

# Unsupervised End-to-End Learning of Discrete Linguistic Units for Voice Conversion

Andy T. Liu<sup>1</sup>, Po-chun Hsu<sup>1</sup>, Hung-yi Lee<sup>1</sup>

<sup>1</sup>College of Electrical Engineering and Computer Science, National Taiwan University

{r07942089, r07942095, hungyilee}@ntu.edu.tw

## Abstract

We present an unsupervised end-to-end training scheme where we discover discrete subword units from speech without using any labels. The discrete subword units are learned under an ASR-TTS autoencoder reconstruction setting, where an ASR-Encoder is trained to discover a set of common linguistic units given a variety of speakers, and a TTS-Decoder trained to project the discovered units back to the designated speech. We propose a discrete encoding method, Multilabel-Binary Vectors (MBV), to make the ASR-TTS autoencoder differentiable. We found that the proposed encoding method offers automatic extraction of speech content from speaker style, and is sufficient to cover full linguistic content in a given language. Therefore, the TTS-Decoder can synthesize speech with the same content as the input of ASR-Encoder but with different speaker characteristics, which achieves voice conversion (VC). We further improve the quality of VC using adversarial training, where we train a TTS-Patcher that augments the output of TTS-Decoder. Objective and subjective evaluations show that the proposed approach offers strong VC results as it eliminates speaker identity while preserving content within speech. In the ZeroSpeech 2019 Challenge, we achieved outstanding performance in terms of low bitrate.

**Index Terms:** acoustic unit discovery, voice conversion, speech disentangled representation, adversarial training

## 1. Introduction

Despite that human speech inherently carries linguistic features that represent textual information, modern text-to-speech training pipelines still require parallel speech and text transcription pairs [1][2][3]. Parallel speech and transcripts may not always be available, and is costly to acquire, however human speech alone can be easily gathered. In this work, we focus on utilizing the advantage of unlabeled speech to discover discrete linguistic units, where machine learns to uncover the linguistic features hidden in human utterance without any supervision. We use these discovered linguistic units for voice conversion (VC) and achieved outstanding results.

Embedding audio signals into latent representations has been a well studied practice [4][5][6][7][8][9][10]. Former studies have also attempted to encode speech content into various representations for voice conversion, including the use of disentangled autoencoders [11], VAEs [12] or GANs [13]. Continuous vectors are the most common approach, however when it comes to encoding speech content, they may not be the best choice. As they are unlike the discrete phonemes that we often used to represent human language. Previous works [14][15][16] also attempt to learn discrete representations from audio, however they did not apply the learned units for VC.

In VC, the state-of-the-art approach [11] requires additional loss in the framework of GAN to guarantee the disentanglement

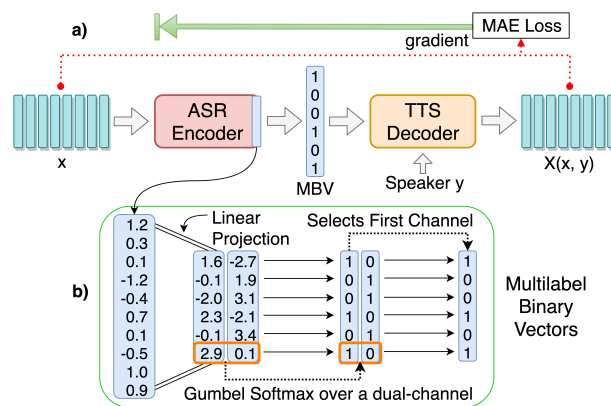


Figure 1: The ASR-TTS autoencoder framework where MBV are learned as discrete linguistic representations.

of learned encodings, whereas the proposed approach naturally possesses the direct ability to separate speaker style from speech content. In [17], VC is achieved through multiple losses together with a K-way quantized embedding space and autoregressive WaveNet, on the other hand, the proposed encoding space does not require additional training constraints.

