# USING CONDITIONAL RESTRICTED BOLTZMANN MACHINES FOR SPECTRAL ENVELOPE MODELING IN SPEECH BANDWIDTH EXTENSION

*Yingxue Wang*<sup>1,3</sup>, *Shenghui Zhao*<sup>1</sup>, *Dan Qu*<sup>2,3</sup>, *Jingming Kuang*<sup>1</sup>

<sup>1</sup>School of Information and Electronics, Beijing Institute of Technology, Beijing 100081, China <sup>2</sup>Zhengzhou Institute of Information Science and Technology, Zhengzhou 450002, China <sup>3</sup>School of Computer Science, Carnegie Mellon University, US

yxwang.bit@gmail.com, shzhao@bit.edu.cn

### ABSTRACT

In this paper, we present a conditional restricted Boltzmann machine (CRBM) based speech bandwidth extension (BWE) method. A CRBM is employed to obtain time information and model deep non-linear relationships between the spectral envelope features of low frequency (LF) and high frequency (HF). Two exclusive CRBMs are adopted to model the distribution of LF's and HF's spectral envelope features. respectively. A neural network (NN) is then used to model the joint distribution of hidden variables extracted from the two CRBMs. The proposed method takes advantage of the strong ability of CRBM in discovering the temporal correlation between adjacent frames and modeling deep non-linear relationships between input and output. Both the objective and subjective evaluations indicate that our proposed method outperforms the conventional Gaussian mixture model based methods and other NN based methods.

*Index Terms*— Speech bandwidth expansion, conditional restricted Boltzmann machine, Gaussian mixture model

### **1. INTRODUCTION**

Speech bandwidth expansion (BWE) is a technique that regenerates the missing high frequency parts in order to improve the quality of speech. Many approaches have been proposed for BWE during the last decades. Among these approaches, mapping method based on Gaussian Mixture Model (GMM) [1] is widely used and a number of improvements [2, 3, 4, 5] have been proposed.

However, the derived mapping function by GMM is a piece-wise linear transformation which is maybe insufficient to describe the complex non-linear relationship between the low frequency (LF) and high frequency (HF). To capture the characteristics of speech more precisely, a deeper non-linear architecture is required. One example of deeper BWE methods was imposed by [6] based on deep neural network (DNN). Previously, we also proposed to use the restricted Boltzmann machines (RBM) [7] to obtain a non-linear relationship between LF and HF. Although these approaches were reported to outperform the GMM-based approaches, they assumed the

speech frames were independent of each other, ignoring the temporal information of speech. The importance of the temporal information as well as its advantage had been investigated appropriately in several literatures. In [8, 9, 10], temporal information was included in GMM-based method by utilizing the delta features of spectral envelopes. In [11], time correlation properties of speech were embedded into spectrum estimation by utilizing the Hidden Markov Model (HMM).

In this paper, we propose to use conditional restricted Boltzmann machine (CRBM) for spectral envelope modeling to avoid the frame-by-frame independence assumption in most of the speech BWE methods. We do this by combining two CRBMs and a concatenation neural network (NN). A CRBM is a non-linear probabilistic model used to capture temporal information of speech and obtain high-order feature space where speech features are converted more easily than in an original acoustic feature space. In our approach, we first train two exclusive CRBMs for the LF and HF, aiming to capture high-order features in an unsupervised manner. Then we train a NN using the projected features produced by two CRBMs. The important advantage of our proposed method is its ability to capture the temporal correlation between adjacent frames of speech and obtain the deep non-linear relationships between the spectral envelope features of LF and HF.

