



Investigation of using disentangled and interpretable representations for one-shot cross-lingual voice conversion

Seyed Hamidreza Mohammadi, Taehwan Kim

ObEN, Inc.

hamid@oben.com, taehwan@oben.com

Abstract

We study the problem of cross-lingual voice conversion in non-parallel speech corpora and one-shot learning setting. Most prior work require either parallel speech corpora or enough amount of training data from a target speaker. However, we convert an arbitrary sentences of an arbitrary source speaker to target speaker's given only one target speaker training utterance. To achieve this, we formulate the problem as learning disentangled speaker-specific and context-specific representations and follow the idea of [1] which uses Factorized Hierarchical Variational Autoencoder (FHVAE). After training FHVAE on multi-speaker training data, given arbitrary source and target speakers' utterance, we estimate those latent representations and then reconstruct the desired utterance of converted voice to that of target speaker. We investigate the effectiveness of the approach by conducting voice conversion experiments with varying size of training utterances and it was able to achieve reasonable performance with even just one training utterance. We also examine the speech representation and show that World vocoder outperforms Short-time Fourier Transform (STFT) used in [1]. Finally, in the subjective tests, for one language and cross-lingual voice conversion, our approach achieved significantly better or comparable results compared to VAE-STFT and GMM baselines in speech quality and similarity.

Index Terms: voice conversion, one-shot learning, cross-lingual, variational autoencoder

1. Introduction

The task of Voice Conversion (VC) [2, 3] is a technique to convert source speaker's spoken sentences into those of a target speaker's voice. It requires to preserve not only the target speaker's identity, but also phonetic context spoken by the source speaker. To tackle this problem, many approaches have been proposed [4, 5, 6]. However, most prior work require parallel spoken corpus and enough amount of data to learn the target speaker's voice. Recently, there were approaches proposed for voice conversion with non-parallel corpus [7, 8, 9]. But they still require that speaker identity was known *priori*, or included in training data for the model.

Recently, Hsu et al. [1] proposed to use disentangled and interpretable representations to overcome these limitations by exploiting Factorized Hierarchical Variation Autoencoder. They achieved reasonable quality with just single utterance from a target speaker but it was still not satisfactory. Nevertheless, most prior work focus on voice conversion *within* one language. But we believe that if we can capture disentangled representations of phonetic or linguistic contexts and speaker identities, the model should be capable for more challenging *cross-lingual* setting, which means that source and target speakers are from different languages. Therefore, we focus on investigating cross-lingual voice conversion, and propose to follow the same spirit from Hsu et al. [1] and improve the performance. Our contri-

butions are:

- We investigate the different feature representations for spoken utterances by considering Mel-cepstrum (MCEP) features and other acoustic features, and achieve better results compared to baselines.
- We examine the effect of the size of training utterances from source and target speakers, and demonstrate that with just a few, or even one, utterances, we are able to achieve the reasonable performance.
- We conduct cross-lingual voice conversion experiments and our approach achieved significantly better or comparable results than baselines in speech quality and similarity in the subjective tests.

2. Related Work

Voice conversion has been an important research problem for over a decade. One popular approach to tackle the problem is spectral conversion such as Gaussian mixture models (GMMs) [4] and deep neural networks (DNN) [5]. However, it requires parallel spoken corpus and dynamic time warping (DTW) is usually used to align source and target utterances. To overcome this limitation, non-parallel voice conversion approaches were proposed, for instance, *eigenvoice* [6], *i-vector* [10], and Variational Autoencoder [7, 9] based models. However, *eigenvoice* based approach [6] still requires reference speaker to train the model, and VAE based approaches [7, 9] require speaker identities to be known *priori* as included in training data for the model. *i-vector* based approach [10] looks promising which remains to be studied further. The *i-vectors* are converted by replacing the source latent variable by the target latent variable. The Gaussian mixture means are then reconstructed from the converted *i-vector*. The Gaussians with adjusted means are then applied to the source vector to perform the acoustic feature conversion. Siamese autoencoder has also been proposed for decomposing speaker identity and linguistic embeddings [11]. However, this approach requires parallel training data to learn the decomposing architecture. This decomposition is achieved by means of applying some similarity and non-similarity costs between the Siamese architectures.