In this work, we use an ASR-TTS autoencoder to discover discrete linguistic units from speech, without any alignment, text label, or parallel data provided. The ASR-Encoder learns to encode speech from different speakers to a common set of small linguistic symbols. These finite set of linguistic symbols are represented by Multilabel-Binary Vectors (MBV), vectors consist of arbitrary number of zeros and ones. The proposed MBV method is differentiable, hence allowing backpropagation of gradients and end-to-end training of an autoencoder reconstruction setting. While the ASR-Encoder learns a many-to-one mapping from speech to discrete subword units, a TTS-Decoder learns a one-to-many mapping from discrete subword units back to speech. The discrete nature of MBV allows it to innately separate linguistic content from speech, removing speaker characteristics. Given an utterance of a source speaker, we were able to encode its speech content using the ASR-Encoder, and perform voice conversion to generate speech with the same linguistic content but style of a target speaker using the TTS-Decoder.

In training, the TTS-Decoder is always given an encoding from an arbitrary speaker and is trained to decode it back to the original speech of that particular speaker. During inference time, the TTS-Decoder has to take encodings of a source speaker and decode to a target speaker, an encoding-speaker pair that it never observed during training. Although speech conversion is already feasible, we propose to use additional adversarial training to compensate the training-testing inconsistency. Under the WGAN [18] setting, we train a TTS-Patcher

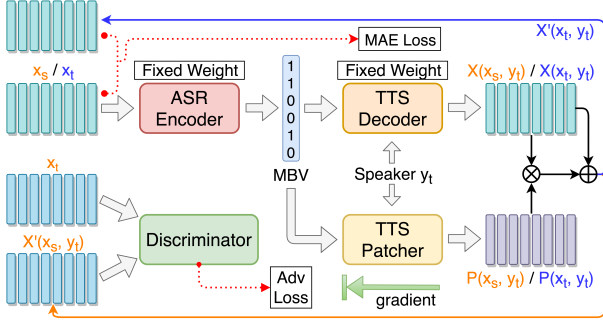


Figure 2: Target guided adversarial learning to improve voice conversion performance. We annotate the adversarial step in orange and the target guided step in blue.

in place of a generator. The TTS-Patcher learns to generate a mask that augments the output of TTS-Decoder. Furthermore, we use a target driven reconstruction loss to guide the generator’s update. As a result, the quality of voice conversion performance is improved. In the ZeroSpeech 2019 Challenge [19], we achieve 2<sup>nd</sup> place in terms of low bitrate under a strong dimension constraint on the Surprise Language [20][21] leaderboard. We further show that the proposed method is capable of generating high quality intelligible speech when the dimension constraint is removed.

## 2. Proposed method

### 2.1. Discrete linguistic units discovery

We present an unsupervised end-to-end learning framework for distinct linguistic units discovery. In this stage we learn a common set of discrete encodings for all the speakers. Let  $x \in X$  be an acoustic feature sequence where  $X$  is a set of all such sequences from all source and target speakers, and  $x$  is a fixed length- $r$  segment randomly sampled from  $X$ . Let  $y \in Y$  be a speaker where  $Y$  is the group of all source and target speakers who produce the sequence collection  $X$ . In a training pair  $(x, y) \in (X, Y)$ , the sequence  $x$  is produced by speaker  $y$ .

#### 2.1.1. Learning Multilabel-Binary Vectors

In Figure 1, an ASR-TTS autoencoder framework is used to learn discrete linguistic units. The ASR-Encoder is trained to map input acoustic feature sequence  $x \in X$  to a latent discrete encoding representation:

$$z = ASR(x), \quad (1)$$

where we propose to use Multilabel-Binary Vectors (MBV) to represent  $z$  generated in (1), the discrete encoding  $z$  is designed to represents the linguistic content of input speech  $x$ . We define MBV as:

$$z = [e_1, e_2, \dots, e_n] \in \mathbb{R}^n, e_i \in \{0, 1\} \quad (2)$$

where  $z$  is a  $n$  dimension binary vector consists of arbitrary number of zeros and ones. The proposed MBV encoding method is differentiable, and allows end-to-end training with direct gradient backpropagation in an autoencoder framework. To obtain the binarized differentiable vector  $z$ , we linearly project a continuous output vector into a  $\mathbb{R}^{n \times 2}$  space, where  $n$  is the

dimension of the wanted MBV. The dual-channel projection allows each dimension in a MBV to symbolize an arbitrary attribute’s presence. We then perform categorical reparameterization trick with Gumbel-Softmax [22] on the 2<sup>nd</sup>-dimension dual-channel, which is equivalent to asking the model to predict whether an attribute is observed in a given input. Since the two channels are linked together by Gumbel-Softmax, simply picking one of them is sufficient, hence we select the first channel as our MBV encoding  $z$ . We illustrate the process of MBV encoding in Figure 1 (b). Former approaches also use Gumbel-Softmax on the output layer to obtain one-hot vectors [22]. However, the resulting one-hot encoding space is too sparse for machines to learn linguistic meanings, which we verified in our experiments.