# 2. SPEECH BANDWIDTH EXTENSION BASED ON GMM

Let  $\boldsymbol{x}_t = [x_1, x_2, \cdots, x_m]^T$  and  $\boldsymbol{y}_t = [y_1, y_2, \cdots, y_n]^T$  be the *m* dimensional LF and *n* dimensional HF feature vectors at frame *t*, respectively. The operator  $[\cdot]^T$  denotes matrix transposition. In GMM based method, the joint probability density function of the joint feature space  $\boldsymbol{z}_t = [\boldsymbol{x}_t^T, \boldsymbol{y}_t^T]^T$  is modeled by a GMM

$$f(z_t|\Theta^z) = \sum_{l=1}^{L} a_l^z \mathcal{N}(\boldsymbol{z}_t; \boldsymbol{\mu}_l^z; \boldsymbol{\Sigma}_l^z), \sum_{l=1}^{L} a_l^z = 1$$
(1)

where  $\mathcal{N}(\cdot)$  denotes a Gaussian distribution of joint vector  $\boldsymbol{z}_t$ .  $\Theta^z$  denotes a parameter set of the GMM, which consists

of weights, mean vectors and the covariance matrices for individual mixture components. The total number of mixture components is L.  $a_l^z$ ,  $\mu_l^z$  and  $\Sigma_l^z$  are the weight, mean vector and covariance matrix of the  $l^{th}$  mixture component, and

$$\boldsymbol{\mu}_{l}^{z} = \begin{bmatrix} \boldsymbol{\mu}_{l}^{x} \\ \boldsymbol{\mu}_{l}^{y} \end{bmatrix} \qquad \boldsymbol{\Sigma}_{l}^{z} = \begin{bmatrix} \boldsymbol{\Sigma}_{l}^{xx} & \boldsymbol{\Sigma}_{l}^{xy} \\ \boldsymbol{\Sigma}_{l}^{yx} & \boldsymbol{\Sigma}_{l}^{yy} \end{bmatrix}$$
(2)

 $\mu_l^x$  and  $\mu_l^y$  are the mean vectors of the  $l^{th}$  component for the LF and HF respectively,  $\Sigma_l^{xx}$  and  $\Sigma_l^{yy}$  are the corresponding covariance matrices,  $\Sigma_l^{xy}$  and  $\Sigma_l^{yx}$  are the cross-covariance matrices. The GMM is trained with the expectation maximization (EM) algorithm using the joint vectors as the input feature vectors.

The idea of BWE is to determine a mapping function that can approximate HF features accurately. A mapping function to convert the LF feature vector  $\boldsymbol{x}_t$  to the HF feature vector  $\boldsymbol{y}_t$  is derived based on the conditional probability density of the HF feature vector, given the LF feature vector. When minimizing the mean square error (MMSE) estimation rule is adopted for parameter generation, the mapping function takes the form:

$$\widetilde{\boldsymbol{y}}_{t} = \int_{\Omega_{y}} \boldsymbol{y}_{t} f_{y|x} \left( \boldsymbol{y}_{t} | \boldsymbol{x}_{t} \right)$$

$$= \sum_{l=1}^{L} p_{l} \left( x \right) \left[ \boldsymbol{\mu}_{l}^{y} + \boldsymbol{\Sigma}_{l}^{yx} \left( \boldsymbol{\Sigma}_{l}^{xx} \right)^{-1} \left( \boldsymbol{x}_{t} - \boldsymbol{\mu}_{l}^{x} \right) \right],$$
(3)

where  $p_l(x)$  is the probability of  $x_t$  belonging to the  $l^{th}$  component, i.e

$$p_{l}(x) = \frac{a_{l}\mathcal{N}\left(\boldsymbol{x}_{t}, \boldsymbol{\mu}_{l}^{x}, \boldsymbol{\Sigma}_{l}^{xx}\right)}{\sum_{m=1}^{L} a_{m}\mathcal{N}\left(\boldsymbol{x}_{t}, \boldsymbol{\mu}_{m}^{x}, \boldsymbol{\Sigma}_{m}^{xx}\right)}.$$
(4)

According to Eq.3, the mapping function is a piece-wise linear transformation. The converted HF feature vector  $\tilde{y}_t$  is determined by the current acoustic vector from LF, which means that the GMM based method ignores the time information.

# **3. PROPOSED SPEECH BANDWIDTH EXTENSION**

Our speech bandwidth extension system uses CRBM to capture high-order features and time-related information. We briefly discuss the main idea of CRBM as well as its parameters estimation and give details of our proposed method in this section.

