



## BANDWIDTH EXTENSION OF CELLULAR PHONE SPEECH BASED ON MAXIMUM LIKELIHOOD ESTIMATION WITH GMM

Wataru Fujitsuru<sup>†</sup>, Hidehiko Sekimoto<sup>†</sup>, Tomoki Toda<sup>†</sup>, Hiroshi Saruwatari<sup>†</sup> and Kiyohiro Shikano<sup>†</sup>

<sup>†</sup>Graduate School of Information Science, Nara Institute of Science and Technology, Nara, 631-0101 Japan  
E-mail: <sup>†</sup> wataru-f, hidehi-s, tomoki, sawatari, shikano @is.naist.jp

### Abstract

Bandwidth extension is a useful technique for reconstructing wideband speech from only narrowband speech. As a typical conventional method, a bandwidth extension algorithm based on minimum mean square error (MMSE) with a Gaussian mixture model (GMM) has been proposed [1]. Although the MMSE-based method has reasonably high conversion-accuracy, there still remain some problems to be solved: 1) inappropriate spectral movements are caused by ignoring a correlation between frames, and 2) the converted spectra are excessively smoothed by the statistical modeling. In order to address these problems, we propose a bandwidth extension algorithm based on maximum likelihood estimation (MLE) considering dynamic features and the global variance (GV) with a GMM. A result of a subjective test demonstrates that the proposed algorithm outperforms the conventional MMSE-based algorithm.

### 1. Introduction

The use of cellular phones enables us to easily communicate with each other through speech. In general, the cellular phone speech is limited to a narrowband signal up to 3.4 kHz. Although narrowband speech is capable of intelligible communication, its sound quality is not high enough. Specifically considerable quality degradation is observable at several phonemes such as fricatives and plosives which have important energy distribution beyond 3.4 kHz. In order to realize higher-quality speech communication, a wideband speech codec has been developed [2]. Essentially it needs more information than the narrowband speech codec. It is no doubtful that to realize wideband speech communication while not increasing information is more convenient.

Bandwidth extension is a useful technique for reconstructing wideband speech from only narrowband speech. Several statistical approaches to bandwidth extension based on a spectral mapping have been studied. A codebook mapping method [3] is an approach based on hard clustering and discrete mapping. A reconstructed feature vector at each frame is determined by quantizing the narrowband speech feature vector to the nearest centroid vector of the narrowband codebook and substituting it with a corresponding wideband centroid vector of the mapping codebook. One of more sophisticated approaches allowing soft clustering and continuous mapping is a probabilistic bandwidth extension method based on a Gaussian mixture model (GMM) [4]. The basic mapping algorithm has originally been proposed for voice conversion [5]. Most of the conventional GMM-based methods perform the mapping based on MMSE criterion [6, 7]. It has been reported that the mapping performance is further improved by modeling dynamic characteristics of a spectral sequence

[8, 9]. Moreover, an approach of combining mapping and coding processes has been studied [10].

Recently the performance of voice conversion with a GMM was significantly improved by maximum likelihood estimation (MLE) considering dynamic features and the global variance (GV) [11]. It is expected that it also causes the performance improvement of bandwidth extension. This paper proposes bandwidth extension based on MLE considering dynamic features and the GV. The effectiveness of considering dynamic features and the GV is demonstrated by a subjective evaluation.

This paper is organized as follows. Section 2 describes the conventional method of the bandwidth extension algorithm based on GMM. Section 3 describes the proposed bandwidth extension based on MLE. Section 4 describes the flow of the bandwidth extension. Section 5 describes an experimental evaluation. Finally, we summarize this paper in Section 6.