Nonetheless, cross-lingual voice conversion is also a challenging task since target language is not known in training time, and only few work has proposed, including GMMs based approach [12] and *eigenvoice* based approach [13], but still have inherent limitations as above.

Recently, deep generative models have been applied and successful for unsupervised learning tasks, and include Variational Autoencoder (VAE) [14], Generative Adversarial Networks (GAN) [15], and auto-regressive models [16, 17]. Among them, VAE can infer latent codes from data and generate data from them by jointly learning inference and generative networks, and VAE has been also applied for voice conversion

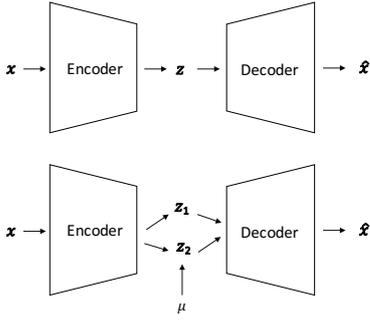


Figure 1: Structures of Variation Autoencoder (upper) and Factorized Hierarchical Variational Autoencoder (lower).

[7, 9]. However, in their models, speaker identities are not inferred from data and instead required to be known in model training time. GAN has been also exploited for non-parallel voice conversion [18] with the cycle consistency constraint [19], but it still has the limitation that it needs to know the target speaker in training time and be trained for each target.

To understand the disentangled and interpretable structure of latent codes, several work were proposed, namely, DC-IGN [20], InfoGAN [21], β -VAE [22], and FHVAE [1]. These approaches to uncover disentangled representation may help voice conversion with very limited resource from target speaker, since it might infer speaker identity information from data without supervision, as illustrated in FHVAE [1]. However, the qualities of converted voices were not good enough, therefore, we focus on the model structure of FHVAE and investigate to improve it, even with more challenging cross-lingual voice conversion setting.

3. Model

Variational autoencoder [14] (VAE) is a powerful model to uncover hidden representation and generate new data samples. Let observations be x and latent variables z . In the variational autoencoder model, the encoder (or inference network) $q_\phi(z|x)$ outputs z given input x , and decoder $p_\Phi(x|z)$ generates data x given z . The encoder and decoder are neural networks. Training is done by maximizing variational lower bound (or also called evidence lower bound):

$$\begin{aligned} \ell(\Phi, \phi) &= \mathbb{E}_q[\log p_\Phi(x, z)] - \mathbb{E}_q[\log q_\phi(z|x)] \\ &= \log p_\Phi(x) - D_{KL}(q_\phi(z|x)||p_\Phi(z|x)). \end{aligned}$$

where D_{KL} is Kullback-Leibler divergence.

However, VAE considers no structure for latent variable z . Assuming structure for z could be beneficial to exploit the inherent structures in data. Here we describe Factorized Hierarchical Variational Autoencoder proposed by Hsu et al [1]. Let a dataset D consist of N_{seg} i.i.d. sequences X^i . For each sequence X^i , it consists of N_{seg}^i observation segments. Then we define factorized latent variables of latent segment variable $Z_1^{i,j}$ and latent sequence variable $Z_2^{i,j}$. In the context of voice conversion, $Z_1^{i,j}$ is responsible for generating phonetic contexts and $Z_2^{i,j}$ is for speaker identity. When generating data $X^{i,j}$, we first sample $Z_2^{i,j}$ from isotropic Gaussian centered at μ^i shared for the entire sequence, and also $Z_1^{i,j}$ independently. Then we generate $X^{i,j}$ conditioned on $Z_1^{i,j}$ and $Z_2^{i,j}$. Thus, joint probability with a sequence X^i is:

$$p_\Phi(X^i, Z_1^i, Z_2^i, \mu^i) = p_\Phi(\mu^i) \prod_{j=1}^{N_{seg}^i} p_\Phi(X^{i,j}|Z_1^{i,j}, Z_2^{i,j}) p_\Phi(Z_1^{i,j}) p_\Phi(Z_2^{i,j}|\mu^i)$$

