total)

In this experiment, four types of mapping techniques were evaluated, which were

- GMM-Static: The GMM-based mapping [2] was applied to the 13-dimensional static MFCC vectors;
- GMM-Complete: Like Complete SPLICE [7], the 13-dimensional static MFCC vectors were augmented with their dynamic features before passing them to the GMM-based mapping. Note that the dynamic features in mapped 39-dimensional vectors were no longer correspond to true dynamic features:
- GMM-Dynamic: The 13-dimensional static MFCC vectors were mapped by GMMs under the constraints between static and dynamic features [3];
- Trajectory GMM: The 13-dimensional static MFCC vectors were transformed by the trajectory GMM-based mapping;

All of the four mapping techniques used a mixture of 256 Gaussian components to model the joint probability of clean and noisy static MFCC vectors for each of the 17 conditions. The diagonal structure was used for  $\Omega_i^{(xx)}$ ,  $\Omega_i^{(xy)}$ ,  $\Omega_i^{(yx)}$ , and  $\Omega_i^{(yy)}$ . The GMMs for GMM-Static were estimated using 26-dimensional joint static feature vector sequences, and the others were estimated using 78-dimensional augmented ones (with dynamic features). Note that

the GMMs for GMM-Complete and GMM-Dynamic were identical because the difference between these two techniques is the use of dynamic feature constraints at the mapping stage only. The trajectory GMMs were re-estimated by the Viterbi training based on the ML criterion [4] using the GMMs for GMM-Dynamic as their initial models.

To select the mapping function in the decoding stage, a GMM-based noise condition classifier was built for the detection of noise type and SNR. In this noise condition classifier, a GMM with four Gaussian components was estimated for each combination of noise type and SNR using the first 10 frames of the utterances from that noise condition in the training data. Before decoding, the likelihood of the first 10 frames of the input utterances was computed against each GMM. The noise condition with the maximum likelihood was chosen as the noise condition of the input utterance, and the mapping of that noise condition was applied to compensate the utterance unless the noise condition was Clean. The dynamic features of mapped static features by GMM-Complete, GMM-Dynamic, and Trajectory GMM were computed on-line during recognition.

# 3.2. Experimental results

Figure 1 is a summary of the full results on the Aurora-2 corpus. Note that the performance of the baseline system (clean condition training, without noise compensation) was the same as that of the ETSI reference system described in [6], and the average accuracy and relative error reduction were computed with the results between 20 and 0 dB as suggested in [6]. It can be seen from the figure that all of four mapping techniques achieved significant improvements over the baseline system. Compared with the other three techniques, the performance of GMM-Static was significantly worse. It demonstrates the effectiveness of augmenting static features by their dynamic features in noise compensation. Although the GMMs of GMM-Dynamic and GMM-Complete were identical, GMM-Dynamic slightly outperformed GMM-Complete on test-set A (seen noise types). On the other hand, on testset B (unseen noise types), GMM-Complete achieved the lower WERs than GMM-Dynamic. Similar tendency can also be observed in Trajectory GMM: Trajectory GMM outperformed GMM-Complete on test-set A but GMM-Complete achieved the lower WERs than Trajectory GMM on test-set B. For seen noise types with unseen SNRs (SNR=0 or SNR=-5dB in testset A), GMM-Dynamic and Trajectory GMM outperformed GMM-Complete. These results suggest that incorporating dynamic feature constraints in the mapping stage improves the mapping performance for seen noise types but degrades the recognition performance for unseen noise types. We expect that the dynamic feature constraints caused over-fitting to the training noise types. In total, Trajectory GMM achieved the lowest WERs among the four mapping techniques.

# 4. CONCLUSIONS

This paper evaluated the performance of stochastic feature mapping technique based on trajectory GMMs on noise compensation. Experimental results showed that the use of dynamic feature constraints in noise compensation worked effective in seen noise types. It also indicated that the use of dynamic features degraded the mapping performance due to over-fitting to the training data. In total, the proposed noise compensation algorithm outperformed the conventional mapping techniques on the AURORA-2 task.

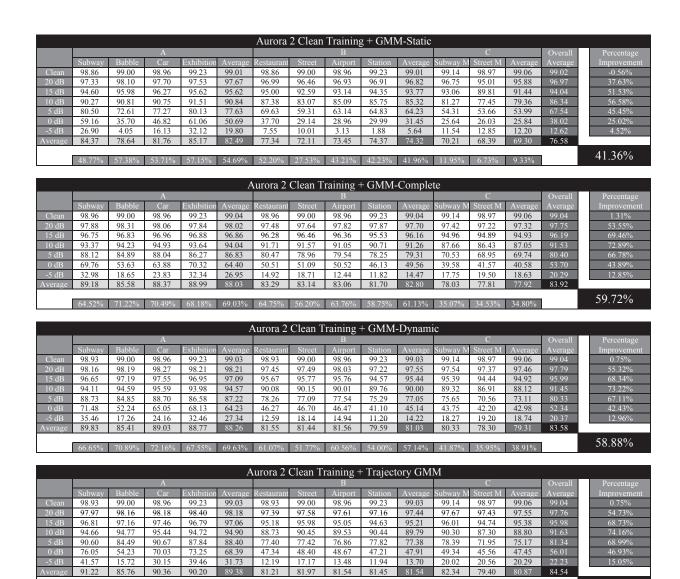


Fig. 1. Summary of the full results by GMM and trajectory GMM-based noise compensation techniques on the AURORA-2 task.

# 5. ACKNOWLEDGMENTS

This work was partly supported by the MEXT e-Society project, the Hori information science promotion foundation, and the Grant-in-Aid for Scientific Research (No. 1880009) of JSPS.

## 6. REFERENCES

- L. Deng, A. Acero, L. Jiang, J. Droppo, and X. Huang, "Highperformance robust speech recognition using stereo training data," in *Proc. ICASSP*, 2001, pp. 301–304.
- [2] X.-D. Cui, M. Afify, and Y. Gao, "MMSE-based stereo feature stochastic mapping for noise robust speech recognition," in *Proc. ICASSP*, 2008, pp. 4077–4080.
- [3] T. Toda, A.W. Black, and K. Tokuda, "Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory," *IEEE Trans. ASLP*, vol. 15, no. 8, pp. 2222–2235, 2007.

[4] H. Zen, K. Tokuda, and T. Kitamura, "Reformulating the HMM as a trajectory model by imposing explicit relationships between static and dynamic features," *Computer Speech & Language*, vol. 21, no. 1, pp. 153–173, 2007.

61.30%

- [5] H. Zen, Y. Nankaku, and K. Tokuda, "Probabilistic feature mapping based on trajectory HMMs," in *Proc. Interspeech*, 2008.
- [6] H.G. Hirsch and D. Pearce, "The AURORA experimental framework for the performance evaluations of speech recognition systems under noisy conditions," in *Proc. ISCA ITRW ASR-2000*, 2000, pp. 181–188.
- [7] J. Droppo, A. Acero, and L. Deng, "Uncertainty decoding with SPLICE for noise robust speech recognition," in *Proc. ICASSP*, 2002, pp. 57–60.