#### 2.1.2. Voice reconstruction and conversion

Given ASR-Encoder’s output, the TTS-Decoder is trained to generate an output acoustic feature sequence defined as:

$$X(x, y) = TTS(ASR(x), y), \quad (3)$$

where  $X(x, y)$  is a reconstruction of  $x$  from  $z = ASR(x)$  given the speaker identity  $y \in Y$ . Mean Absolute Error (MAE) is employed to evaluate the ASR-TTS autoencoder reconstruction loss, as MAE is reported to be able to generate sharper outputs than Mean Square Error [23]. The reconstruction loss is given by:

$$L_{rec}(\theta_{asr}, \theta_{tts}) = E_{(x,y) \sim (X,Y)} \|x - X(x, y)\|_1, \quad (4)$$

where  $\theta_{asr}$  and  $\theta_{tts}$  are the parameters of the ASR-Encoder and TTS-Decoder, respectively. We uniformly sample  $(x, y)$  for training in (4). Because the speaker identity is provided to the TTS-Decoder, the proposed MBV encodings  $z$  is able to learn an abstract space that is invariant to speaker identity and only encodes the content of speech, without using any form of linguistic supervision.

At inference time, given a source speech  $x_s$ , and target speaker  $y_t$ , the TTS-Decoder can generate the voice of the target speaker  $y_t$  using the linguistic content  $ASR(x_s)$  from  $x_s$ :

$$X(x_s, y_t) = TTS(ASR(x_s), y_t), \quad (5)$$

$X(x_s, y_t)$  is the output of TTS-Decoder, which has the linguistic content of  $x_s$  but style of speaker  $y_t$ .

### 2.2. Target guided adversarial learning

With the learned speaker invariant ASR-Encoder mapping in (1), we successfully represent speech content with discrete binary vectors MBVs. In this section, we describe how adversarial training is used to boost VC quality based on the discrete linguistic units learned in Section 2.1. We train a TTS-Patcher in an unsupervised manner, in which the TTS-Patcher generates a spectrum mask that residually augments the output of equation (5), resulting in a more precise voice conversion result.

We define two sets of speakers, source speaker set and target speaker set, where we aim to convert the speech of a source speaker into a target speaker’s style while preserving its content. Let  $x_s \in X_S$  and  $x_t \in X_T$  be sequences where  $X_S$  and  $X_T$  are the set of sequences from source speakers and target speakers, respectively. Let  $y_s \in Y_S$  be a speaker from the set of all source speakers  $Y_S$  who produce  $X_S$ , and let  $y_t \in Y_T$  be a speaker from the set of all target speakers  $Y_T$  who produce  $X_T$ .

### 2.2.1. Adversarial learning step

In Figure 2, a TTS-Patcher is trained as a generator under an adversarial learning setting. The TTS-Patcher takes  $ASR(x_s)$  and  $y_t$  as input, and generates a spectrum mask  $P(x_s, y_t)$  that ranges from zero to one, and modifies  $X(x_s, y_t)$  through a residual augmentation:

$$X'(x_s, y_t) = (P(x_s, y_t) \otimes X(x_s, y_t)) \oplus X(x_s, y_t). \quad (6)$$

The  $\oplus$  and  $\otimes$  symbols indicate element-wise addition and multiplication, respectively. In (6),  $x_s$  is a randomly sampled input speech segment from source speaker  $y_s$ , where  $y_t$  is a randomly sampled target speaker.  $X(x_s, y_t)$  is the voice conversion utterance obtained in (5),  $P(x_s, y_t)$  is the output of TTS-Patcher, and finally  $X'(x_s, y_t)$  the augmented spectrum.

