
Investigation of using disentangled and interpretable representations with language conditioning for cross-lingual voice conversion

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the problem of cross-lingual voice conversion in non-parallel speech
2 corpora and one-shot learning setting. Most prior work require either parallel
3 speech corpora or enough amount of training data from a target speaker. However,
4 we convert an arbitrary sentences of an arbitrary source speaker to target speaker's
5 given only one target speaker training utterance. To achieve this, we formulate the
6 problem as learning disentangled speaker-specific and context-specific represen-
7 tations and follow the idea of [1] which uses Factorized Hierarchical Variational
8 Autoencoder (FHVAE). After training FHVAE on multi-speaker training data,
9 given arbitrary source and target speakers' utterance, we estimate those latent
10 representations and then reconstruct the desired utterance of converted voice to that
11 of target speaker. We use multi-language speech corpus to learn a universal model
12 that works for all of the languages. We investigate the use of a one-hot language
13 embedding to condition the model on the language of the utterance being queried
14 and show the effectiveness of the approach. We also investigate the effect of using
15 or not using the language conditioning. Furthermore, we visualize the embeddings
16 of the different languages and sexes. Finally, in the subjective tests, for one lan-
17 guage and cross-lingual voice conversion, our approach achieved moderately better
18 or comparable results compared to the baseline in speech quality and similarity.

19 1 Introduction

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

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

- 36 • We build a voice model which is trained on utterances from 5 different languages to let the
37 model observe as much speaker and phonetic variations as possible.
- 38 • We conduct cross-lingual voice conversion experiments and our approach achieved mod-
39 erately better or comparable results than baselines in speech quality and similarity in the
40 subjective tests.
- 41 • We examine the effect of using additional one-hot embedding along with speaker embedding
42 that determines the input utterance language.

43 2 Related Work

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

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

62 Recently, deep generative models have been applied and successful for unsupervised learning tasks,
63 and include Variational Autoencoder (VAE) [14], Generative Adversarial Networks (GAN) [15], and
64 auto-regressive models [16, 17]. Among them, VAE can infer latent codes from data and generate
65 data from them by jointly learning inference and generative networks, and VAE has been also applied
66 for voice conversion [7, 9]. However, in their models, speaker identities are not inferred from data and
67 instead required to be known in model training time. GAN has been also exploited for non-parallel
68 voice conversion [18] with the cycle consistency constraint [19], but it still has the limitation that it
69 needs to know the target speaker in training time and be trained for each target.

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

77 3 Model

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

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

83 where D_{KL} is Kullback-Leibler divergence.

84 However, VAE considers no structure for latent variable z . Assuming structure for z could be
85 beneficial to exploit the inherent structures in data. Here we describe Factorized Hierarchical

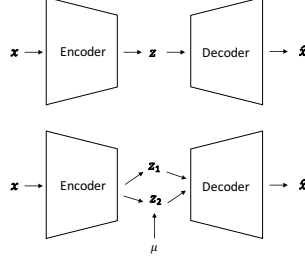


Figure 1: Structures of Variation Autoencoder (upper) and Factorized Hierarchical Variational Autoencoder (lower).

86 Variational Autoencoder proposed by Hsu et al [1]. Let a dataset D consist of N_{seg} i.i.d. sequences
 87 X^i . For each sequence X^i , it consists of N_{seg}^i $X^{i,j}$ observation segments. Then we define factorized
 88 latent variables of latent segment variable $Z_1^{i,j}$ and latent sequence variable $Z_2^{i,j}$. In the context of
 89 voice conversion, $Z_1^{i,j}$ is responsible for generating phonetic contexts and $Z_2^{i,j}$ is for speaker identity.
 90 When generating data $X^{i,j}$, we first sample $Z_2^{i,j}$ from isotropic Gaussian centered at μ^i shared for
 91 the entire sequence, and also $Z_1^{i,j}$ independently. Then we generate $X^{i,j}$ conditioned on $Z_1^{i,j}$ and
 92 $Z_2^{i,j}$. Thus, joint probability with a sequence X^i is:

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

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

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

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

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

97 where $\tilde{\mu}^i$ is the posterior mean of μ^i . Please refer to Hsu et al. [1] for more details. Additionally,
 98 Hsu et al. also proposed discriminative segment variational lower bound to encourage Z_2^i to be more
 99 sequence-specific by adding the additional term of inferring the sequence index i from $Z_2^{i,j}$. For our
 100 experiments, we exploit this FHVAE model and sequence-to-sequence model [23] as the structure of
 101 encoder-decoder for sequential data. We propose adding an input language embedding to the input of
 102 the model. This language embedding will be used to determine the input utterance language using a
 103 one-hot representation of the in-training languages.