### 2. MMSE-based bandwidth extension [12]

#### 2.1. Training

Let  $\mathbf{x}_t$  and  $\mathbf{y}_t$  be  $D_x$ -dimensional narrowband and  $D_y$ -dimensional wideband feature vectors at frame  $t$ , respectively. Joint probability density of the narrowband and wideband feature vectors is modeled by a GMM as follows:

$$P(\mathbf{z}_t|\theta) = \sum_{m=1}^M \omega_m \mathcal{N}(\mathbf{z}_t; \boldsymbol{\mu}_m^{(z)}, \boldsymbol{\Sigma}_m^{(z)}), \quad (1)$$

where  $\mathbf{z}_t$  is a joint vector  $[\mathbf{x}_t^T, \mathbf{y}_t^T]^T$ . The notation  $\top$  denotes transposition of the vector. The number of mixture components is  $M$ . The weight of the  $m$ -th mixture component is  $\omega_m$ . The normal distribution with  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  is denoted as  $\mathcal{N}(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ . A parameter set of the GMM is  $\theta$ , which consists of weights, mean vectors and the covariance matrices for individual mixture components. The mean vector  $\boldsymbol{\mu}_m^{(z)}$  and the covariance matrix  $\boldsymbol{\Sigma}_m^{(z)}$  of the  $m$ -th mixture component are written as

$$\boldsymbol{\mu}_m^{(z)} = \begin{bmatrix} \boldsymbol{\mu}_m^{(x)} \\ \boldsymbol{\mu}_m^{(y)} \end{bmatrix}, \boldsymbol{\Sigma}_m^{(z)} = \begin{bmatrix} \boldsymbol{\Sigma}_m^{(xx)} & \boldsymbol{\Sigma}_m^{(xy)} \\ \boldsymbol{\Sigma}_m^{(yx)} & \boldsymbol{\Sigma}_m^{(yy)} \end{bmatrix}. \quad (2)$$

where  $\boldsymbol{\mu}_m^{(x)}$  and  $\boldsymbol{\mu}_m^{(y)}$  are the mean vector of the  $m$ -th mixture component for the narrowband and that for the wideband, respectively, and  $\boldsymbol{\Sigma}_m^{(xx)}$  and  $\boldsymbol{\Sigma}_m^{(yy)}$  are the covariance matrix of the  $m$ -th mixture component for the narrowband and that for the wideband, respectively. The matrix  $\boldsymbol{\Sigma}_m^{(yx)}$  is the cross-covariance matrix of the  $m$ -th mixture component for the narrowband and wideband.

The GMM is trained with the EM algorithm using the joint vectors in a training set. This training method robustly estimates model parameters compared with the least squares estimation [1].

## 2.2. MMSE-based conversion

The conventional method performs the conversion based on MMSE as follows:

$$\begin{aligned}\hat{y}_t &= E[y_t|x_t] \\ &= \int P(y_t|x_t, \theta) y_t dy_t \\ &= \sum_{m=1}^M P(m|x_t, \theta) E_{m,t},\end{aligned}\quad (3)$$

where

$$P(m|x_t, \theta) = \frac{\omega_m \mathcal{N}(x_t; \mu_m^{(y)}, \Sigma_m^{(xx)})}{\sum_{j=1}^M \omega_j \mathcal{N}(x_t; \mu_j^{(x)}, \Sigma_j^{(xx)})}, \quad (4)$$

$$E_{m,t} = \mu_m^{(y)} + \Sigma_m^{(yx)} (\Sigma_m^{(xx)})^{-1} (x_t - \mu_m^{(x)}). \quad (5)$$

## 2.3. Problems

Although MMSE has reasonably high conversion-accuracy, there still remain some problems to be solved: 1) inappropriate spectral movements are caused by ignoring a correlation between frames, and 2) the converted spectra are excessively smoothed by the statistical modeling.

## 3. MLE-based bandwidth extension

In order to solve two main problems of the conventional method, we proposed the bandwidth extension based on MLE considering dynamic features and the GV. Inter-frame correlation is considered for realizing the converted parameter trajectory with appropriate dynamic characteristics. The over-smoothing effects are alleviated by considering the GV capturing one of characteristics of the trajectory.