A discriminator  $D$  is trained to distinguish whether an input acoustic feature sequence is real or reconstructed by machine. Since naive GAN [24] is notoriously hard to train, we minimize the Wasserstein-1 distance between real and fake distributions instead, as proposed in the WGAN [18] formulation:

$$L_{WGAN} = \mathbb{E}_{x_t \sim X_T} [D(x_t)] - \mathbb{E}_{x_s \sim X_S, y_t \sim Y_T} [D(X'(x_s, y_t))]. \quad (7)$$

The discriminator computes the Wasserstein-1 distance of two distributions: real data  $x_t$  sampled from the target speaker set  $X_T$ , and augmented voice conversion outputs from (6). We use the alternative WGAN-GP [25] to enforce the 1-Lipschitz constraint required by  $D$ , where weight clipping is replaced with gradient penalty. On the last layer of the discriminator, we stretch an additional layer that learns a classifier to predict speaker from a given speech. This allows the discriminator to consider input spectrum’s fidelity and speaker identity at the same time [26].

### 2.2.2. Target guided training step

The decoupled learning [27] of ASR-TTS autoencoder and TTS-Patcher stabilizes the GAN [24] training process. However we found that under the adversarial learning scheme, the TTS-Patcher can easily learn to deceive the discriminator by over-adding style, this greatly compromises the original speech content. This is caused by the discriminator’s inability to discriminate utterances with incorrect or ambiguous content, the discriminator only learns to focus on speaker style. As a solution, we propose to add an additional target guided training step, we apply additional reconstruction loss after every adversarial step, as shown in Figure 2. Instead of converting  $x_s$ , the ASR-Encoder now takes a segment of target speech  $x_t$  as input, equation (6) then becomes:

$$X'(x_t, y_t) = (P(x_t, y_t) \otimes X(x_t, y_t)) \oplus X(x_t, y_t), \quad (8)$$

and we minimize MAE between  $x_t$  and  $X'(x_t, y_t)$ :

$$L_{rec}(\theta_p) = E_{(x_t, y_t) \sim (X_T, Y_T)} \left\| x_t - X'(x_t, y_t) \right\|_1, \quad (9)$$

where  $\theta_p$  is the parameter of the TTS-Patcher. This loss effectively guides the TTS-Patcher’s update under adversarial settings, as the added style is regularized to preserve intelligibility.

## 3. Implementation

The ASR-Encoder is inspired by the CBHG module [1], where the linear output of ASR-Encoder is fed to the MBV encoding

module. We add noise in training by adding dropout layers in the ASR-Encoder as suggested in [28]. The TTS-Decoder and TTS-Patcher have identical model architectures, where we use pixel shuffle layers to generate high resolution spectrum [29]. We add speaker embedding on the feature map of all layers, where a distinct embedding is learned for all different layers as different information may be needed for each layer. The discriminator is consist of 2D-convolution blocks for temporal texture capturing, and convolutions projection layers followed by fully-connected output layers. We trained the network using Adam [30] optimizer and a batch size of 16. In the discrete linguistic units discovery stage, we train the ASR-TTS autoencoder for 200k mini-batches. In the target guided adversarial learning stage we train the model for 50k mini-batches, in one batch we train a step of adversarial learning including 5 iterations of discriminator update and 1 iteration of generator update, followed by a target guided reconstruction step.

We train and evaluate our model using the ZeroSpeech 2019 English dataset [19]. In particular, we use the “Train Unit Dataset” as our source speaker set  $X_S$ , the “Train Voice Dataset” as our target speaker set  $X_T$ , and we evaluate models with the “Test Dataset”. We used log-magnitude spectrograms as acoustic features, the detailed settings are in Table 1. During training, our model is trained to process 128 consecutive overlapping frames of spectrogram, where we uniformly sample from the dataset. At inference time, for a given input with more than 128 frames, the model process them as segments and concatenate the outputs on the time-axis. Source code are publicly available<sup>1</sup>.

Table 1: Acoustic feature preprocess settings

Pre-emphasis	0.97	Sample rate	16k
Frame length	50 ms	Mel-spec	80
Frame shift	12.5 ms	Linear-spec	1024
Window type	Hann	Vocoder	Griffin-Lim

## 4. Experiments

We compare the proposed MBV encodings with one-hot encodings, continuous encodings, and continuous encodings with additional loss [11], all of which under the same autoencoder training setting as described in Section 3.