This is illustrated in Figure 1. For inference, we use variational inference to approximate the true posterior and have:

$$q_\phi(Z_1^i, Z_2^i, \mu^i|X^i) = q_\phi(\mu^i) \prod_{j=1}^{N_{seg}^i} q_\phi(Z_1^{i,j}|X^{i,j}, Z_2^{i,j}) q_\phi(Z_2^{i,j}|X^{i,j})$$

Since sequence variational lower bound can be decomposed to segment variational lower bound, we can use batches of segment instead of sequence level to maximize:

$$\begin{aligned} \ell(\Phi, \phi; X^{i,j}) &= \ell(\Phi, \phi; X^{i,j}|\tilde{\mu}^i) + \frac{1}{N_{seg}^i} \log p_\Phi(\tilde{\mu}^i) + const \\ \ell(\Phi, \phi; X^{i,j}|\tilde{\mu}^i) &= \mathbb{E}_{q_\phi(Z_1^{i,j}, Z_2^{i,j}|X^{i,j})}[\log p_\Phi(X^{i,j}|Z_1^{i,j}, Z_2^{i,j})] \\ &\quad - \mathbb{E}_{q_\phi(Z_2^{i,j}|X^{i,j})}[D_{KL}(q_\phi(Z_1^{i,j}|X^{i,j}, Z_2^{i,j})||p_\Phi(Z_1^{i,j}))] \\ &\quad - D_{KL}(q_\phi(Z_2^{i,j}|X^{i,j})||p_\Phi(Z_2^{i,j}|\tilde{\mu}^i)) \end{aligned}$$

where $\tilde{\mu}^i$ is the posterior mean of μ^i . Please refer to Hsu et al. [1] for more details. Additionally, Hsu et al. also proposed discriminative segment variational lower bound to encourage Z_2^i to be more sequence-specific by adding the additional term of inferring the sequence index i from $Z_2^{i,j}$. For our experiments, we exploit this FHVAE model and sequence-to-sequence model [23] as the structure of encoder-decoder for sequential data.

For performing the voice conversion, we compute the average Z_2 from the training utterance(s) of source and target speakers. For a given input utterance, we compute Z_1 and Z_2 of the input utterance. There are two ways to perform voice conversion. First, we can replace Z_2 values of the source speaker with the average Z_2 from the target speaker. This approach resulted in too muffled generated result. Second, we compute a difference vector between source and target average $Z_2^{diff} = Z_2^{trg} - Z_2^{src}$. This difference vector is added to Z_2 from the input utterance as $Z_2^{converted} = Z_2 + Z_2^{diff}$ and then decoded using FHVAE to achieve the speech features. In an informal listening test, we decided to the second approach since it resulted in significantly higher quality generated speech.

4. Experiments

4.1. Datasets

We used the TIMIT corpus [24] which is a multi-speaker speech corpus as the training data for FHVAE model. We used the training speakers as suggested by the corpus to train the model. For English test speakers, we select speakers from TIMIT testing part of the corpus. We also use a proprietary Chinese speech corpus (hereon referred to as CH) with 5200 speakers each uttering one sentence. We consider using the combination of TIMIT and Chinese corpus for training the model as well. For Chinese test speakers, we utilize speakers from the THCHS-30 speech corpus [25]. To observe the effect of having more utterances per speaker but less speakers we also train the model on VCTK corpus [26]. Finally, for objective testing (which requires availability of parallel data), we utilized four CMU-arctic voices (BDL, SLT, RMS, CLB)[27]. As speech features, we used 40th-order MCEPs (excluding the energy coefficient, dimensionality $D=39$), extracted using the World toolkit [28] with

a 5ms frame shift. All audio files are transformed to 16kHz and 16 bit before any analysis.

4.2. Experimental setting

For the encoder and decoder in FHVAE model, we use Long Short Term Memory (LSTM) [29] as the first layer with 256 hidden units with a fully-connected layer on top. We use 32 dimensions for each latent variable Z_1 and Z_2 . The models were trained with stochastic gradient descent. We use a mini-batch size of 256. The Adam optimizer [30] is used with $\beta_1 = 0.95$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$, and initial learning rate of 10^{-4} . The model is trained for 500 epochs and select the model best performing on the development set.