### 3.1. MLE considering dynamic features (ML) [1]

We use  $2D_x$ -dimensional narrowband speech feature vector  $X_t = [x_t^T, \Delta x_t^T]^T$  and  $2D_y$ -dimensional wideband speech feature vector  $Y_t = [y_t^T, \Delta y_t^T]^T$ , consisting of static and dynamic features at frame  $t$ . Their time sequences are written as  $X = [X_1^T, X_2^T, \dots, X_T^T]^T$  and  $Y = [Y_1^T, Y_2^T, \dots, Y_T^T]^T$ , respectively. A time sequence of the converted static feature vectors  $\hat{y} = [\hat{y}_1^T, \dots, \hat{y}_T^T]^T$  is determined as follows:

$$\begin{aligned}\hat{y} &= \arg \max_{\mathbf{y}} P(Y|X, \Theta) \\ &\text{subject to } Y = W\mathbf{y},\end{aligned}\quad (6)$$

where  $W$  is a conversion matrix that extends a sequence of the static feature vectors  $\mathbf{y}$  into that of the static and dynamic feature vectors.

In order to effectively reduce the computational cost, we approximate the likelihood as follows,

$$P(Y|X, \Theta) \approx P(m|X, \Theta) P(Y|X, m, \Theta), \quad (7)$$

where  $m$  is a mixture component sequence  $[m_1, m_2, \dots, m_T]$ . After determining the sub-optimum mixture sequence  $\hat{m}$  written as

$$\hat{m} = \arg \max_{\mathbf{m}} P(m|X, \Theta), \quad (8)$$

we determine  $\hat{y}$  that maximizes the approximated likelihood function as follows:

$$\begin{aligned}\hat{y} &= \arg \max_{\mathbf{y}} P(\hat{m}|X, \Theta) P(Y|X, \hat{m}, \Theta) \\ &= \left( W^T D_{\hat{m}}^{(Y)^{-1}} W \right)^{-1} W^T D_{\hat{m}}^{(Y)^{-1}} E_{\hat{m}}^{(Y)},\end{aligned}\quad (9)$$

where

$$E_{\hat{m}}^{(Y)} = [E_{\hat{m}_1,1}^{(Y)}, \dots, E_{\hat{m}_t,t}^{(Y)}, \dots, E_{\hat{m}_T,T}^{(Y)}], \quad (10)$$

$$D_{\hat{m}}^{(Y)^{-1}} = \text{diag} [D_{\hat{m}_1}^{(Y)^{-1}}, \dots, D_{\hat{m}_t}^{(Y)^{-1}}, \dots, D_{\hat{m}_T}^{(Y)^{-1}}]. \quad (11)$$

A GMM parameter set  $\Theta$  is estimated in advance with training data in the same manner as the converted method.

### 3.2. MLE considering GV (MLGV) [1]

The GV of the static feature vectors in each utterance is written as

$$\mathbf{v}(\mathbf{y}) = [v(1), v(2), \dots, v(D)]^T \quad (12)$$

$$v(d) = \frac{1}{T} \sum_{t=1}^T \left( y_t(d) - \frac{1}{T} \sum_{\tau=1}^T y_{\tau}(d) \right)^2 \quad (13)$$

where  $y_t^{(d)}$  is the  $d$ -th component of the target static feature at frame  $t$ .

The following likelihood function consisting of two probability densities for a sequence of the wideband feature vectors and for the GV of the wideband static feature vectors is maximized,

$$\begin{aligned}L &= \alpha \log P(Y|X, \hat{m}, \Theta) + \log P(\mathbf{v}(\mathbf{y})|\Theta_v), \\ &\text{subject to } Y = W\mathbf{y}\end{aligned}\quad (14)$$

where  $P(\mathbf{v}(\mathbf{y})|\Theta_v)$  is modeled by the normal distribution. A set of model parameters  $\Theta_v$  consists of the mean vector  $\mu^{(v)}$  and the covariance matrix  $\Sigma^{(vv)}$  for the GV vector  $\mathbf{v}(\mathbf{y})$ , which is also estimated in advance with training data. The constant  $\alpha$  is a likelihood weight. In this paper, it is set to  $\frac{1}{2T}$ . Note that the GV likelihood works as a penalty term for the over-smoothing.