104 For performing the voice conversion, we compute the average Z_2 from the training utterance(s) of
 105 source and target speakers. For a given input utterance, we compute Z_1 and Z_2 of the input utterance.
 106 There are two ways to perform voice conversion. First, we can replace Z_2 values of the source speaker
 107 with the average Z_2 from the target speaker. This approach resulted in too muffled generated result.
 108 Second, we compute a difference vector between source and target average $Z_2^{diff} = Z_2^{trg} - Z_2^{src}$.

109 This difference vector is added to Z_2 from the input utterance as $Z_2^{converted} = Z_2 + Z_2^{diff}$ and then
110 decoded using FHVAE to achieve the speech features. In an informal listening test, we decided to the
111 second approach since it resulted in significantly higher quality generated speech.

112 4 Experiments

113 4.1 Datasets

114 We used the TIMIT corpus [24] which is a multi-speaker speech corpus as the training data for
115 FHVAE model. We used the training speakers as suggested by the corpus to train the model. For
116 English test speakers, we select speakers from TIMIT testing part of the corpus. We also use a
117 proprietary Chinese speech corpus (hereon referred to as CH) with 5200 speakers each uttering one
118 sentence. We use Microsoft’s Indian Language Speech Corpus for Indian language. We also use
119 proprietary Korean and Japanese multi-speaker speech corpora. Additionally, we consider using the
120 combination of all languages corpus for training the model. For Korean, Japanese, and Indian corpora,
121 we randomly exclude 10 percent of the speakers from each corpus for training purposes. For Chinese
122 test speakers, we utilize speakers from the THCHS-30 speech corpus [25]. To observe the effect of
123 having more utterances per speaker but less speakers we also train the model on VCTK corpus [26].
124 Finally, for objective testing (which requires availability of parallel data), we utilized four CMU-arctic
125 voices (BDL, SLT, RMS, CLB)[27]. As speech features, we used 40th-order MCEPs (excluding the
126 energy coefficient, dimensionality $D=39$), extracted using the World toolkit [28] with a 5ms frame
127 shift. All audio files are transformed to 16kHz and 16 bit before any analysis.

128 4.2 Experimental setting

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

135 From now on, we use the abbreviation VAE for FHVAE model. In our experiments, we consider
136 two models: VAE-UNC (unconditioned) [31] and VAE-CND (conditioned) which mean models
137 trained either with or without the language conditioning input. We consider four gender conversions
138 (F: female, M: male): F2F, F2M, M2F, M2M. We also consider 25 cross-language conversions: all
139 permutations of 5 languages considered as source and target. The voice conversion samples are
140 available at: <https://shamidreza.github.io/nips2018samples>

141 4.3 Visualizing embeddings

142 In this experiment, we investigate the speaker embeddings Z_2 by visualizing them in Figure 2. For
143 visualizing the speaker embeddings, we use the 10 test speakers (5 male, 5 female) from each language
144 test corpus. We use VAE-UNC and VAE-CND. In Figures 2, we show the speaker embeddings
145 computed from 1 utterance where the 2D plot of the speaker embeddings (computed using PCA)
146 are shown. In all subplots, the female and male embedding cluster locations are clearly separated.
147 Furthermore, the plot shows that the speaker embeddings of unique speakers fall near the same
148 location. One phenomenon that we notice is that the speaker embeddings for different languages and
149 gender fall to different locations for VAE-UNC, however, they fall closer to each other in VAE-CND.
150 This might be due to the conditioning on language improving the representation ability of the model.
151 Furthermore, we investigate the phonetic context embedding Z_1 for a sentence for four English test
152 speakers on VAE-UNC. The phonetic context matrix over the computed utterances (compressed using
153 PCA) is shown in Figures 3. Ideally, we want the matrices should be close to each other since the
154 phonetic context embedding is supposed to be speaker-independent. The figure show the closeness
155 of the embeddings at the similar time frames. There is still some minor discrepancy between the
156 embeddings which shows room for further improvement of model architecture and/or larger speech
157 corpus.