From now on, we use the abbreviation VAE for FHVAE model. In our experiments, we consider three models: GMM (GMM MAP Adaptation [4]), VAE-STFT (uses STFT as speech analysis/synthesis[1]), VAE (uses World as speech analysis/synthesis[28]). We consider four gender conversions (F: female, M: male): F2F, F2M, M2F, M2M. We also consider four cross-language conversions (E: English, Z: Chinese): E2E, E2Z, Z2E, Z2Z. The voice conversion samples are available at <https://shamidreza.github.io/is18samples>

4.3. Visualizing embeddings

In this experiment, we investigate the speaker embeddings Z_2 by visualizing them in Figure 2. For visualizing the speaker embeddings, we use the 10 test speakers from TIMIT test set (red data points for males and blue for female) and 10 test speakers from THCHS-30 (orange/greenish data points for males and light blue for female). We also use VAE models trained on TIMIT (top), CH corpus (middle), and TIMIT+CH corpus (bottom). In Figures 2, we show the speaker embedding from 1 sentence in the left subplots and from 5 sentences in the right subplots, where the 2D plot of the speaker embeddings (computed using PCA) are shown. In all subplots, the female and male embedding cluster locations are clearly separated. Furthermore, the plot shows that the speaker embeddings of unique speakers fall near the same location. Although when 5 utterances are used to compute the embedding value, the variation is visibly less compared to when merely one sentence is used. This shows sensitiveness of the speaker embedding computed from the model to sentence variations. Also it is interesting to note that when both TIMIT+CH corpus are used for training, the speaker embeddings are further apart suggesting a better model property. One phenomenon that we notice is that the speaker embeddings for different languages and gender fall to different locations. This shows the embeddings are still language dependent, which might suggest the network learn to use the phonetic information to learn some language embedding within Z_2 as an additional factor. Furthermore, we investigate the phonetic context embedding Z_1 for a sentence for four test speakers on TIMIT-trained VAE. The phonetic context matrix over the computed utterances (compressed using PCA) is shown in Figures 3. Ideally, we want the matrices should be close to each other since the phonetic context embedding is supposed to be speaker-independent. The figure show the closeness of the embeddings at the similar time frames. There is still some minor discrepancy between the embeddings which shows room for further improvement of model architecture and/or larger speech corpus.

4.4. Effect of training data size

We investigate the effect of VC training data size on the performance of the system. In order to be able to do objective test using mel-CD [4], we require parallel data from the speakers. We use 20 parallel CMU-acric utterance from each speaker for

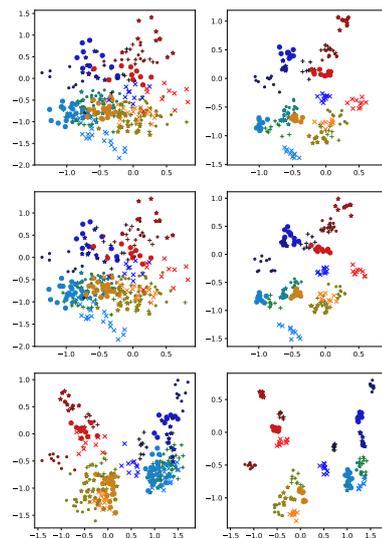


Figure 2: Visualization of speaker embedding. Each point represents single utterance and different colors represent different speaker/languages; blueish dots are English females and light blueish are Chinese females; and reddish dots are English males and orange dots are Chinese males. See Section 4.3 for details.

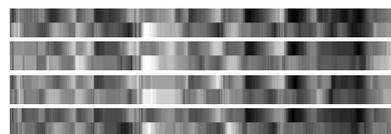


Figure 3: Visualization of phonetic context embedding sequence of the sentence "She had your dark suit in greasy wash water all year" aligned to each other for two female speakers (top) and two female speakers (bottom). The embeddings are transformed to 2D using PCA.

computing the objective score. We vary non-parallel sentence numbers from source and target speaker that is used to compute the speaker embeddings. The results are shown in Figure 4. As can be seen, VAE performs better with less than 10 sentences, however, with more than 10 sentences, GMM starts achieving lower mel-CD. This might be due to VAE only having one degree of freedom (speaker identity vector) to convert the voice, compared to GMM which is able to use all the training data to adapt the background GMM model to better match the speaker data distribution.