In order to maximize the likelihood  $L$  with respect to  $\mathbf{y}$ , we employ a steepest descent algorithm using the first derivative,

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{y}} &= \alpha \left( -W^T \overline{D^{(Y)^{-1}}} W \mathbf{y} + W^T \overline{D^{(Y)^{-1}}} E^{(Y)} \right) \\ &\quad + [v_1^T, v_2^T, \dots, v_t^T, \dots, v_T^T]^T,\end{aligned}\quad (15)$$

$$v_t^T = [v_t'(1), v_t'(2), \dots, v_t'(d), \dots, v_t'(D)]^T, \quad (16)$$

$$v_t'(d) = -\frac{2}{T} P_v^{(d)^T} (\mathbf{v}(\mathbf{y}) - \mu_v) \left( y_t(d) - \frac{1}{T} \sum_{\tau=1}^T y_{\tau}(d) \right). \quad (17)$$

The vector  $P_v^{(d)}$  is the  $d$ -th column vector of  $P_v = \Sigma^{(vv)^{-1}}$ .

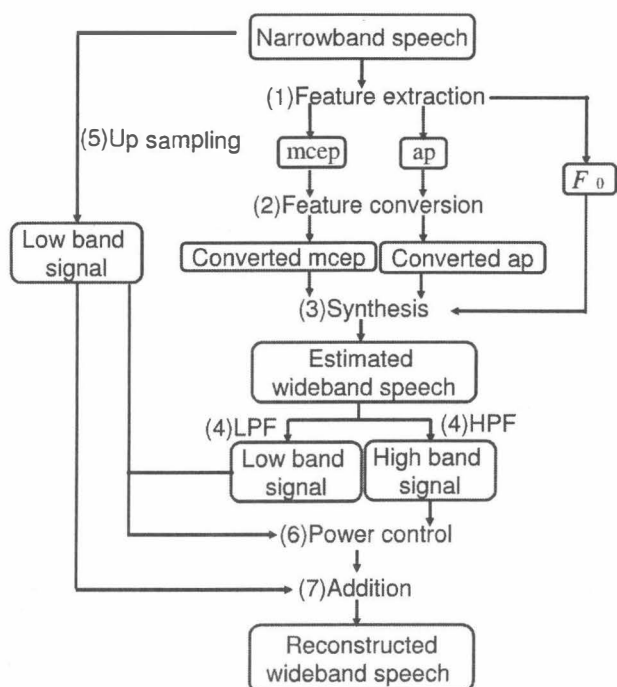


Figure 1: Bandwidth extension system. “mcep” denotes the mel-cepstrum and “ap” denotes the aperiodic component.

#### 4. Details of bandwidth extension process

Figure 1 shows a flow of the bandwidth extension.

- Step 1** Extracting  $F_0$ , mel-cepstrum and aperiodic component [ ] as speech features from the narrowband speech signal.
- Step 2** Converting mel-cepstrum and aperiodic component of the narrowband speech into those of the wideband speech.
- Step 3** Generating STRAIGHT mixed excitation [ ] using the extracted  $F_0$  and the converted aperiodic component, and then synthesizing the estimated wideband speech with MLSA filter [ ] based on the converted mel-cepstrum.
- Step 4** The estimated wideband speech is separated into a low-band signal and a high-band signal with a low-pass filter (LPF) and a high-pass filter (HPF).
- Step 5** The input narrowband speech is converted into a low-band signal with up-sampling.
- Step 6** Power of the estimated high-band signal is adjusted so that power of the estimated low-band signal is equal to that of the input low-band signal.
- Step 7** Reconstructing the wideband speech by adding the resulting high-band signal and the low-band signal.

#### 5. Experimental evaluations

In order to demonstrate the effectiveness of the proposed method, we conducted a subjective evaluation.

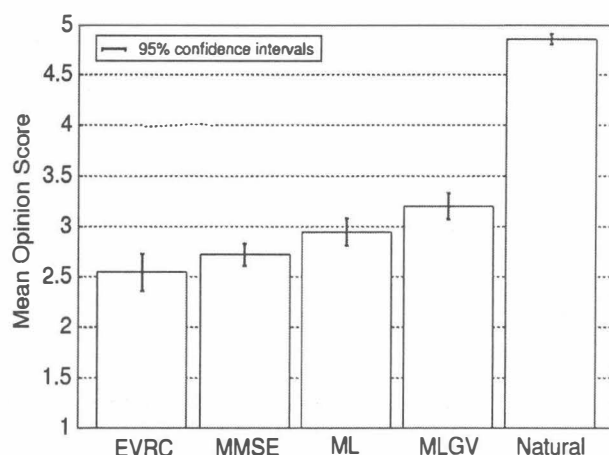


Figure 2: Result of subjective evaluation.