### **3.1. CRBM**

CRBM is the conditional form of RBM proposed by Taylor et al [12]. In this model, short-term temporal structures can be captured by making the visible and hidden units receive additional input from past and future states of visible units dynamically. Given a hidden vector  $\boldsymbol{h}^t = [h_1^t, h_2^t, \cdots , h_J^t]^T$ ,  $h_j^t \in \{0, 1\}$  a visible vector  $\boldsymbol{v}^t = [v_1^t, v_2^t, \cdots , v_I^t]^T$ ,  $v_i^t \in$   $\{0,1\}$  and a conditional vector  $v^{t-r}(r)$  is the number of previous frames from the current frame taken in account, here we choose r = 1 for simplicity) at the current frame t, the conditional probability could be defined as follows:

$$p(\boldsymbol{v}^{t} \mid \boldsymbol{v}^{t-1}) = \frac{1}{Z} \sum_{\boldsymbol{h}^{t}} \exp(-E(\boldsymbol{v}^{t}, \boldsymbol{h}^{t} \mid \boldsymbol{v}^{t-1}))$$
(5)

$$E(\boldsymbol{v}^{t},\boldsymbol{h}^{t} \mid \boldsymbol{v}^{t-1}) = -\boldsymbol{b}^{\mathrm{T}}\boldsymbol{v}^{t} - \boldsymbol{c}^{\mathrm{T}}\boldsymbol{h}^{t} - (\boldsymbol{v}^{t})^{\mathrm{T}}\boldsymbol{W}^{\boldsymbol{v}^{t}\boldsymbol{h}^{t}}\boldsymbol{h}^{t}$$
$$- (\boldsymbol{v}^{t-1})^{\mathrm{T}}\boldsymbol{W}^{\boldsymbol{v}^{t-1}\boldsymbol{v}^{t}}\boldsymbol{v}^{t} \qquad (6)$$
$$- (\boldsymbol{v}^{t-1})^{\mathrm{T}}\boldsymbol{W}^{\boldsymbol{v}^{t-1}\boldsymbol{h}^{t}}\boldsymbol{h}^{t}$$

$$Z = \sum_{\boldsymbol{v}^{t}, \boldsymbol{h}^{t}} \exp(-E(\boldsymbol{v}^{t}, \boldsymbol{h}^{t} \mid \boldsymbol{v}^{t-1}))$$
(7)

where **b** and **c** are a bias vector of visible units and a bias vector of hidden units respectively.  $W^{v^{t}h^{t}}, W^{v^{t-1}v^{t}}$  and  $W^{v^{t-1}v^{t}}$  are the weight matrices between  $v^{t}$  and  $h^{t}, v^{t-1}$  and  $v^{t}, v^{t-1}$  and  $h^{t}$  respectively.

In this model, there are five parameters to be estimated:  $\boldsymbol{b}, \boldsymbol{c}, \boldsymbol{W}^{v^{t-1}h^t}, \boldsymbol{W}^{v^{t}h^t}, \boldsymbol{W}^{v^{t-1}v^t}$ . These parameters are estimated by maximizing the log-likelihood  $\mathcal{L} = \log \prod_t p(\boldsymbol{v}^t | \boldsymbol{v}^{t-1})$ . Differentiating partially with respect to each parameter, we obtain

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{W}_{ij}^{v^t h^t}} = \left\langle v_i^t h_j^t \right\rangle_{data} - \left\langle v_i^t h_j^t \right\rangle_{model} \tag{8}$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{W}_{i'i}^{v^{t-1}v^{t}}} = \left\langle \boldsymbol{v}_{i}^{t}\boldsymbol{v}_{i'}^{t-1} \right\rangle_{data} - \left\langle \boldsymbol{v}_{i}^{t}\boldsymbol{v}_{i'}^{t-1} \right\rangle_{model} \tag{9}$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{W}_{i'j}^{v^{t-1}h^t}} = \left\langle v_{i'}^{t-1}h_j^t \right\rangle_{data} - \left\langle v_{i'}^{t-1}h_j^t \right\rangle_{model} \tag{10}$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{b}_i} = \left\langle \boldsymbol{v}_i^t \right\rangle_{data} - \left\langle \boldsymbol{v}_i^t \right\rangle_{model} \tag{11}$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{c}_j} = \left\langle h_j^t \right\rangle_{data} - \left\langle h_j^t \right\rangle_{model} \tag{12}$$

where  $\langle \cdot \rangle_{data}$  and  $\langle \cdot \rangle_{model}$  indicate the expectations of the input data and the inner model. Because  $\langle \cdot \rangle_{model}$  is extremely expensive to compute exactly, the contrastive divergence approximation to the gradient is used, where  $\langle \cdot \rangle_{model}$  is replaced by running the Gibbs sampler initialized at the data for one full step [13].