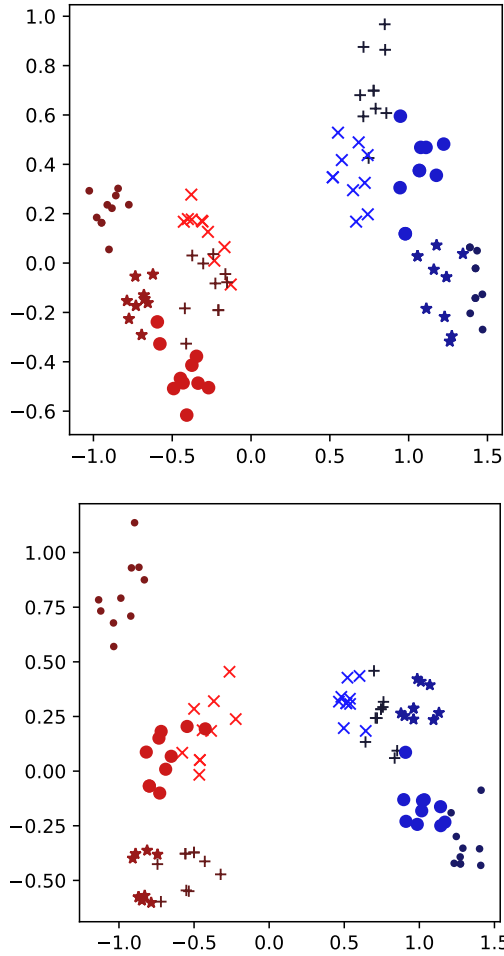


Figure 2: Visualization of speaker embeddings: unconditioned (top) versus conditioned (bottom). Red represents female speakers and blue represents male speakers. Each dot type represents a language

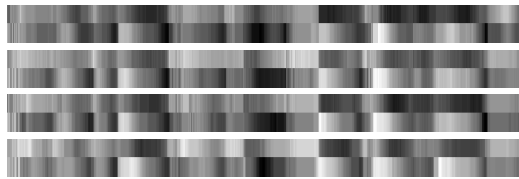


Figure 3: Visualization of phonetic context embedding sequence of a sentences aligned to each other for four English speakers. The embeddings are transformed to 2D using PCA.

158 **4.4 Subjective evaluation**

159 To subjectively evaluate voice conversion performance, we performed two perceptual tests. The
 160 first test measured speech quality, designed to answer the question “how natural does the converted
 161 speech sound”?, and the second test measured speaker similarity, designed to answer the question
 162 “how accurate does the converted speech mimic the target speaker”?. The listening experiments were
 163 carried out using Amazon Mechanical Turk, with participants who had approval ratings of at least
 164 90% and were located in North America. Both perceptual tests used three trivial-to-judge trials,

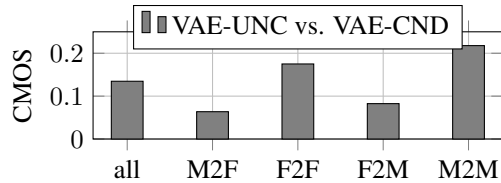


Figure 4: Speech Quality average score with gender and language break-down. Positive scores favor VAE-CND. (confidence intervals for all is close to 0.14)

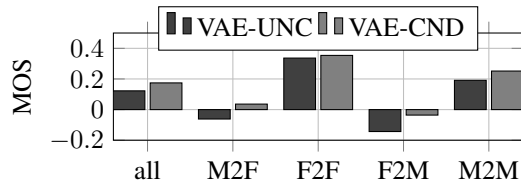


Figure 5: Speech Similarity average score with conversion break-down. Positive scores are desirable.

165 added to the experiment to exclude unreliable listeners from statistical analysis. No listeners were
 166 flagged as unreliable in our experiments.

167 4.4.1 Speech quality

168 To evaluate the speech quality of the converted utterances, we conducted a Comparative Mean Opinion
 169 Score (CMOS) test. In this test, listeners heard two stimuli A and B with the same content, generated
 170 using the same source speaker, but in two different processing conditions, and were then asked to
 171 indicate whether they thought B was better or worse than A, using a five-point scale comprised
 172 of +2 (much better), +1 (somewhat better), 0 (same), -1 (somewhat worse), -2 (much worse). We
 173 randomized the order of stimulus presentation, both the order of A and B, as well as the order of the
 174 comparison pairs. We utilized two processing conditions: VAE-UNC, VAE-CND. We assessed the
 175 VC approach effect by directly comparing VAE-UNC vs. VAE-CND utterances. The experiment was
 176 administered to 50 listeners with each listener judging 100 sentence pairs. We achieved $+0.15 \pm 0.14$
 177 mean score towards VAE-CND. Although this is a positive difference, we did not find statistically
 178 significant difference between the quality of VAE-UNC and VAE-CND. The language-breakdown of
 179 the results are shown in Figure 4.