4.5. Subjective evaluation

To subjectively evaluate voice conversion performance, we performed two perceptual tests. The first test measured speech quality, designed to answer the question "how natural does the converted speech sound"?, and the second test measured speaker similarity, designed to answer the question "how accurate does the converted speech mimic the target speaker"?. The listening experiments were carried out using Amazon Mechanical Turk, with participants who had approval ratings of at least 90% and were located in North America. Both perceptual tests used three trivial-to-judge trials, added to the experiment to exclude unreliable listeners from statistical analysis. No listeners were flagged as unreliable in our experiments. In this subjective experiment, we focus on VAE train on TIMIT. We provide samples of VAE trained on TIMIT, CH, TIMIT+CH and VCTK in the samples webpage demo as well. In informal listening tests, we found out that VAE trained on TIMIT performs better

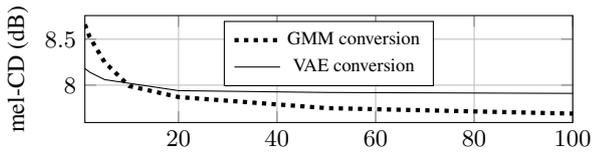


Figure 4: Effects of varying number of training sentences from 1 to 100

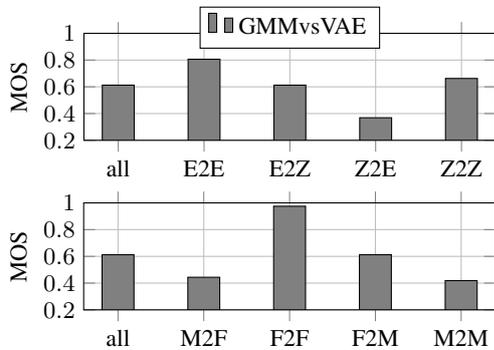


Figure 5: Speech Quality average score with gender and language break-down. Positive scores favor VAE. (confidence intervals for all is close to 0.13, and all scores are statistically significant)

than VAE trained on VCTK. Also VAE trained on TIMIT+CH generate better quality than on TIMIT or CH alone.

4.5.1. Speech quality

To evaluate the speech quality of the converted utterances, we conducted a Comparative Mean Opinion Score (CMOS) test. In this test, listeners heard two stimuli A and B with the same content, generated using the same source speaker, but in two different processing conditions, and were then asked to indicate whether they thought B was better or worse than A, using a five-point scale comprised of +2 (much better), +1 (somewhat better), 0 (same), -1 (somewhat worse), -2 (much worse). We randomized the order of stimulus presentation, both the order of A and B, as well as the order of the comparison pairs. We utilized three processing conditions: GMM, VAE-STFT, VAE. We ran two separate experiments. First, to assess the effect of using World vocoder instead of STFT, we directly compared VAE-STFT vs. VAE. We only limited this experiment to English to English conversion. The experiment was administered to 40 listeners with each listener judging 16 sentence pairs. The results shows a very significant preference of VAE over VAE-STFT, achieving $+1.25 \pm 0.12$ mean score towards VAE. We performed planned one-sample t-tests with a mean of zero and achieved $p < 0.0001$. Second, we assessed the VC approach effect by directly comparing GMM vs. VAE utterances. The experiment was administered to 40 listeners with each listener judging 80 sentence pairs. The results shows VAE has a statistically significant quality improvement over GMM, achieving $+0.61 \pm 0.14$ mean score towards VAE. The language-breakdown of the results are shown in Figure 5. We performed planned one-sample t-tests with a mean of zero and achieved $p < 0.05$ for all language and gender conversion pairs separately, showing statistically significant improvements for all break-downs of gender and language. The Z2E conversion is achieving lowered quality compared to other conversion pairs. We speculate the reason is the slight noise present in THCHS-30 recordings which cause some distortion during vocoding.

