

Acoustic Compensation Method for Accepting Different Recording Devices in Body-Conducted Voice Conversion

Daisuke Deguchi*, Hironori Doi*, Tomoki Toda *, Hiroshi Saruwatari* and Kiyohiro Shikano*

* Graduate School of Information Science, Nara Institute of Science and Technology, Ikoma Nara 630-0192 Japan
E-mail:{daisuke-d, hironori-d, tomoki, sawatari, shikano}@is.naist.jp Tel:+81-743-72-5288

Abstract—This paper presents an acoustic compensation method in body-conducted speech conversion that automatically compensates for acoustic differences caused by changes in recording conditions. An enhancement process for body-conducted speech recorded with a Non-Audible Murmur (NAM) microphone has successfully applied a statistical voice conversion technique. Speech waveforms are generated from acoustic parameters of normal speech estimated from those of body-conducted speech with a conversion model previously trained using stereo data of those two types of speech. This framework suffers from mismatched conditions between training and conversion processes. To alleviate this issue, an unsupervised acoustic compensation method based on constrained maximum likelihood linear regression (CMLLR) has been proposed and its effectiveness has been reported in the compensation of acoustic differences caused by attachment location of the NAM microphone. This paper further applies the CMLLR-based acoustic compensation method to the compensation of acoustic differences caused by different recording devices and evaluates its effectiveness. The experimental results demonstrate that the proposed method effectively reduces quality degradation and converted speech caused by differences in recording devices as well as attachment location of the NAM microphone.

I. INTRODUCTION

Recently cellular phones have enabled us to communicate with each other conveniently. However, noisy environments such as a crowd make it difficult to smoothly convey speech. To address this issue, the use of body-conducted speech for speech communication [1], [2] has been proposed. A Non-Audible Murmur (NAM) microphone[2], that is a body-conductive microphone, is placed on the neck below the ear and detects various types of speech such as non-audible murmurs, whispers, and normal speech through the soft tissue of the head. This approach is robust against external noise due to the noise-proof structure of body-conductive microphones. However, body-conducted speech causes severe quality degradation due to essential mechanisms of body conduction such as the lack of radiation characteristics from lips and the influence of the low-pass characteristics of soft tissue.

To improve the quality of body-conducted speech, body-conducted voice conversion has been proposed [3]. This technique is based on statistical voice conversion techniques [4], [5]. A conversion model is trained in advance using a parallel data set consisting of utterance-pairs of body-conducted speech and normal speech (i.e., air-conducted speech) uttered by the same speaker. A Gaussian mixture model (GMM) of joint probability density of speech parameters of those two types of speech is effectively used as the conversion model.

The trained model is capable of converting body-conducted speech parameters into normal speech parameters without the use of linguistic information.

One weakness of body-conducted voice conversion is that severe quality degradation of the converted speech is caused by mismatched acoustic conditions between training and conversion. The recording of body-conducted speech with a NAM microphone is sensitive to various conditions such as the attachment location of the NAM microphone, settings of the amplifier, the type of NAM microphone, and so on. In the practical use of a NAM microphone, it is almost impossible to keep these conditions consistent. Moreover, the NAM microphone is still under development and various types of recording devices would may be developed. Therefore, it techniques for compensating for acoustic differences of body-conducted speech caused by changes of various recording conditions must be developed.

In our previous work [6], an acoustic compensation method based on constrained maximum likelihood linear regression (CMLLR) [7] has been proposed to compensate for acoustic differences caused by a change in the attachment location of the NAM microphone. This method estimates feature-space transforms for compensating for the acoustic differences in a completely unsupervised manner using only the acoustics of body-conducted speech. It is worthwhile to investigate the effectiveness of this method in various recording conditions beyond simply changing the location of the NAM microphone.

In this paper, we apply the CMLLR-based acoustic compensation method to the compensation of acoustic differences caused by different recording devices. Three different types of NAM microphone and their own amplifiers are used for recording body-conducted speech. The experimental results demonstrate that the proposed compensation method is effective for compensating the acoustic differences caused by the different recording devices as well as the different attachment location of NAM microphone.

II. BODY-CONDUCTED SPEECH

Body-conducted speech has a different spectral structure from that of air-conducted speech [6]. In particular, higher frequency components of body-conducted speech are severely attenuated. Consequently body-conducted speech usually sounds muffled.