180 4.4.2 Speaker similarity

181 To evaluate the speaker similarity of the converted utterances, we conducted a same-different speaker
 182 similarity test [32]. In this test, listeners heard two stimuli A and B with different content, and
 183 were then asked to indicate whether they thought that A and B were spoken by the same, or by two
 184 different speakers, using a five-point scale comprised of +2 (definitely same), +1 (probably same), 0
 185 (unsure), -1 (probably different), and -2 (definitely different). One of the stimuli in each pair was
 186 created by one of the two conversion methods, and the other stimulus was a purely MCEP-vocoded
 187 condition, used as the reference speaker. The listeners were explicitly instructed to disregard the
 188 language of the stimuli and merely judge based on the fact whether they think the utterances are from
 189 the same speaker regardless of the language. Half of all pairs were created with the reference speaker
 190 identical to the target speaker of the conversion (expecting listeners to reply "same", ideally); the
 191 other half were created with the reference speaker being the same gender, but not identical to the
 192 target speaker of the conversion (expecting listeners to reply different). We only report "same" scores.
 193 The experiment was administered to 50 listeners, with each listener judging 100 sentence pairs. The
 194 results are shown in Figure 5. We did find a consistent improvement of VAE-CND performance over
 195 VAE-UNC, however these differences were not statistically significant.

196 5 Conclusions

197 We proposed to exploit FHVAE model for challenging non-parallel and cross-lingual voice conversion,
198 even with very small number of training utterances such as only one target speaker’s utterance. We use
199 multi-language corpus to learn disentangled representations from speech. We also introduce a one-hot
200 language embedding to the model in order to improve the model performance. We perform several
201 visualizations to show the effect of the model. In the subjective tests, we found some improvement in
202 the quality and similarity performance of the system, although the differences were not statistically
203 significant.

204 References

- 205 [1] W.-N. Hsu, Y. Zhang, and J. Glass, “Unsupervised learning of disentangled and interpretable
206 representations from sequential data,” in *Advances in neural information processing systems*,
207 pp. 1876–1887, 2017.
- 208 [2] Y. Stylianou, “Voice transformation: a survey,” in *Acoustics, Speech and Signal Processing,
209 2009. ICASSP 2009. IEEE International Conference on*, pp. 3585–3588, IEEE, 2009.
- 210 [3] S. H. Mohammadi and A. Kain, “An overview of voice conversion systems,” *Speech Communi-
211 cation*, vol. 88, pp. 65–82, 2017.
- 212 [4] T. Toda, A. W. Black, and K. Tokuda, “Voice conversion based on maximum-likelihood
213 estimation of spectral parameter trajectory,” *IEEE Transactions on Audio, Speech, and Language
214 Processing*, vol. 15, no. 8, pp. 2222–2235, 2007.
- 215 [5] S. Desai, A. W. Black, B. Yegnanarayana, and K. Prahallad, “Spectral mapping using artificial
216 neural networks for voice conversion,” *IEEE Transactions on Audio, Speech, and Language
217 Processing*, vol. 18, no. 5, pp. 954–964, 2010.
- 218 [6] Z. Wu, T. Kinnunen, E. S. Chng, and H. Li, “Mixture of factor analyzers using priors from
219 non-parallel speech for voice conversion,” *IEEE Signal Processing Letters*, vol. 19, no. 12,
220 pp. 914–917, 2012.
- 221 [7] C.-C. Hsu, H.-T. Hwang, Y.-C. Wu, Y. Tsao, and H.-M. Wang, “Voice conversion from non-
222 parallel corpora using variational auto-encoder,” in *Signal and Information Processing Association
223 Annual Summit and Conference (APSIPA), 2016 Asia-Pacific*, pp. 1–6, IEEE, 2016.
- 224 [8] P. Song, W. Zheng, and L. Zhao, “Non-parallel training for voice conversion based on adaptation
225 method,” in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International
226 Conference on*, pp. 6905–6909, IEEE, 2013.
- 227 [9] C.-C. Hsu, H.-T. Hwang, Y.-C. Wu, Y. Tsao, and H.-M. Wang, “Voice conversion from unaligned
228 corpora using variational autoencoding wasserstein generative adversarial networks,” *arXiv
229 preprint arXiv:1704.00849*, 2017.
- 230 [10] T. Kinnunen, L. Juvela, P. Alku, and J. Yamagishi, “Non-parallel voice conversion using i-vector
231 plda: Towards unifying speaker verification and transformation,” in *Acoustics, Speech and
232 Signal Processing (ICASSP), 2017 IEEE International Conference on*, pp. 5535–5539, IEEE,
233 2017.
- 234 [11] S. H. Mohammadi and A. Kain, “Siamese autoencoders for speech style extraction and switching
235 applied to voice identification and conversion,” *Proceedings of Interspeech*, pp. 1293–1297,
236 2017.
- 237 [12] B. Ramani, M. A. Jeeva, P. Vijayalakshmi, and T. Nagarajan, “A multi-level gmm-based
238 cross-lingual voice conversion using language-specific mixture weights for polyglot synthesis,”
239 *Circuits, Systems, and Signal Processing*, vol. 35, no. 4, pp. 1283–1311, 2016.
- 240 [13] M. Charlier, Y. Ohtani, T. Toda, A. Moinet, and T. Dutoit, “Cross-language voice conversion
241 based on eigenvoices,” *Proceedings of Interspeech*, 2009.