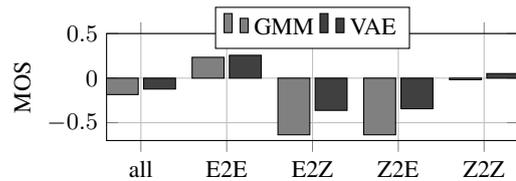


Figure 6: Speech Similarity average score with language break-down. Positive scores are desirable. (Comparison of scores between GMM vs. VAE do not show statistical significance)

4.5.2. Speaker similarity

To evaluate the speaker similarity of the converted utterances, we conducted a same-different speaker similarity test [31]. In this test, listeners heard two stimuli A and B with different content, and were then asked to indicate whether they thought that A and B were spoken by the same, or by two different speakers, using a five-point scale comprised of +2 (definitely same), +1 (probably same), 0 (unsure), -1 (probably different), and -2 (definitely different). One of the stimuli in each pair was created by one of the two conversion methods, and the other stimulus was a purely MCEP-vocoded condition, used as the reference speaker. The listeners were explicitly instructed to disregard the language of the stimuli and merely judge based on the fact whether they think the utterances are from the same speaker regardless of the language. Half of all pairs were created with the reference speaker identical to the target speaker of the conversion (expecting listeners to reply “same”, ideally); the other half were created with the reference speaker being the same gender, but not identical to the target speaker of the conversion (expecting listeners to reply different). We only report “same” scores. The experiment was administered to 40 listeners, with each listener judging 64 sentence pairs. The results are shown in Figure 6. The results show GMM and VAE achieving -0.18 ± 0.15 and -0.12 ± 0.16 , respectively. We did not find any statistical significance between GMM vs VAE systems for average, or any of the language/gender conversion break-downs of the stimuli. For both VAE and GMM, we noticed that for E2E case, we achieve the highest average score and the only case that is able to transform identity, achieving $P < 0.05$ in one-sample t-test compared to chance. This is reasonable given the training was done only on an English corpus. Furthermore, Z2Z achieves higher score compared to E2Z and Z2E. This might be due to the listener’s bias toward not rating different language utterances as high score as same language utterances.

5. Conclusions

We proposed to exploit FHVAE model for challenging non-parallel and cross-lingual voice conversion, even with very small number of training utterances such as only one target speaker’s utterance. We investigate the importance of speech representations and found that World vocoder outperformed STFT which was used in [1] in experimental evaluation, both speech quality and similarity. We also examined the effect of the size of training utterances from target speaker for VC, and our approach outperformed baseline with less than 10 sentences, and achieve reasonable performance even with only one training utterance. In the subjective tests, our approach achieved significantly better results than both VAE-STFT and GMM in speech quality, and outperformed VAE-STFT and comparable to GMM in speech similarity. As future work, we are interested in and working on joint end-to-end learning with Wavenet [17] or Wavenet-vocoder [32, 33], and also building models trained on multi-language corpora with or without explicit modeling of different languages such as providing language coding.