In this paper, three types of different NAM microphone and their own amplifiers shown in Figure 1 are used for recording

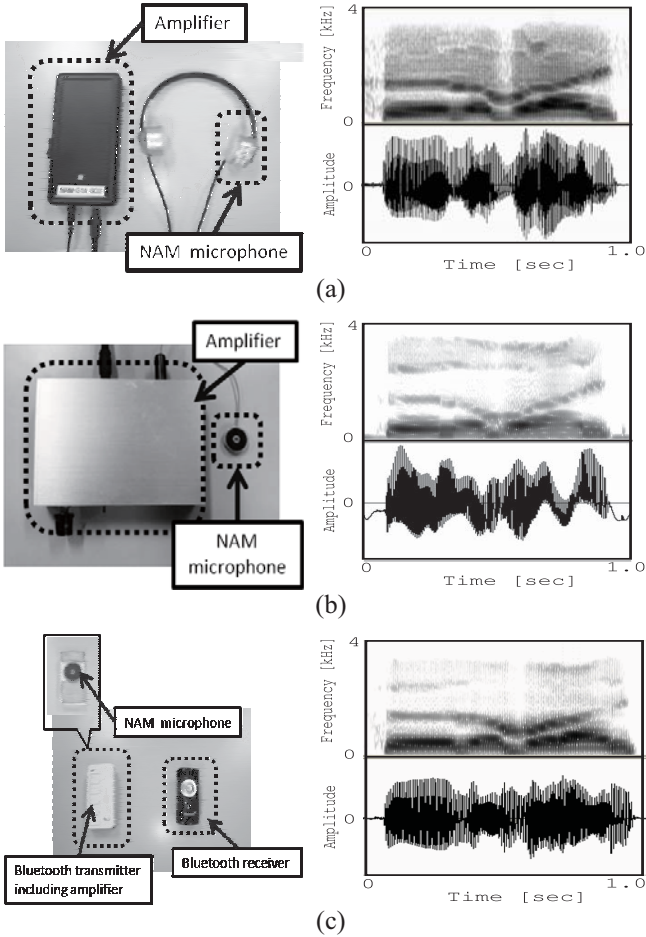


Fig. 1. Three types of NAM microphone, (a) wired-type with neckband, (b) wired-type without neckband, and (c) wireless-type, and an example of spectrogram of body-conducted speech recorded with each NAM microphone.

body-conducted speech. The use of different types of NAM microphone and their amplifiers causes noticeable changes in the acoustic features of body-conducted speech. These acoustic changes are caused by differences in various conditions such as the attachment location of the NAM microphone and the way of attaching the NAM microphone on the body as well as differences the recording device.

III. BODY-CONDUCTED SPEECH CONVERSION

A. Feature Extraction of Body-Conducted Speech

As a source feature, we employ a spectral segment vector. Let \mathbf{x}_t is a mel-cepstral vector at a frame t . We construct a concatenated vector $\mathbf{c}_t = [\mathbf{x}_{t-n}^\top \cdots \mathbf{x}_t^\top \cdots \mathbf{x}_{t+n}^\top]^\top$ over the current $\pm n$ frames, where the symbol \top indicates transpose. And then, the spectral segment vector \mathbf{X}_t at frame t is extracted by PCA as follows:

$$\mathbf{X}_t = \mathbf{D}\mathbf{c}_t - \mathbf{d}, \quad (1)$$

where \mathbf{D} is the transformation matrix of PCA, and $\mathbf{d} = \mathbf{D}\bar{\mathbf{c}}$. The vector $\bar{\mathbf{c}}$ is the mean vector of \mathbf{c}_t within all training data for PCA.

As a target feature, we employ the joint static and dynamic feature vector $\mathbf{Y}_t = [\mathbf{y}_t^\top \ \Delta\mathbf{y}_t^\top]^\top$, where \mathbf{y}_t is the static mel-cepstral vector, and $\Delta\mathbf{y}_t$ is the delta mel-cepstral vector of the target data at frame t .

