

# Error Reduction Network for DBLSTM-based Voice Conversion

Mingyang Zhang<sup>\*†</sup>, Berrak Sisman<sup>†</sup>, Sai Sirisha Rallabandi<sup>†</sup>, Haizhou Li<sup>†</sup>, Li Zhao<sup>\*</sup>

<sup>\*</sup> Key Laboratory of Underwater Acoustic Signal Processing of Ministry of Education, Southeast University, Nanjing, China

E-mail: zhangmy@seu.edu.cn, zhaoli@seu.edu.cn

<sup>†</sup> National University of Singapore, Singapore

E-mail: berraksisman@u.nus.edu, siri.gene@gmail.com, haizhou.li@nus.edu.sg

**Abstract**—So far, many of the deep learning approaches for voice conversion produce good quality speech by using a large amount of training data. This paper presents a Deep Bidirectional Long Short-Term Memory (DBLSTM) based voice conversion framework that can work with a limited amount of training data. We propose to implement a DBLSTM based average model that is trained with data from many speakers. Then, we propose to perform adaptation with a limited amount of target data. Last but not least, we propose an error reduction network that can improve the voice conversion quality even further. The proposed framework is motivated by three observations. Firstly, DBLSTM can achieve a remarkable voice conversion by considering the long-term dependencies of the speech utterance. Secondly, DBLSTM based average model can be easily adapted with a small amount of data, to achieve a speech that sounds closer to the target. Thirdly, an error reduction network can be trained with a small amount of training data, and can improve the conversion quality effectively. The experiments show that the proposed voice conversion framework is flexible to work with limited training data and outperforms the traditional frameworks in both objective and subjective evaluations.

## I. INTRODUCTION

Voice Conversion (VC) is a technology that modifies the speech of the source speaker to make it sounds like the target speaker. The voice conversion technology has been applied to many tasks, such as Text-to-Speech (TTS) system [1], speech enhancement [2] and speaking assistance [3].

Voice conversion can be formulated as a regression problem of estimating a mapping function between the source and target features. Many state-of-the-art approaches for voice conversion are including Gaussian Mixed Model (GMM) [4–6] which is based on the maximum-likelihood estimation of spectral parameter trajectory. Dynamic Kernel Partial Least Squares (DKPLS) [7] integrates a kernel transformation into partial least squares to model nonlinearities as well as to capture the dynamics in the data. Sparse representation [8–11] can be seen as a data-driven, non-parametric approach as an alternative to the traditional parametric approaches to voice conversion. Frequency warping based approaches [12–14] aim to modify the frequency axis of source spectra towards that of the target. There are also some post-filter approaches for voice conversion to improve the speech quality [15, 16].

Recently, deep learning approaches became popular in the field of voice conversion. For example, Deep Neural Network (DNN) based approaches [17–19] focus on spectrum

conversion under the constraint of parallel training data and achieve high-quality speech by using a large amount of parallel training data. In addition, there have been some researches on variational autoencoder [20, 21] that effectively improve the conversion performance.

The above-mentioned voice conversion frameworks consider the frame features as individual components, and do not concern about the long-term dependencies of the speech sequences. Standard Recurrent Neural Networks (RNNs) can be used to solve this problem [22, 23], but it has limited ability in modeling context because of the vanishing gradient problem [24]. Moreover, the standard RNNs can only capture the information from the former sequences and not the latter sequences.

To alleviate these problems, Deep Bidirectional Long Short-Term Memory (DBLSTM) has been proposed to perform voice conversion [25–27], and achieves remarkable performance over the traditional DNN-based voice conversion framework [27]. Moreover, DBLSTM has been successfully used in various tasks in the field of speech and language processing, such as Automatic Speech Recognition (ASR) [28–30], speech synthesis [31] and emotion recognition [32, 33].