Once the parameters are estimated, the conditional probability of  $h^t$  given  $v^t$  and  $v^{t-1}$  and the conditional probability of  $v^t$  given  $h^t$  and  $v^{t-1}$  are respectively written as:

$$p\left(\boldsymbol{h}_{j}^{t}=1 \mid \boldsymbol{v}^{t}, \boldsymbol{v}^{t-1}\right) = \sigma\left(c_{j} + (\boldsymbol{v}^{t})^{\mathrm{T}} \boldsymbol{W}_{:j}^{\boldsymbol{v}^{t}h^{t}} + (\boldsymbol{v}^{t-1})^{\mathrm{T}} \boldsymbol{W}_{:j}^{\boldsymbol{v}^{t-1}h^{t}}\right)$$

$$(13)$$

$$p\left(\boldsymbol{v}_{i}^{t}=1 \mid \boldsymbol{h}^{t}, \boldsymbol{v}^{t-1}\right) = \sigma\left(b_{i} + (\boldsymbol{h}^{t})^{\mathrm{T}} (\boldsymbol{W}_{i:}^{\boldsymbol{v}^{t}h^{t}})^{\mathrm{T}} + (\boldsymbol{v}^{t-1})^{\mathrm{T}} \boldsymbol{W}_{:j}^{\boldsymbol{v}^{t-1}v^{t}}\right)$$

$$(14)$$

where  $W_{i:}$  and  $W_{:j}$  denote the column vector and the row vector in W respectively, and  $\sigma$  indicates a sigmoid function; i.e.  $\sigma(x) = \frac{1}{1+e^x}$ .



**Fig. 1**. (a) CRBMs for the low frequency (below) and high frequency (above) (b)our proposed speech bandwidth extension architecture combining two CRBMs and a NN

#### 3.2. Speech bandwidth extension using CRBM

Figure 1 shows an overview of our proposed speech BWE system. In our approach, two CRBMs (one for the low frequency and the other for the high frequency) are adopted to describe the distribution of  $x^t$ ,  $x^{t-1}$  and  $y^t$ ,  $y^{t-1}$  respectively.  $x^t$  and  $x^{t-1}$  are spectral envelop features of the LF at frame t and t-1 respectively.  $y^t$  and  $y^{t-1}$  are spectral envelop features of the HF at frame t and t-1 respectively. Then a NN is employed to model the joint distribution of  $h_x^t$  and  $h_y^t$ .  $h_x^t$  and  $h_y^t$  are the hidden variables extracted from two CRBMs respectively. The extracted hidden variables can be considered as the high-order binary representation of the raw spectral envelopes. Therefore, the parameter set of our proposed method is given by

$$\boldsymbol{\Theta} = \{\boldsymbol{\theta}_x, \boldsymbol{\theta}_y, \boldsymbol{\theta}_n\} \tag{15}$$

where  $\theta_x = \left\{ W^{x^t x^{t-1}}, W^{x^{t-1}h^t}, W^{x^{t}h^t}, b_x, c_x \right\}$  are the parameters of the CRBM using the training spectral envelopes of LF,  $\theta_y = \left\{ W^{y^t y^{t-1}}, W^{y^{t-1}h^t}, W^{y^{t}h^t}, b_y, c_y \right\}$  are the parameter of the CRBM using the training spectral envelopes of HF,  $\theta_n = \{ W^1, \dots, W^L, d^1, \dots, d^L \}$   $(L-1) = \{0, 1, 2\}$  is the number of the hidden layers,  $d^l$  is the bias vector of the  $l^{th}$  layer,  $W^l$  is the weight matrix from the  $(l-1)^{th}$  layer to the  $l^{th}$  layer) are the parameters of the NN trained using the extracted hidden variables. The weight matrices  $W^{x^t x^{t-1}}, W^{x^{t-1}h^t}, W^{y^t y^{t-1}}, W^{y^{t-1}h^t}$  can absorb time-related information.