B. Feature Conversion Based on Maximum Likelihood Estimation of Parameter Trajectory [5]

The joint probability density of the source and target feature vectors is modeled by a GMM as follows:

$$P(\mathbf{Z}_t|\boldsymbol{\lambda}) = \sum_{m=1}^M w_m \mathcal{N}(\mathbf{Z}_t; \boldsymbol{\mu}_m^{(Z)}, \boldsymbol{\Sigma}_m^{(ZZ)}), \quad (2)$$

where \mathbf{Z}_t is the joint feature vector $\mathbf{Z}_t = [\mathbf{X}_t^\top \ \mathbf{Y}_t^\top]^\top$. The symbol $\mathcal{N}()$ indicates the normal distribution. The number of mixture components is M . $\boldsymbol{\lambda}$ is the model parameter set including w_m , $\boldsymbol{\mu}_m^{(Z)}$, and $\boldsymbol{\Sigma}_m^{(ZZ)}$, which are the weight, mean vector, and covariance matrix of the m -th mixture component, respectively. $\boldsymbol{\mu}_m^{(Z)}$ and $\boldsymbol{\Sigma}_m^{(ZZ)}$ are represented by

$$\boldsymbol{\mu}_m^{(Z)} = \begin{bmatrix} \boldsymbol{\mu}_m^{(X)} \\ \boldsymbol{\mu}_m^{(Y)} \end{bmatrix}, \quad (3)$$

$$\boldsymbol{\Sigma}_m^{(ZZ)} = \begin{bmatrix} \boldsymbol{\Sigma}_m^{(XX)} & \boldsymbol{\Sigma}_m^{(XY)} \\ \boldsymbol{\Sigma}_m^{(YX)} & \boldsymbol{\Sigma}_m^{(YY)} \end{bmatrix}, \quad (4)$$

where the matrices $\boldsymbol{\Sigma}_m^{(XX)}$ and $\boldsymbol{\Sigma}_m^{(YY)}$ are the covariance matrix of the m -th mixture component of the source and that of the target, respectively. The matrices $\boldsymbol{\Sigma}_m^{(XY)}$ and $\boldsymbol{\Sigma}_m^{(YX)}$ are the cross covariance matrices of the m -th mixture component between the source and target. These covariance matrices are completely full.

Let $\mathbf{X} = [\mathbf{X}_1^\top; \cdots; \mathbf{X}_T^\top]^\top$ and $\mathbf{Y} = [\mathbf{Y}_1^\top; \cdots; \mathbf{Y}_T^\top]^\top$ be time sequences of the source and the target features, respectively. The converted static feature vector sequence is determined so that the following approximated conditional probability density is maximized.

$$P(\mathbf{Y}|\mathbf{X}, \boldsymbol{\lambda}) \simeq P(\mathbf{m}|\mathbf{X}, \boldsymbol{\lambda})P(\mathbf{Y}|\mathbf{X}, \mathbf{m}, \boldsymbol{\lambda}), \quad (5)$$

where $\mathbf{m} = \{m_1, \cdots, m_T\}$ is a mixture component sequence. First, suboptimum mixture component sequence $\hat{\mathbf{m}}$ is determined by

$$\hat{\mathbf{m}} = \arg \max_{\mathbf{m}} P(\mathbf{m}|\mathbf{X}, \boldsymbol{\lambda}). \quad (6)$$

And then, the converted static feature vector $\hat{\mathbf{y}}$ is obtained by

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} P(\mathbf{Y}|\mathbf{X}, \hat{\mathbf{m}}, \boldsymbol{\lambda}), \quad (7)$$

subject to $\mathbf{Y} = \mathbf{E}\mathbf{y}$,

where \mathbf{E} is a window matrix to extend the static feature vector sequence into the joint static and dynamic feature vector sequence. Furthermore, the quality of the converted voice is dramatically improved by considering the global variance of the converted feature [5].

IV. UNSUPERVISED ACOUSTIC COMPENSATION METHOD BASED ON CMLLR

To compensate for the acoustic differences caused by changes in recording conditions, the CMLLR-based unsupervised adaptation method [6] is adopted. In this method, to effectively reduce the mismatch between the model and the adaptation data, the CMLLR transform is estimated using only body-conducted speech and the previously trained GMM of joint probability density.

We apply the CMLLR transformation to the source features to compensate for their acoustic variations. The transformed source feature vector is given by

$$\hat{\mathbf{X}}_t = \mathbf{A}\mathbf{X}_t + \mathbf{b} = \mathbf{W}\boldsymbol{\xi}(t), \quad (8)$$

where \mathbf{W} is the extended transform, $[\mathbf{b} \ \mathbf{A}]$, and $\boldsymbol{\xi}(t)$ is the extended source feature vector, $[1 \ \mathbf{X}_t^\top]^\top$.