Although the DBLSTM and DNN based voice conversion frameworks can achieve good voice conversion performance, they still suffer from the dependency of a large amount of training data which is not practical in real life. The remaining issue is to find a way to make a good use of limited data. Different from the previous studies, in this paper, we take advantage of the powerful deep learning framework DBLSTM, and propose a voice conversion framework that can produce high-quality speech under the constraint of limited parallel data. Specifically, we make the following contributions: 1) due to DBLSTM can achieve a remarkable voice conversion by considering the long-term dependencies of the speech utterance, we build a DBLSTM-based average model by using data from many speakers; 2) since the DBLSTM-based average model can be easily adapted with a small amount of data, we perform adaptation to the DBLSTM-based average model by using limited target data, to achieve a converted sound that is more similar to that of target; 3) an error reduction network can be trained with a small amount of training data from source and target, so we propose an error reduction network for the adapted DBLSTM framework, that can improve the

voice conversion quality even further. Overall, we propose a DBLSTM-based voice conversion framework that can produce high-quality speech with a small amount of training data.

The rest of this paper is organized as follows: Section II explains the traditional DBLSTM-based approach of voice conversion. Section III describes our proposed voice conversion framework that is based on an error reduction network for the adapted DBLSTM-based approach. We report the objective and subjective results in Section IV. Section V concludes this paper.

## II. DBLSTM-BASED VOICE CONVERSION

The network architecture of BLSTM-based Voice Conversion is a combination of bidirectional RNNs and LSTM memory block, which can learn long-range contextual in both forward and backward directions. By stacking multiple hidden layers, a deep network architecture is built to capture the high-level representation of voice features. For bidirectional RNNs, the iteration functions for the forward sequence  $\vec{h}$  and backward sequence  $\overleftarrow{h}$  are as follows:

$$\vec{h}_t = \mathcal{H}(W_{x\vec{h}}x_t + W_{\vec{h}\vec{h}}\vec{h}_{t-1} + b_{\vec{h}}) \quad (1)$$

$$\overleftarrow{h}_t = \mathcal{H}(W_{x\overleftarrow{h}}x_t + W_{\overleftarrow{h}\overleftarrow{h}}\overleftarrow{h}_{t+1} + b_{\overleftarrow{h}}) \quad (2)$$

$$y_t = W_{\vec{h}y}\vec{h}_t + W_{\overleftarrow{h}y}\overleftarrow{h}_t + b_y \quad (3)$$

where  $x, y, h, t$  are the input, output, hidden state and time sequence respectively. A LSTM network consists of recurrently connected blocks, known as memory block. Every memory block contains self-connected memory cells and three adaptive and multiplicative gate units. The recurrent hidden layer function  $\mathcal{H}$  of the LSTM network is implemented according to the following equations:

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \quad (4)$$

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \quad (5)$$

$$c_t = f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \quad (6)$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \quad (7)$$

$$h_t = o_t \tanh(c_t) \quad (8)$$

where  $i, f, o, c$  refer to the input gate, forget gate, output gate and the element of cells  $C$  respectively.  $\sigma$  is the logistic sigmoid function.

The overall framework of a DBLSTM-based voice conversion is shown in Fig. 1. In this framework, the three feature streams including the spectrum feature,  $\log(F_0)$  and the aperiodic component are converted separately. The spectrum feature is converted by the DBLSTM model.  $\log(F_0)$  is converted by equalizing the mean and the standard deviation of the source and target speech. And the aperiodic component is directly copied to synthesize the converted speech. The whole utterance is treated as input so that the system can access the long-range context in both forward and backward directions. In this paper, we propose to use DBLSTM to perform voice conversion under limited training data.

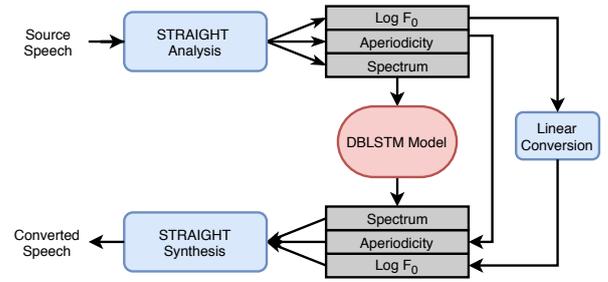


Fig. 1. DBLSTM-based voice conversion framework.