At the training phase of the NN, the input vectors and the output vectors are written as

$$\boldsymbol{h}_{x}^{t} = \sigma \left( c_{x} + \boldsymbol{W}^{x^{t}h^{t}} \boldsymbol{x}^{t} + \boldsymbol{W}^{x^{t-1}h^{t}} \boldsymbol{x}^{t-1} \right)$$
(16)

$$\boldsymbol{h}_{y}^{t} = \sigma \left( c_{y} + \boldsymbol{W}^{y^{t}h^{t}} \boldsymbol{y}^{t} + \boldsymbol{W}^{y^{t-1}h^{t}} \boldsymbol{y}^{t-1} \right)$$
(17)

Weight parameters  $\theta_n$  of NN are estimated by minimizing the error between the output  $\tilde{h}_y^t$  of NN and the target vector  $h_y^t$ . Once the parameters of NN are obtained, the input vector  $h_x^t$  can be converted to:

$$\widetilde{\boldsymbol{h}}_{y}^{t} = \boldsymbol{o}^{l} = \sigma\left(\boldsymbol{s}^{l}\right), l \ge 1$$
(18)

where  $s^{l} = W^{l}o^{l-1} + d^{l}$  and  $o^{l-1} = \sigma\left(s^{l-1}\right), o^{0} = h_{x}^{t}$ .

To map the output of the NN to the spectral envelope features of the HF, we can use backward inference of CRBM using Eq.14,

$$p\left(\boldsymbol{y}^{t} \mid \widetilde{\boldsymbol{h}}_{y}^{t}, \boldsymbol{y}^{t-1}\right) = \sigma\left(b_{y} + (\boldsymbol{W}^{y^{t}h^{t}})^{\mathrm{T}}\widetilde{\boldsymbol{h}}_{y}^{t} + \boldsymbol{W}^{y^{t-1}y^{t}}\boldsymbol{y}^{t-1}\right)$$
(19)

According to Eq.18 and Eq.19, the mapping function of our method from a LF's spectral envelope feature vector  $x^t$ to a HF's spectral envelope feature vector  $y^t$  at frame t, given the previous feature vectors  $x^{t-1}$  and  $y^{t-1}$  is given as:

$$\boldsymbol{y}^{t} = \sigma \left( b_{y} + (\boldsymbol{W}^{y^{t}h^{t}})^{\mathrm{T}} \sigma \left( \boldsymbol{s}^{l} \right) + \boldsymbol{W}^{y^{t-1}y^{t}} \boldsymbol{y}^{t-1} \right)$$
(20)

As Eq.20 indicates, the mapping function is a composite function of multiple different non-linear functions. We need a current acoustic vector from LF, and previous vectors from both LF and HF to reconstruct the HF's current acoustic vector. Thus, additional information from adjacent signal frames can be captured and deep non-linear relationships between spectral envelope features of LF and HF can be discovered by our proposed method.

#### 4. EXPERIMENTS

### **4.1. setup**

We conducted speech bandwidth extension experiments using two Chinese speech databases, comparing our method (CRB-M) with the traditional GMM-based method (GMM) and our previous work (RBM). The first Chinese speech database is from the NTT Advanced Technology Corporation (NTT-AT). The second database is from Ericsson and Beijing Institute of Technology (EBIT). The data in the two databases is sampled at a 16-kHz sampling rate with 16-bits resolution. Each utterance in the two databases lasts 8s. A high-pass filtering supplied the high frequency signal. The low frequency signal resulted from a 0.3 to 3.4 kHz band-pass filtering followed by a down-sampling and up-sampling with a factor 2. We used 64 utterances randomly selected from all speech sound classes in NTT and 64 utterances in EBIT as our training set. The test set consisted of the 32 utterances in NTT and the 32 utterances in EBIT that were not included in the training data.