To perform unsupervised compensation, the CMLLR transform is estimated so that the likelihood of the marginal distribution for the adaptation source data $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_T]$ is maximized as follows:

$$\hat{\mathbf{W}} = \arg \max_{\mathbf{W}} \prod_{t=1}^T \int P(\mathbf{X}_t, \mathbf{Y}_t | \mathbf{W}, \lambda) d\mathbf{Y}_t. \quad (9)$$

Because the probability density is modeled by a GMM, EM algorithm is employed. In each M-step, an iterative row-by-row update is performed for determining the updated transformation matrix [7], [8]. The i -th row vector of $\hat{\mathbf{W}}$ is given by

$$\mathbf{w}_i = (\alpha \mathbf{c}_i + \mathbf{k}^{(i)}) \mathbf{G}^{(ii)^{-1}}, \quad (10)$$

where \mathbf{c}_i is the extended cofactor row vector of \mathbf{A} and α is found by solving a quadratic equation [7]. $\mathbf{k}^{(i)}$ and $\mathbf{G}^{(i)}$ are given by

$$\mathbf{G}^{(ij)} = \sum_{m=1}^M p_m(i, j) \sum_{t=1}^T \gamma_m(t) \boldsymbol{\xi}_m(t) \boldsymbol{\xi}_m(t)^\top, \quad (11)$$

$$\mathbf{k}^{(i)} = \sum_{m=1}^M p_m(i) \boldsymbol{\mu}_m \sum_{t=1}^T \gamma_m(t) \boldsymbol{\xi}_m(t)^\top - \sum_{j=1, j \neq i}^d \mathbf{w}_j \mathbf{G}^{(ij)}, \quad (12)$$

where $p_m(i)$ and $p_m(i, j)$ are the i -th row vector and the (i, j) -th element of the inverse covariance matrix $\boldsymbol{\Sigma}_m^{(XX)^{-1}}$, respectively. $\gamma_m(t)$ is the posterior probability of the m -th mixture component given \mathbf{X}_t .

When applying the CMLLR transformation in the model-space, the adapted model parameters are given by

$$\hat{\boldsymbol{\mu}}_m^{(Z)} = \begin{bmatrix} \mathbf{A}' \boldsymbol{\mu}_m^{(X)} - \mathbf{b}' \\ \boldsymbol{\mu}_m^{(Y)} \end{bmatrix}, \quad (13)$$

$$\hat{\boldsymbol{\Sigma}}_m^{(Z)} = \begin{bmatrix} \mathbf{A}' \boldsymbol{\Sigma}_m^{(XX)} \mathbf{A}'^\top & \mathbf{A}' \boldsymbol{\Sigma}_m^{(XY)} \\ \boldsymbol{\Sigma}_m^{(YX)} \mathbf{A}'^\top & \boldsymbol{\Sigma}_m^{(YY)} \end{bmatrix}, \quad (14)$$

where $\hat{\boldsymbol{\mu}}_m^{(Z)}$ and $\hat{\boldsymbol{\Sigma}}_m^{(Z)}$ are the adapted mean vector and covariance matrix of the m -th mixture component, respectively. Note that $\mathbf{A}' = \mathbf{A}^{-1}$ and $\mathbf{b}' = \mathbf{A}'\mathbf{b}$. A global transform is used in this paper.

V. EXPERIMENTAL EVALUATIONS

A. Experimental Conditions

Body-conducted speech and air-conducted speech uttered by one Japanese male speaker were simultaneously recorded with a NAM microphone and a headset microphone. Using each type of NAM microphone and its amplifier shown in Figure 1, 100 phoneme-balanced sentences were recorded. Consequently, three sets of parallel data of body-conducted speech (a, b, and c) and air-conducted speech (A, B, and C) were developed. Because each data set was recorded in a different recording session, the voice quality of the recorded speech samples were slightly different from each other. Fifty

TABLE I
SEVERAL CONDITIONS EVALUATED IN EXPERIMENTS. BODY-CONDUCTED SPEECH (a, b, AND c) AND AIR-CONDUCTED SPEECH (A, B, AND C) ARE SIMULTANEOUSLY RECORDED IN THREE DIFFERENT CONDITIONS USING DIFFERENT RECORDING DEVICES ((a), (b), AND (c) IN FIGURE 1)

Condition	Parallel data set in training	Body-conducted speech in conversion
Fully-matched	a-A	a
	b-B	b
	c-C	c
Mismatched	a-A	b or c
	b-B	a or c
	c-C	a or b
Adapted	a-A	b or c (adapted to a)
	b-B	a or c (adapted to b)
	c-C	a or b (adapted to c)
Matched	a-B	a
	a-C	a
	b-A	b
	b-C	b
	c-A	c
	c-B	c

sentences were used for training or adaptation and the remaining fifty sentences were used for conversion in each data set.