## III. ERROR REDUCTION NETWORK FOR ADAPTED DBLSTM-BASED VOICE CONVERSION

The DBLSTM-based voice conversion has a good performance while it needs a large amount of parallel data from source speaker and target speaker, which is expensive to collect in practice. To solve this problem, we propose an error reduction network for adapted DBLSTM-based approach.

### A. Training Phase

As illustrated in Fig. 2, the proposed approach can be divided into three training phases. In training phase 1, an average DBLSTM model is trained for the one-hot phoneme label to Mel-cepstral coefficients (MCEPs) mapping, using the data from many speakers except the source speaker and the target speaker. MCEPs are the Mel Log Spectral Approximation (MLSA) parameters which approximate Mel-Frequency Cepstral Coefficients (MFCCs). A trained ASR system is used to extract the phoneme information of the input speech. The input of the ASR model is MFCC feature of the speech frame. The output is a one-hot phoneme label vector that indicates the phoneme information of the speech frame. Then a DBLSTM-based model is trained to get the mapping relationship between the one-hot phoneme label and the corresponding MCEPs which are extracted by STRAIGHT [34]. We call this framework as Average Model, it can only generate MCEPs of an average voice of the speakers whose data are used.

In training phase 2, the average model is adapted using a small amount of data from the target speaker. The adaptation is similar to the training of the average model, except the initialized network is the trained average model and the training data is the target speech. After the adaptation, the output of the adapted model will be closer to the target speaker. We call this framework as adapted average model. However, there always exists an error between the converted features and the target features. This error degrades both quality and similarity of the converted speech [35]. To reduce such error, we propose an error reduction network after the adapted average model.

Training phase 3 involves the error reduction network which is essentially an additional DBLSTM network, used to map the converted MCEPs to the target MCEPs. The error reduction network brings the final output MCEPs features closer to the

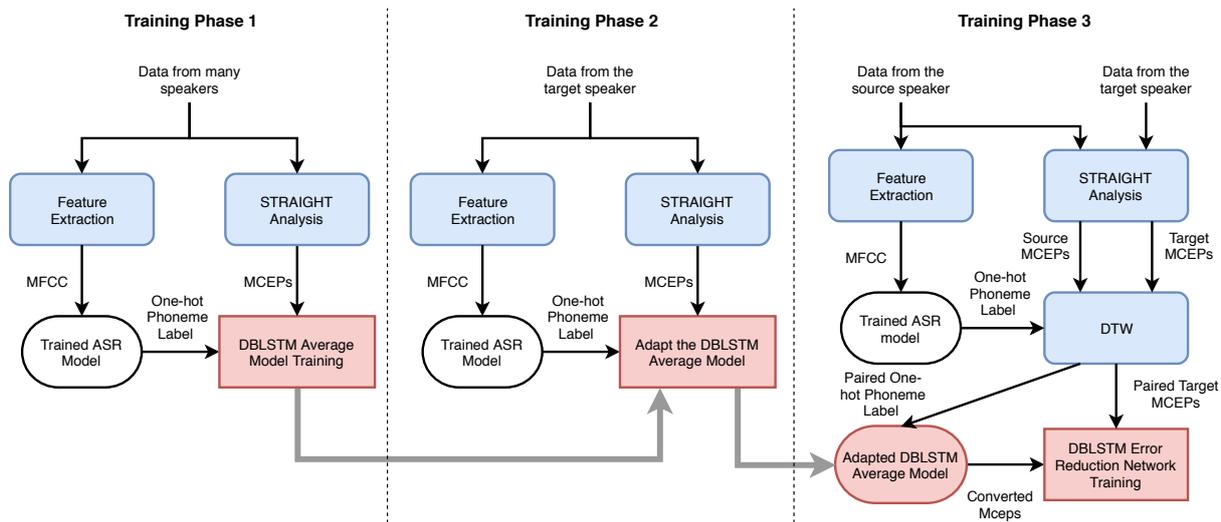


Fig. 2. The proposed DBLSTM based VC framework. In training phase 1, we exclude the data from the source and the target speakers. In training phase 2, we only use the data from the target speaker. In training phase 3, we use the same sentences from both source and target speakers.