#### 5.1. Experimental conditions

We used 16 kHz sampled natural speech of 4 Japanese speakers (2 males, 2 females) as the wideband speech. The 3.4 kHz narrowband speech was prepared by down-sampling the wideband speech and then passing it through EVRC (Enhanced Variable Rate Codec) [ ]. The training data was 50 sentences from subset A of ATR’s phonetically balanced sentence database. The evaluation data was 50 sentences from subset B of ATR’s phonetically balanced sentence database. For narrowband, we used the 16-dimensional mel-cepstral coefficients from the mel-cepstral analysis [ ]. For wideband, we used the 24-dimensional mel-cepstral coefficients from the STRAIGHT analysis [ ]. We used the averaged aperiodic components [ ] on three frequency bands (0 to 1, 1 to 2, 2 to 4 kHz) for narrowband and those on five frequency bands (0 to 1, 1 to 2, 2 to 4, 4 to 6, 6 to 8 kHz) for wideband. The frame shift was 5 ms. The number of mixture components of the GMM for mel-cepstral conversion was set to 64. The number of mixture components of the GMM for aperiodic components conversion was set to 4. Speaker dependent models were evaluated.

We conducted an opinion test on speech quality. An opinion score was set to a 5-point scale (5: excellent, 4: good, 3: fair, 2: poor, 1: bad). The evaluated speech samples consisted of EVRC, MMSE, ML, MLGV, and wideband natural speech (Natural). The listeners were eight Japanese adult man and woman.

#### 5.2. Experimental results

Figure 2 shows a result of the opinion test. There is no significant difference between EVRC and MMSE. On the other hand, the proposed method ML is significantly better than both EVRC and MMSE. The reconstructed wideband speech with the highest speech quality was obtained by considering the GV as well in the proposed method. An example of spectral sequences of the narrowband speech, the reconstructed speech and the natural wideband speech is shown in Figure 3. The proposed method estimates spectral envelopes considering inter-frame correlation while alleviating the over-smoothing effects.

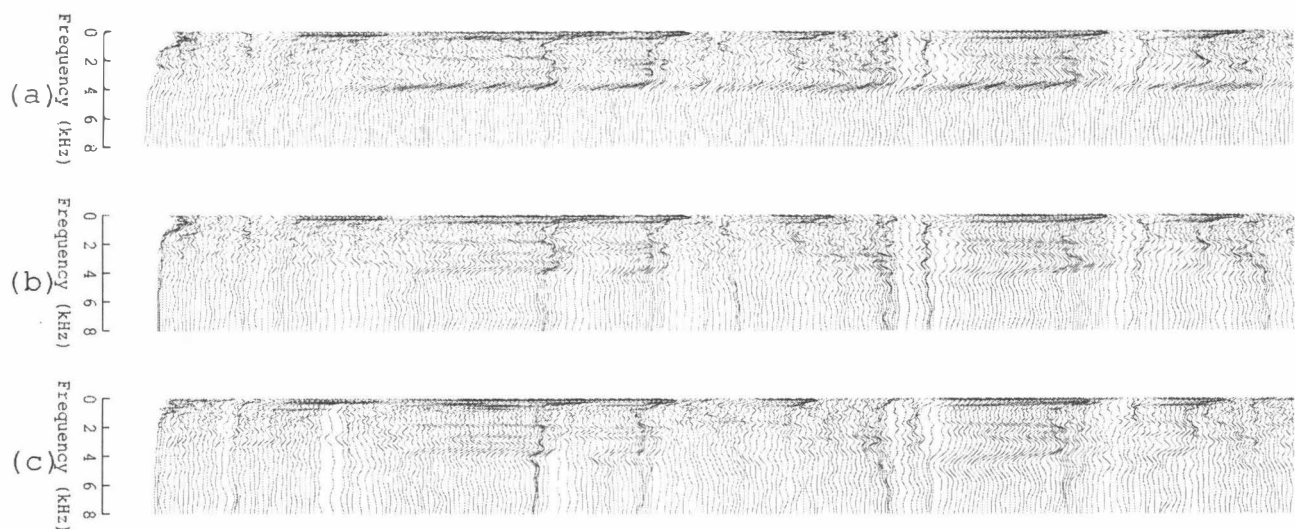


Figure 3: An example of spectral sequences of narrowband speech, (a) converted speech by the ML using the GV, (b) spectra of natural speech, (c) for a sentence fragment, “/jy u: j i t s U sh l t e i k u f’”.