By changing data sets used in training and in conversion, various conditions were evaluated as shown in Table I. Three conversion models were independently built using three parallel data sets (a-A, b-B, and c-C). In the "Fully-matched" condition, the same data sets used in training were also used in conversion. Therefore, not only the recording conditions but also the recording sessions were consistent between training and conversion. In this condition, three sets of converted speech were generated. In the "Mismatched" condition, different data sets were used in conversion. Consequently, six sets of converted speech were generated. Acoustic differences caused by changes in various recording conditions such as different recording devices and different attachment locations of the NAM microphone were observed in this condition. In the "Adapted" condition, the proposed method was used to compensate for these differences. Again, six sets of converted speech were generated. Moreover, to show ideal results in the proposed compensation framework, the other conversion models were also trained using parallel data sets consisting of body-conducted speech and air-conducted speech recorded in different sessions. The other models were then used to convert the body-conducted speech recorded in the same session as used in training. This was called the "Matched" condition and again six sets of converted speech were generated. In this condition, there were some differences of voice quality in target air-conducted speech between training and conversion due to the different recording sessions.

The 0-th through 16-th mel-cepstral coefficients were adopted as a spectral parameter. As the input feature, a 34-dimensional spectral segment feature vector was extracted by PCA from a concatenated vector consisting of current and ± 4 frames of mel-cepstrum vectors. The sampling frequency was set to 8 kHz. The number of mixture components of GMM was set to 32. The frame shift was set to 5 ms.

In the objective evaluation, we evaluated spectral conversion accuracy with mel-cepstral distortion calculated from the first through 16-th mel-cepstral coefficients between the converted and target mel-cepstra. Converted speech samples in four con-

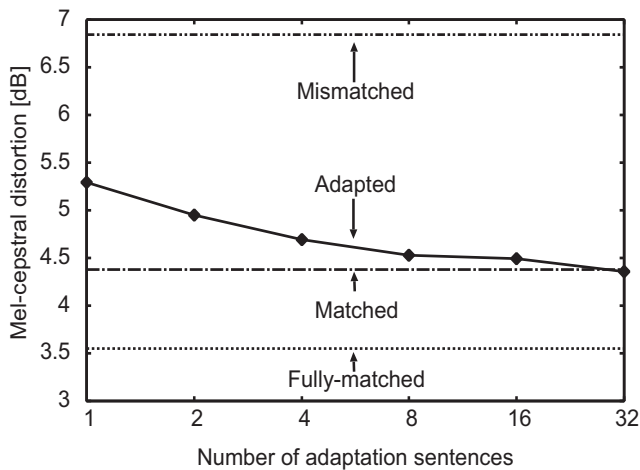


Fig. 2. Mel-cepstral distortion as a function of the number of adaptation sentences.

ditions as shown in Table I were evaluated. In the "Adapted" condition, the number of adaptation sentences was set to 1, 2, 4, 8, 16, or 32, which were used for estimating the CMLLR transform.

In subjective evaluation, we conducted an opinion test of speech quality. An opinion score was set to a 5-point scale (1: bad - 5: excellent). Seven listeners evaluated the sound quality of six types of converted speech: under the "Mismatched" condition; under the "Adapted" condition when using 2, 8, or 16 adaptation sentences; under the "Matched" condition, and under the "Fully-matched" condition. For each listener, 18 sentences were randomly selected from the test set and 108 samples of the converted speech were evaluated in total.

B. Experimental Results

Figure 2 shows the results of the objective evaluation. A large mel-cepstral distortion is observed in the converted speech under the "Mismatched" condition. This result shows that body-conducted voice conversion is very sensitive to the recording conditions. This degradation of conversion accuracy is effectively alleviated by the proposed acoustic compensation method, i.e., the "Adapted" condition. Even if only one sentence is used for adaptation, mel-cepstral distortion significantly decreases compared with the "Mismatched" condition. Mel-cepstral distortion gradually decreases according to an increase in the amount of adaptation data up to around 8 sentences which is close to an ideal result shown in the "Matched" condition. Note that the difference between "Matched" and "Fully-matched" is caused by the voice quality difference between different recording sessions.