target speaker. The same utterances have been used in the adapted average model of the target speaker and the parallel data of the source speaker are used in the error reduction network. The same ASR system is used to generate the one-hot phoneme label of the target speech. MCEPs features from the same sentences of the source speech and the target speech are aligned by dynamic time warping (DTW), and the alignment information is also used to get the paired one-hot phoneme label. Then feed the label to the adapted average model to generate the paired converted MCEPs. For the training of the error reduction network, the input is the paired converted MCEPs, and the output is the original feature of the target speech. The error reduction network can reduce the error that created by the previous training part.

### B. Run-time Conversion Phase

In the conversion stage, the input is a whole sentence of the source speaker. The conversion of  $\log(F_0)$  and aperiodicity is the same as that of the DBLSTM-based system mentioned in Section II. MFCC features of the source speech are used by the trained ASR model to obtain the one-hot phoneme labels. Next, the one-hot phoneme labels are converted to MCEPs by the trained adapted average model. Then, the converted MCEPs are fed into the error reduction network to get the final result. Finally, the converted MCEPs together with the converted  $\log(F_0)$  and aperiodicity are used by the STRAIGHT vocoder to synthesize the output speech.

## IV. EXPERIMENTS

### A. Experimental Setup

We conduct listening experiments to assess the performance of our proposed framework that is error reduction network for adapted DBLSTM-based voice conversion. We compare this framework with the baseline DBLSTM [27] that is explained

in Section II, and DBLSTM-based adapted average model that is explained in Section III and given in Fig. 3. We note that the adapted average model is an intermediate step of our proposed algorithm. Fig. 4 also shows the differences between of our proposed framework and the adapted average model.

The database used in the experiments is CMU ARCTIC corpus [36]. As it is the most challenging work in voice conversion, we conduct the cross-gender voice conversion experiments. The speech signals are sampled at 16kHz with mono channel, windowed by 25ms, and the frameshift is 5ms. For the DBLSTM-based average model training, data from four male speakers (awb, jmk, ksp, rms) are used. 4433 and 489 sentences are used as training data and validation data. In training phase 2, 45 and 5 sentences from the target speaker (slt) are used as training data and validation data to adapt the average model. For the error reduction training, the same sentences of the target speaker that have been used to adapt the average model and the parallel data of the source speaker (bdl) are used. A DNN-HMM based ASR system [37] is used to get the 171-dimension one-hot phoneme label. 40-dimension MCEPs are extracted by STRAIGHT to train the model.

In our proposed approach, to train a DBLSTM-based model, we prefer to use four hidden layers, the number of units in each layer is [171 128 256 256 128 40] respectively. Each bidirectional LSTM hidden layer contains one forward LSTM layer and one backward LSTM layer. The training samples are normalized to zero mean and unit variance for each dimension before training. For the error reduction network, in order to take advantage of context information, three frames of converted MCEPs i.e. current frame, one left frame and 1 right frame are used as input features. In addition, there are three hidden layers in the error reduction network, the number of units in each layer is [120 128 256 128 40] respectively.