## 6. Conclusions

We proposed bandwidth extension based on maximum likelihood estimation (MLE) with a Gaussian mixture model (GMM) considering dynamic features and the global variance (GV). A result of the subjective evaluation demonstrated that the speech quality of the narrowband speech is significantly improved by the proposed method. We plan to deal with a speaker independent model, online processing and noise robustness.

## 7. Acknowledgments

This research was supported in part by e-Society project and KDDI collaborative research.

### References

- [1] B. Bessette, R. Salami, R. Lefebvre, M. Jelinek, J. Potola-Pukkila, J. Vainio, H. Mikkola and K. Jarvinen, “The adaptive multirate wideband speech codec (AMR-WB),” *IEEE Trans.*, Vol. 10, No. 8, pp. 620–636, 2002.
- [2] Y. Yoshida, and M. Abe, “An Algorithm to Reconstruct Wideband Speech From Narrowband Speech Based on Codebook Mapping,” *Proc. ICSLP94*, pp. 1591–1593, 1994.
- [3] K.Y. Park and H.S. Kim. “Narrowband to wideband conversion of speech using GMM based transformation,” *Proc. ICSLP*, pp. 1847–1850, Istanbul, June, 2000.
- [4] Y. Stylianou, O. Cappe, E. Moulines, “Continuous probabilistic transform for voice conversion,” *IEEE Trans.*, Speech and Audio Processing, Vol. 6, No. 2, pp. 131–142, 1998.
- [5] M.L. Seltzer, A. Acero, and J. Droppo, “Robust Bandwidth Extension of Noise-corrupted Narrowband Speech,” *Proc. ICSLP*, pp. 1509–1512, 2005.
- [6] S. Yao and C.F. Chan, “Block-based Bandwidth Extension of Narrowband Speech Signal by using CDHMM,” *Proc. ICASSP*, pp. 1793–1796, 2005.
- [7] S.Y. Yao and C.F. Chan, “Speech bandwidth enhancement using state space speech dynamics,” *Proc. ICASSP2006*, pp. 1489–1492, 2006.
- [8] Y. Agiomyrgiannakis and Y. Stylianou, “Combined Estimation/coding of Highband Spectral Envelopes for Speech Spectrum Expansion,” *Proc. ICASSP2004*, pp. 469–472, 2004.
- [9] 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. 2222–2234, 2007.
- [10] A. Kain and M.W. Macon. “Spectral voice conversion for text-to-speech synthesis,” *Proc. ICASSP*, Seattle, U.S.A., pp. 285–288, May 2004.
- [11] H. Kawahara, Jo Estill and O. Fujimura, “Aperiodicity extraction and control using mixed mode excitation and group delay manipulation for a high quality speech analysis, modification and synthesis system STRAIGHT,” *Proc. MAVEBA*, Sep. 13–15, Firenze Italy, 2001.
- [12] H. Kawahara, I. Masuda-Katsuse, and A. de Cheveigné, “Restructuring speech representations using a pitch-adaptive time-frequency smoothing and an instantaneous-frequency-based  $F_0$  extraction: Possible role of a repetitive structure in sounds,” *Speech Commun.*, Vol. 27, No. 3–4, pp. 187–207. 1999.
- [13] T. Fukada, K. Tokuda, T. Kobayashi, and S. Imai, “An adaptive algorithm for mel-cepstral analysis of speech,” *Proc. ICASSP*, Vol 1, pp. 137–140, San Francisco, USA, Mar. 1992.
- [14] T.V. Ramabadran, J.P. Ashley and M.J. McLaughlin, “Background Noise Suppression for Speech Enhancement and Coding,” *IEEE Workshop on Speech Coding and Tel.*, Pocono Manor, PA, pp. 43–44, 1997.