A GMM with 128 components was trained for the baseline system. The 16-order and 10-order line spectral frequencies (LSFs) [14] were adopted as the spectral envelope features for the LF and HF respectively. LSFs are first normalized to have zero mean and unit variance, and then converted to binary using a sigmoid function before feeding them into CRBMs. The frame size and the frame shift for calculating spectral envelopes was set to 20ms and 10ms respectively. We investigated on three neural network architectures (arc.1: no hidden layer, arc.2: 1 hidden layer and 128 hidden units in the hidden layer, arc.3: 2 hidden layers and 128 hidden units in each hidden layer) for the following experiments. For each architecture, the contrastive divergence (CD) learning with 10-step Gibbs sampling was employed to train two CRBMs. The stochastic batch gradient descent algorithm was adopted to update the model parameters. The size of each mini-batch was set to 64. The learning rate and momentum were set to 0.0001 and 0.9. The number of epochs of CRBMs and NNs were set to 500 and 300 respectively. The number of hidden units of two CRBMs was fixed to 64. We used the LF excitation signal as HF excitation signal. In order to adjust the power of the extended HF excitation signal, we builded a one-to-one codebook from the LSFs of LF to gains between the HF signal and the synthesized signal that the LF excitation signal filtered through the high frequency synthesis filter.

To assess the overall quality of reconstructed speech, we conducted objective evaluation and subjective evaluation. For the objective evaluation, we used root mean square log-spectral distortion (RMS-LSD) to measure how close the reconstructed speech is to the original one. We calculated the RMS-LSD for each frame and averaged the RMS-LSD values. For the subjective evaluation, mean opinion score (MOS) listening tests were conducted.

# 4.2. objective evaluation

We measured the RMS-LSD in the missing high frequency (4-8 kHz). The definition of RMS-LSD is as follows,

$$D = \sqrt{\frac{1}{N} \sum_{n=1}^{N} \frac{1}{w_2 - w_1} \int_{w_1}^{w_2} \left[ 20 \log \left( \mathbf{G}_c \frac{\mathbf{A}_n(w)}{\hat{\mathbf{A}}_n(w)} \right) \right]^2 dw} \quad (21)$$

$$G_{c} = \frac{1}{w_{2} - w_{1}} \int_{w_{1}}^{w_{2}} \left[ 20 \log \left( \frac{\hat{A}_{n}(w)}{A_{n}(w)} \right) \right] dw,$$
(22)

where  $A_n(w)$  and  $\widehat{A}_n(w)$  denote the original and the reconstructed power spectrum of  $n^{th}$  frame of high frequency respectively,  $w_1, w_2$  are the lower and higher bound of the missing high frequency, compensating gain factor  $G_c$  has the effect of removing the mean difference between the two log envelopes. The smaller the value of RMS-LSD is, the closer the reconstructed high frequency is to the original high frequency, the better the speech quality is. The RMS-LSD results are shown in Figure 2. Experimental results show that CRB-M based method has lower log spectral distortion than GMM and RBM based method.

# 4.3. subjective evaluation

For MOS tests, 12 participants were asked to listen to the original speech signals in the test set and the reconstructed



Fig. 2. Results of Objective evaluation



Fig. 3. Results of Subjective evaluation

speech signals for each method, and to select how close the reconstructed speech is to the original one on a 5-point scale (1-bad,2-poor,3-fair,4-good,5-excellent). Figure 2 and figure 3 summarize the objective and subjective experimental results. As shown in these figures, "CRBM" based method outperforms other conversional methods (GMM, RBM) in both criteria. The reason for the improvement is attributed to the fact that CRBM can capture time-related information and nonlinear relationship between spectral envelope features of LF and HF.

# **5. CONCLUSION**

In this paper, we proposed a new speech bandwidth extension method that combined two conditional restricted Boltzmann machines and a neural network to construct a non-linear relationships between the low frequency's and high frequency's spectral envelope features and capture the time-related information. The objective evaluation and subjective evaluation showed improvement of the proposed method when compared with the conventional GMM based method and other network based method.