In order to evaluate our proposed approach, 100 parallel

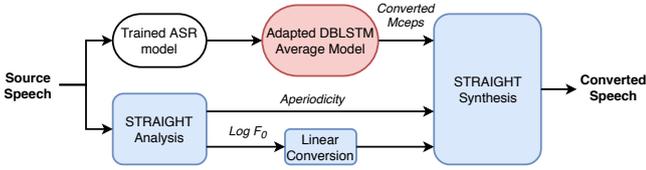


Fig. 3. The run-time phase of the adapted average model.

utterances from the source speaker and the target speaker are used to train a DBLSTM-based parallel voice conversion system. This system is developed as the baseline approach. There are four hidden layers in the baseline model where the number of units in each layer is same as the training of the adapted average model that is [40 128 256 256 128 40] respectively. We use an open-source CUDA recurrent neural network toolkit CURRENNT [38] to train the DBLSTM model with a learning rate of  $10^{-5}$  and a momentum of 0.9.

### B. Objective Evaluation

Mel-cepstral distortion (MCD) [39] is used as objective measure of the spectral distance from converted to target speech, which is denoted as:

$$MCD[dB] = \frac{10}{\ln 10} \sqrt{2 \sum_{d=1}^D (C_d^{target} - C_d^{converted})^2} \quad (9)$$

where  $C_d^{target}$  and  $C_d^{converted}$  are the  $d^{th}$  dimension of the original target MCEPs and the converted MCEPs, respectively. We expect a good system to report a low MCD value.

The MCD scores of the different systems for the cross-gender voice conversion are summarized in Table I. We can see that our proposed approach outperforms the baseline model and the adapted average model. We can also note that the MCD scores of the adapted average model is not as good as the baseline model, because there is no parallel data in the training of the adapted average model. But after the error reduction network with only 50 utterances of parallel data, the performance can be improved obviously, and outperform both adapted average model and baseline model with 100 utterances of parallel training data.

TABLE I  
THE MCD OF DIFFERENT SYSTEMS.

Source-Target	Baseline	Adapted Average Model	Proposed Approach
9.3197	6.3042	6.7378	<b>6.1989</b>

### C. Subjective Evaluation

To evaluate the quality and similarity of the converted speech from the different systems, we conduct a subjective listening test and 10 listeners are invited to evaluate 10 sentences in each system.

We carry out Mean Opinion Score (MOS) test for evaluating speech quality and naturalness. In the MOS test, comparing

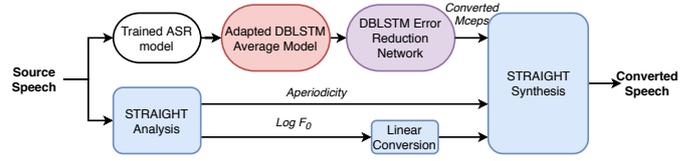


Fig. 4. The run-time phase of the proposed framework.

with target speech, the grades of the converted speech are: 5 = excellent, 4 = good, 3 = fair, 2 = poor, and 1 = bad. The listeners are asked to rate the speech according to this regulation. In this experiment, we conduct the MOS test among three systems: 1) baseline approach, the parallel DBLSTM-based voice conversion training with 100 utterances of parallel data; 2) adapted average model that explained in Section III; 3) our proposed approach. The results of the MOS test and the 95% confidence intervals are shown in Fig. 5. The scores of the baseline, adapted average Model and the proposed approach are 2.62, 2.71 and 3.41 respectively.

ABX preference test is adopted to evaluate speaker similarity of the converted speech generated by two different systems. The listeners are asked to choose either A or B that sounds closer to the target speaker's speech X. We conduct the ABX preference test between the baseline approach and our proposed approach. The preference bars for speaker similarity are shown in Fig. 6.

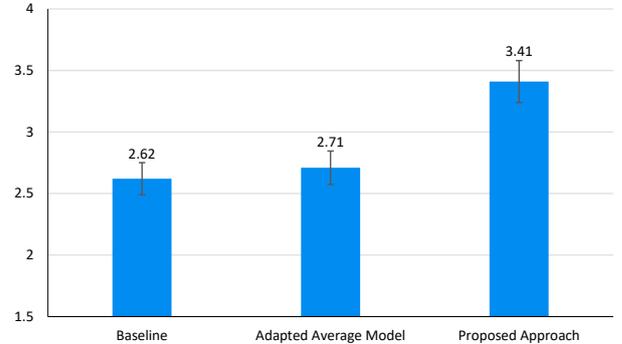


Fig. 5. The result of the MOS test with the 95% confidence intervals for speech quality and naturalness among the three systems.