6. References

- [1] W.-N. Hsu, Y. Zhang, and J. Glass, "Unsupervised learning of disentangled and interpretable representations from sequential data," in *Advances in neural information processing systems*, 2017, pp. 1876–1887.
- [2] Y. Stylianou, "Voice transformation: a survey," in *Acoustics, Speech and Signal Processing, 2009. ICASSP 2009. IEEE International Conference on*. IEEE, 2009, pp. 3585–3588.
- [3] S. H. Mohammadi and A. Kain, "An overview of voice conversion systems," *Speech Communication*, vol. 88, pp. 65–82, 2017.
- [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. Desai, A. W. Black, B. Yegnanarayana, and K. Prahallad, "Spectral mapping using artificial neural networks for voice conversion," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 5, pp. 954–964, 2010.
- [6] Z. Wu, T. Kinnunen, E. S. Chng, and H. Li, "Mixture of factor analyzers using priors from non-parallel speech for voice conversion," *IEEE Signal Processing Letters*, vol. 19, no. 12, pp. 914–917, 2012.
- [7] 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*. IEEE, 2016, pp. 1–6.
- [8] P. Song, W. Zheng, and L. Zhao, "Non-parallel training for voice conversion based on adaptation method," in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*. IEEE, 2013, pp. 6905–6909.
- [9] 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.
- [10] T. Kinnunen, L. Juvela, P. Alku, and J. Yamagishi, "Non-parallel voice conversion using i-vector plda: Towards unifying speaker verification and transformation," in *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*. IEEE, 2017, pp. 5535–5539.
- [11] S. H. Mohammadi and A. Kain, "Siamese autoencoders for speech style extraction and switching applied to voice identification and conversion," *Proceedings of Interspeech*, pp. 1293–1297, 2017.
- [12] B. Ramani, M. A. Jeeva, P. Vijayalakshmi, and T. Nagarajan, "A multi-level gmm-based cross-lingual voice conversion using language-specific mixture weights for polyglot synthesis," *Circuits, Systems, and Signal Processing*, vol. 35, no. 4, pp. 1283–1311, 2016.
- [13] M. Charlier, Y. Ohtani, T. Toda, A. Moinet, and T. Dutoit, "Cross-language voice conversion based on eigenvoices," *Proceedings of Interspeech*, 2009.
- [14] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [15] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [16] A. v. d. Oord, N. Kalchbrenner, and K. Kavukcuoglu, "Pixel recurrent neural networks," *arXiv preprint arXiv:1601.06759*, 2016.
- [17] A. Van Den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," *arXiv preprint arXiv:1609.03499*, 2016.
- [18] T. Kaneko and H. Kameoka, "Parallel-data-free voice conversion using cycle-consistent adversarial networks," *arXiv preprint arXiv:1711.11293*, 2017.
- [19] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," *arXiv preprint arXiv:1703.10593*, 2017.
- [20] T. D. Kulkarni, W. F. Whitney, P. Kohli, and J. Tenenbaum, "Deep convolutional inverse graphics network," in *Advances in Neural Information Processing Systems*, 2015, pp. 2539–2547.
- [21] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever, and P. Abbeel, "Infogan: Interpretable representation learning by information maximizing generative adversarial nets," in *Advances in Neural Information Processing Systems*, 2016, pp. 2172–2180.
- [22] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, "beta-vae: Learning basic visual concepts with a constrained variational framework," 2016.
- [23] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Advances in neural information processing systems*, 2014, pp. 3104–3112.
- [24] J. S. Garofolo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, and D. S. Pallett, "Darpa timit acoustic-phonetic continuous speech corpus cd-rom. nist speech disc 1-1.1," *NASA STI/Recon technical report n*, vol. 93, 1993.
- [25] D. Wang and X. Zhang, "Thchs-30: A free chinese speech corpus," *arXiv preprint arXiv:1512.01882*, 2015.
- [26] C. Veaux, J. Yamagishi, K. MacDonald *et al.*, "Cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit," 2017.
- [27] J. Kominek and A. W. Black, "The cmu arctic speech databases," in *Fifth ISCA Workshop on Speech Synthesis*, 2004.
- [28] M. Morise, F. Yokomori, and K. Ozawa, "World: a vocoder-based high-quality speech synthesis system for real-time applications," *IEICE TRANSACTIONS on Information and Systems*, vol. 99, no. 7, pp. 1877–1884, 2016.
- [29] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [30] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [31] A. B. Kain, "High resolution voice transformation," 2001.
- [32] A. Tamamori, T. Hayashi, K. Kobayashi, K. Takeda, and T. Toda, "Speaker-dependent wavenet vocoder," in *Proceedings of Interspeech*, 2017, pp. 1118–1122.
- [33] T. Hayashi, A. Tamamori, K. Kobayashi, K. Takeda, and T. Toda, "An investigation of multi-speaker training for wavenet vocoder," in *Automatic Speech Recognition and Understanding Workshop (ASRU), 2017 IEEE*. IEEE, 2017, pp. 712–718.