- 242 [14] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” *arXiv preprint*
243 *arXiv:1312.6114*, 2013.
- 244 [15] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and
245 Y. Bengio, “Generative adversarial nets,” in *Advances in neural information processing systems*,
246 pp. 2672–2680, 2014.
- 247 [16] A. v. d. Oord, N. Kalchbrenner, and K. Kavukcuoglu, “Pixel recurrent neural networks,” *arXiv*
248 *preprint arXiv:1601.06759*, 2016.
- 249 [17] A. Van Den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner,
250 A. Senior, and K. Kavukcuoglu, “Wavenet: A generative model for raw audio,” *arXiv preprint*
251 *arXiv:1609.03499*, 2016.
- 252 [18] T. Kaneko and H. Kameoka, “Parallel-data-free voice conversion using cycle-consistent adver-
253 sarial networks,” *arXiv preprint arXiv:1711.11293*, 2017.
- 254 [19] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, “Unpaired image-to-image translation using
255 cycle-consistent adversarial networks,” *arXiv preprint arXiv:1703.10593*, 2017.
- 256 [20] T. D. Kulkarni, W. F. Whitney, P. Kohli, and J. Tenenbaum, “Deep convolutional inverse graphics
257 network,” in *Advances in Neural Information Processing Systems*, pp. 2539–2547, 2015.
- 258 [21] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever, and P. Abbeel, “Infogan: Inter-
259 pretable representation learning by information maximizing generative adversarial nets,” in
260 *Advances in Neural Information Processing Systems*, pp. 2172–2180, 2016.
- 261 [22] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Ler-
262 chner, “beta-vae: Learning basic visual concepts with a constrained variational framework,”
263 2016.
- 264 [23] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,”
265 in *Advances in neural information processing systems*, pp. 3104–3112, 2014.
- 266 [24] J. S. Garofolo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, and D. S. Pallett, “Darpa timit acoustic-
267 phonetic continous speech corpus cd-rom. nist speech disc 1-1.1,” *NASA STI/Recon technical*
268 *report n*, vol. 93, 1993.
- 269 [25] D. Wang and X. Zhang, “Thchs-30: A free chinese speech corpus,” *arXiv preprint*
270 *arXiv:1512.01882*, 2015.
- 271 [26] C. Veaux, J. Yamagishi, K. MacDonald, *et al.*, “Cstr vctk corpus: English multi-speaker corpus
272 for cstr voice cloning toolkit,” 2017.
- 273 [27] J. Kominek and A. W. Black, “The cmu arctic speech databases,” in *Fifth ISCA Workshop on*
274 *Speech Synthesis*, 2004.
- 275 [28] M. Morise, F. Yokomori, and K. Ozawa, “World: a vocoder-based high-quality speech synthesis
276 system for real-time applications,” *IEICE TRANSACTIONS on Information and Systems*, vol. 99,
277 no. 7, pp. 1877–1884, 2016.
- 278 [29] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9,
279 no. 8, pp. 1735–1780, 1997.
- 280 [30] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint*
281 *arXiv:1412.6980*, 2014.
- 282 [31] S. H. Mohammadi and T. Kim, “Investigation of using disentangled and interpretable represen-
283 tations for one-shot cross-lingual voice conversion,” *Proc. Interspeech*, 2018.
- 284 [32] A. B. Kain, “High resolution voice transformation,” 2001.