Overall, the results of both MOS test and ABX preference test show that our proposed error reduction network for adapted DBLSTM-based voice conversion with a limited amount of parallel data outperforms the baseline approach with a large amount of parallel data in both speech quality and similarity. Probably because the average model is trained with a large amount of data to achieve a better speech quality than the baseline approach, improve the performance of the following portions of the system.

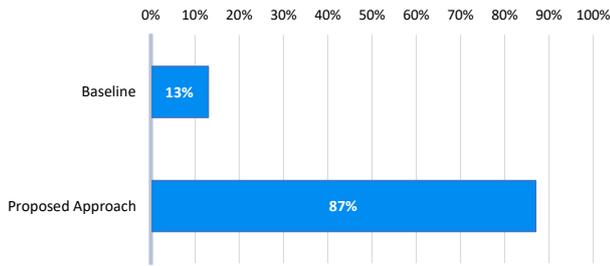


Fig. 6. The result of the ABX preference test for speaker similarity between the baseline and our proposed approach.

## V. CONCLUSIONS

This paper presents an error reduction network for adapted DBLSTM-based voice conversion approach, which can achieve a good performance with limited parallel data of the source speaker and the target speaker. Firstly, we propose to train an average model for the one-hot phoneme label to MCEPs mapping with data from many speakers exclude the source speaker and the target speaker. Then, we propose to adapt the average model with a limited amount of target data. Furthermore, we implement an error reduction network that can improve the voice conversion quality. Experiment results from both objective and subjective evolution show that our proposed approach can make a good use of limited data, and outperforms the baseline approach. In the future, we will investigate to use the WaveNet Vocoder, which is a convolutional neural network that can generate raw audio waveform sample by sample, to improve the quality and naturalness of the converted speech. Some samples for the listening test are available through this link: <https://arkhamimp.github.io/ErrorReductionNetwork/>

## ACKNOWLEDGMENT

This research is supported by Ministry of Education, Singapore AcRF Tier 1 NUS Start-up Grant FY2016. Mingyang Zhang is also supported by the China Scholarship Council (Grant No.201706090063). Berrak Sisman is also funded by SINGA Scholarship under A\*STAR Graduate Academy.

## REFERENCES

- [1] A. Kain and M. W. Macon, "Spectral voice conversion for text-to-speech synthesis," in *Acoustics, Speech and Signal Processing, 1998. Proceedings of the 1998 IEEE International Conference on*, vol. 1, pp. 285–288 vol.1, May 1998.
- [2] T. Toda, M. Nakagiri, and K. Shikano, "Statistical voice conversion techniques for body-conducted unvoiced speech enhancement," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, pp. 2505–2517, Nov 2012.
- [3] K. Nakamura, T. Toda, H. Saruwatari, and K. Shikano, "Speaking-aid systems using gmm-based voice conver-