Figure 3 shows the results of the subjective evaluation. We can see that changes of the recording devices and the attachment location of NAM microphone cause significant quality degradation. Our proposed compensation method effectively improves the converted speech quality. The use of only 8 adaptation sentences makes the converted speech quality equivalent to that in the "Matched" condition. We can also see a significant quality difference between the "Matched" and "Fully-matched" conditions. These results are consistent with those observed in the objective evaluation.

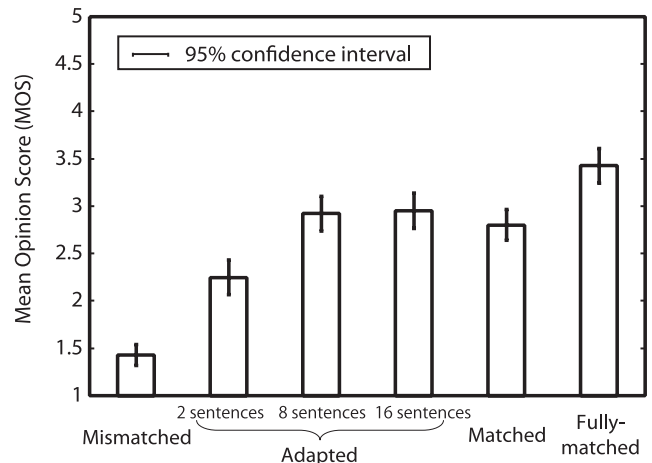


Fig. 3. Result of opinion test on speech quality.

VI. CONCLUSION

This paper described an unsupervised acoustic compensation method based on CMLLR in body-conducted voice conversion to compensate for acoustic differences caused by different recording conditions. The proposed method has been applied to body-conducted voice conversion in mismatched conditions caused by the use of different recording devices as well as a change in the attachment location of a NAM microphone. The experimental results from both objective and subjective tests have demonstrated that the proposed compensation method is capable of effectively alleviating performance degradation in body-conducted voice conversion under mismatched conditions using only a few adaptation sentences.

VII. ACKNOWLEDGMENT

This work was supported in part by MIC SCOPE.

REFERENCES

- [1] Y. Zheng, Z. Liu, M. Sinclair, J. Droppo, L. Deng, A. Acero, and X. Huang, "Air- and Bone-Conductive Integrated Microphones for Robust Speech Detection and Enhancement", *Proc. ASRU*, pp. 249 – 254, 2003
- [2] Y. Nakajima, H. Kashioka, N. Campbell and K. Shikano, "Non-Audible Murmur (NAM) Recognition", *IEICE Trans. Information and Systems*, Vol. E89-D, No. 1, pp. 1 – 8, 2006.
- [3] T. Toda, K. Nakamura, H. Sekimoto and K. Shikano, "Voice Conversion for Various Types of Body Transmitted Speech", *Proc. ICASSP*, pp. 3061 – 3064, 2009.
- [4] Y. Stylianou, O. Cappé and E. Moulines, "Continuous Probabilistic Transform for Voice Conversion", in *Proc. IEEE Trans. SAP*, Vol. 6, No. 2, pp. 131 – 142, 1998.
- [5] T. Toda, A. W. Black, and K. Tokuda, "Voice Conversion Based on Maximum Likelihood Estimation of Spectral Parameter Trajectory", *Trans. ASLP*, Vol. 15, No. 8, pp. 2222 – 2235, 2007.
- [6] D. Miyamoto, K. Nakamura, T. Toda, H. Saruwatari, and K. Shikano, "Acoustic Compensation Methods for Body Transmitted Speech Conversion", *Proc. ICASSP*, pp. 3901 – 3904, 2009.
- [7] M. J. F. Gales, "Maximum Likelihood Linear Transformations for HMM-based Speech Recognition", *Computer Speech and Language*, vol. 12, no. 2, pp. 75 – 98, 1998.
- [8] K. C. Sim and M. J. F. Gales, "Adaptation of Precision Matrix Models on Large Vocabulary Continuous Speech Recognition", *Proc. ICASSP*, 2005.