# 6. ACKNOWLEDGEMENT

This work was supported by National Natural Science Foundation of China (Grant No: 61175067).

# 7. REFERENCES

- Kun-Youl Park and Hyung Soon Kim, "Narrowband to wideband conversion of speech using gmm based transformation," in Acoustics, Speech, and Signal Processing, 2000. ICASSP'00. Proceedings. 2000 IEEE International Conference on. IEEE, 2000, vol. 3, pp. 1843– 1846.
- [2] Xin Liu, Chang-Chung Bao, Mao-shen Jia, and Yongtao Sha, "A harmonic bandwidth extension based on gaussian mixture model," in *Signal Processing (ICSP)*, 2010 IEEE 10th International Conference on. IEEE, 2010, pp. 474–477.
- [3] Hannu Pulakka, Ulpu Remes, Kalle Palomäki, Mikko Kurimo, and Paavo Alku, "Speech bandwidth extension using gaussian mixture model-based estimation of the highband mel spectrum," in Acoustics, Speech and Signal Processing (ICASSP), 2011 IEEE International Conference on. IEEE, 2011, pp. 5100–5103.
- [4] Murali D Mohan, Dileep B Karpur, Manoj Narayan, and J Kishore, "Artificial bandwidth extension of narrowband speech using gaussian mixture model," in *Communications and Signal Processing (ICCSP)*, 2011 International Conference on. IEEE, 2011, pp. 410–412.
- [5] Hannu Pulakka, Ulpu Remes, Santeri Yrttiaho, Kalle Palomäki, Mikko Kurimo, and Paavo Alku, "Bandwidth extension of telephone speech to low frequencies using sinusoidal synthesis and a gaussian mixture model," *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 20, no. 8, pp. 2219–2231, 2012.
- [6] Kehuang Li and Chin-Hui Lee, "A deep neural network approach to speech bandwidth expansion," in Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on. IEEE, 2015, pp. 4395–4399.
- [7] Yingxue Wang, Shenghui Zhao, Wenbo Liu, Ming Li, and Jingming Kuang, "Speech bandwidth expansion based on deep neural networks," in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.
- [8] Amr H Nour-Eldin, Turaj Z Shabestary, and Peter Kabal, "The effect of memory inclusion on mutual information between speech frequency bands," in Acoustics, Speech and Signal Processing, 2006. ICASSP 2006 Proceedings. 2006 IEEE International Conference on. IEEE, 2006, vol. 3, pp. III–III.
- [9] Amr H Nour-Eldin and Peter Kabal, "Combining frontend-based memory with mfcc features for bandwidth extension of narrowband speech," in Acoustics, Speech and Signal Processing, 2009. ICASSP 2009.

*IEEE International Conference on*. IEEE, 2009, pp. 4001–4004.

- [10] Amr H Nour-Eldin and Peter Kabal, "Memory-based approximation of the gaussian mixture model framework for bandwidth extension of narrowband speech.," in *INTERSPEECH*, 2011, pp. 1185–1188.
- [11] Peter Jax and Peter Vary, "Artificial bandwidth extension of speech signals using mmse estimation based on a hidden markov model," in *Acoustics, Speech, and Signal Processing, 2003. Proceedings.(ICASSP'03). 2003 IEEE International Conference on.* IEEE, 2003, vol. 1, pp. I–680.
- [12] Graham W Taylor, Geoffrey E Hinton, and Sam T Roweis, "Modeling human motion using binary latent variables," in *Advances in neural information processing systems*, 2006, pp. 1345–1352.
- [13] Geoffrey E Hinton, "Training products of experts by minimizing contrastive divergence," *Neural computation*, vol. 14, no. 8, pp. 1771–1800, 2002.
- [14] Samir Chennoukh, A Gerrits, G Miet, and R Sluijter, "Speech enhancement via frequency bandwidth extension using line spectral frequencies," in Acoustics, Speech, and Signal Processing, 2001. Proceedings.(ICASSP'01). 2001 IEEE International Conference on. IEEE, 2001, vol. 1, pp. 665–668.