- sion for electrolaryngeal speech," *Speech Communication*, vol. 54, no. 1, pp. 134 – 146, 2012.
- [4] T. Toda, A. W. Black, and K. Tokuda, "Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 15, no. 8, pp. 2222–2235, 2007.
- [5] S. Takamichi, T. Toda, A. W. Black, and S. Nakamura, "Modulation spectrum-constrained trajectory training algorithm for gmm-based voice conversion," in *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on*, pp. 4859–4863, IEEE, 2015.
- [6] K. Tanaka, S. Hara, M. Abe, M. Sato, and S. Minagi, "Speaker dependent approach for enhancing a glossectomy patient's speech via gmm-based voice conversion," *Proc. Interspeech 2017*, pp. 3384–3388, 2017.
- [7] E. Helander, H. Silen, T. Virtanen, and M. Gabouj, "Voice conversion using dynamic kernel partial least squares regression," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 3, pp. 806–817, 2012.
- [8] Z. Wu, T. Virtanen, E. S. Chng, and H. Li, "Exemplar-based sparse representation with residual compensation for voice conversion," *IEEE/ACM Transactions on Audio, Speech and Language Processing*, vol. 22, no. 10, pp. 1506–1521, 2014.
- [9] R. Takashima, T. Takiguchi, and Y. Ariki, "Exemplar-based voice conversion using sparse representation in noisy environments," *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, vol. 96, no. 10, pp. 1946–1953, 2013.
- [10] B. Çiřman, H. Li, and K. C. Tan, "Sparse representation of phonetic features for voice conversion with and without parallel data," in *Automatic Speech Recognition and Understanding Workshop (ASRU), 2017 IEEE*, pp. 677–684, IEEE, 2017.
- [11] B. Sisman, H. Li, and K. C. Tan, "Transformation of prosody in voice conversion," in *2017 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pp. 1537–1546, Dec 2017.
- [12] X. Tian, Z. Wu, S. W. Lee, and E. S. Chng, "Correlation-based frequency warping for voice conversion," in *Chinese Spoken Language Processing (ISCSLP), 2014 9th International Symposium on*, pp. 211–215, IEEE, 2014.
- [13] X. Tian, Z. Wu, S. W. Lee, N. Q. Hy, M. Dong, and E. S. Chng, "System fusion for high-performance voice conversion," in *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, vol. 2015-January, pp. 2759–2763, 2015.
- [14] X. Tian, S. W. Lee, Z. Wu, E. S. Chng, S. Member, and H. Li, "An Exemplar-based Approach to Frequency Warping for Voice Conversion," pp. 1–10, 2016.
- [15] S. Takamichi, T. Toda, A. W. Black, G. Neubig, S. Sakti, and S. Nakamura, "Postfilters to modify the modulation

- spectrum for statistical parametric speech synthesis,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, pp. 755–767, April 2016.
- [16] N. Xu, X. Yao, A. Jiang, X. Liu, and J. Bao, “High quality voice conversion by post-filtering the outputs of gaussian processes,” in *2016 24th European Signal Processing Conference (EUSIPCO)*, pp. 863–867, Aug 2016.
- [17] L.-h. Chen, Z.-h. Ling, L.-j. Liu, and L.-r. Dai, “Voice Conversion Using Deep Neural Networks With Layer-Wise Generative Training,” *IEEE Transactions on Audio, Speech and Language Processing*, vol. 22, no. 12, pp. 1859–1872, 2014.
- [18] T. Nakashika, R. Takashima, T. Takiguchi, and Y. Ariki, “Voice conversion in high-order eigen space using deep belief nets,” in *INTERSPEECH*, no. August, pp. 369–372, 2013.
- [19] S. H. Mohammadi and A. Kain, “Voice conversion using deep neural networks with speaker-independent pre-training,” in *Spoken Language Technology Workshop (SLT), 2014 IEEE*, pp. 19–23, IEEE, 2014.
- [20] C.-C. Hsu, H.-T. Hwang, Y.-C. Wu, Y. Tsao, and H.-M. Wang, “Voice conversion from non-parallel corpora using variational auto-encoder,” in *Signal and Information Processing Association Annual Summit and Conference (APSIPA), 2016 Asia-Pacific*, pp. 1–6, IEEE, 2016.
- [21] C.-C. Hsu, H.-T. Hwang, Y.-C. Wu, Y. Tsao, and H.-M. Wang, “Voice conversion from unaligned corpora using variational autoencoding wasserstein generative adversarial networks,” *arXiv preprint arXiv:1704.00849*, 2017.
- [22] T. Nakashika, T. Takiguchi, and Y. Ariki, “High-order sequence modeling using speaker-dependent recurrent temporal restricted Boltzmann machines for voice conversion,” in *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, no. September, pp. 2278–2282, 2014.
- [23] T. Nakashika, T. Takiguchi, and Y. Ariki, “Voice conversion using rnn pre-trained by recurrent temporal restricted Boltzmann machines,” *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)*, vol. 23, no. 3, pp. 580–587, 2015.
- [24] Y. Bengio, P. Simard, and P. Frasconi, “Learning long-term dependencies with gradient descent is difficult,” *IEEE transactions on neural networks*, vol. 5, no. 2, pp. 157–166, 1994.
- [25] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [26] A. Graves and J. Schmidhuber, “Framewise phoneme classification with bidirectional lstm and other neural network architectures,” *Neural Networks*, vol. 18, no. 5-6, pp. 602–610, 2005.
- [27] L. Sun, S. Kang, K. Li, and H. Meng, “Voice conversion using deep bidirectional Long Short-Term Memory based Recurrent Neural Networks,” in *ICASSP*, no. 1, pp. 4869–4873, 2015.
- [28] A. Graves, N. Jaitly, and A.-r. Mohamed, “Hybrid speech recognition with deep bidirectional lstm,” in *Automatic Speech Recognition and Understanding (ASRU), 2013 IEEE Workshop on*, pp. 273–278, IEEE, 2013.
- [29] M. Wöllmer, Z. Zhang, F. Weninger, B. Schuller, and G. Rigoll, “Feature enhancement by bidirectional lstm networks for conversational speech recognition in highly non-stationary noise,” in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pp. 6822–6826, IEEE, 2013.
- [30] A. Zeyer, P. Doetsch, P. Voigtlaender, R. Schlüter, and H. Ney, “A comprehensive study of deep bidirectional lstm rnns for acoustic modeling in speech recognition,” in *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*, pp. 2462–2466, IEEE, 2017.
- [31] Y. Fan, Y. Qian, F.-L. Xie, and F. K. Soong, “Tts synthesis with bidirectional lstm based recurrent neural networks,” in *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- [32] M. Wöllmer, A. Metallinou, F. Eyben, B. Schuller, and S. Narayanan, “Context-sensitive multimodal emotion recognition from speech and facial expression using bidirectional lstm modeling,” in *Proc. INTERSPEECH 2010, Makuhari, Japan*, pp. 2362–2365, 2010.
- [33] L. He, D. Jiang, L. Yang, E. Pei, P. Wu, and H. Sahli, “Multimodal affective dimension prediction using deep bidirectional long short-term memory recurrent neural networks,” in *Proceedings of the 5th International Workshop on Audio/Visual Emotion Challenge*, pp. 73–80, ACM, 2015.
- [34] H. Kawahara, I. Masuda-Katsuse, and A. de Cheveigne, “Restructuring speech representations using a pitch-adaptive timefrequency smoothing and an instantaneous-frequency-based f0 extraction: Possible role of a repetitive structure in sounds,” *Speech communication*, pp. 187–207, 1999.
- [35] J. Wu, D. Huang, L. Xie, and H. Li, “Denoising recurrent neural network for deep bidirectional lstm based voice conversion,” *Proc. Interspeech 2017*, pp. 3379–3383, 2017.
- [36] J. Kominek and A. W. Black, “The cmu arctic speech databases,” in *Fifth ISCA Workshop on Speech Synthesis*, pp. 223–224, 2004.
- [37] D. Povey, A. Ghoshal, N. Goel, M. Hannemann, Y. Qian, P. Schwarz, J. Silovsk, and P. Motl, “The Kaldi Speech Recognition Toolkit,” in *IEEE ASRU*, 2011.
- [38] F. Weninger, J. Bergmann, and B. Schuller, “Introducing currennt: The munich open-source cuda recurrent neural network toolkit,” *The Journal of Machine Learning Research*, vol. 16, no. 1, pp. 547–551, 2015.
- [39] R. Kubichek, “Mel-cepstral distance measure for objective speech quality assessment,” *Communications, Computers and Signal Processing*, pp. 125–128, 1993.