### 4.1. Degree of disentanglement

To evaluate the degree of disentanglement for the proposed MBV with respect to speaker characteristics, we trained a speaker verification classifier that takes spectrum as input and predicts speaker identity. We encode speech from source speaker and convert it to a target speaker. With the pre-trained classifier we measure target speaker classification accuracy on the converted results. A disentangled representation should produce voice similar to the target speakers and leads to higher classification accuracy. The results are shown in Table 2, where the *Dim* column indicates the dimension  $n$  of encoding vectors  $z \in \mathbb{R}^n$ . One-hot encodings are insufficient to encode speech, resulting in poor conversion results. continuous encodings are incapable of disentangling content from style, resulting in lower classification accuracy. Where as the proposed MBV encodings has the ability to preserve speech content while removing

<sup>1</sup><https://github.com/andi611/ZeroSpeech-TTS-without-T>

Table 2: Comparison of different latent representations.

Types of encodings	Dim	Acc
One-hot	1024	43.3%
continuous	1024	84.1%
continuous	128	79.9%
continuous (with add'l loss)	1024	78%
continuous (with add'l loss)	128	81.3%
Ours (MBV)	1024	<b>92.3%</b>
Ours (MBV)	128	<b>93.9%</b>

speaker style. Although both one-hot vectors and MBV are discrete, each dimension of one-hot vector corresponds to a linguistic unit, while each dimension of MBV may correspond to a pronunciation attribute. This makes MBV more data efficient than one-hot vectors.

#### 4.2. Subjective and objective evaluation

We perform subjective human evaluation on the converted voices. We use 20 subjects to grade each method on a 1 to 5 scale under two measures: the naturalness of speech and the similarity in speaker characteristics to the target speaker. In Table 3 we show the result of our evaluation, the proposed method results in significant increase of target similarity with a slight degrade of naturalness. We easily achieve comparable speech intelligibility as ordinary continuous methods, while achieving better voice conversion quality with more disentanglement (Table 2). Subjective and objective evaluations suggest that the proposed MBV method eliminates speaker identity while reserving content within speech, and is suitable for voice conversion.

Table 3: Results of subjective human evaluation. All methods used an encoding dimension of 1024 if not specified otherwise.

Types of encodings	naturalness	similarity
continuous	3.80	2.14
continuous (with add'l loss)	3.21	2.58
Ours (MBV with dim 6)	1.61	1.51
Ours (MBV)	3.36	<b>3.06</b>
Ours (with adv. training)	2.57	<b>3.15</b>

#### 4.3. Encoding dimension analysis

We use several objective measures to determine the quality of an encoding, these measurements are shown in Table 4. The *CER* column is the output Character Error Rate (CER) from a pre-trained ASR, where we use the ASR results of real input voice as ground truth, which measures intelligibility of the converted speech. The *BR* column measures the bitrate (amount of information) that encodings carry in average with respect to the testing set, as suggested in [19]. The *ABX* column measures the machine ABX score, which indicates the goodness of encoding quality [31][19]. The *distinct* column indicates the number of unique symbols used to encode speech in the test set. Lower values suggest a better performance for all the measures described above. In Table 4, we compare the proposed method with other approaches along side with the baseline model [32][33] demonstrated in [19]. When compared to other approaches, the proposed method achieves lower *BR* and *distinct* values with comparable *ABX* scores. Due to the dis-

Table 4: Performance of different encoding dimensions.

Method	Dim	CER	BR	ABX	distinct
Baseline	200	1.000	71.98	35.90	65
Cont.	1024	0.036	138.45	31.83	16849
	128	0.040	138.45	33.96	16849
Ours	1024	0.196	138.45	32.02	16849
	512	0.313	138.45	32.82	16849
	256	0.430	138.45	32.52	16849
	128	0.629	138.45	31.58	16849
	64	0.717	138.35	32.57	16772
	32	0.797	134.80	31.82	14591
	16	0.887	105.96	35.62	3723
	8	0.998	61.79	38.10	146
	7	0.998	55.97	37.71	94
	6	1.000	<b>48.78</b>	39.60	<b>51</b>
5	1.000	<b>41.32</b>	41.79	<b>28</b>	

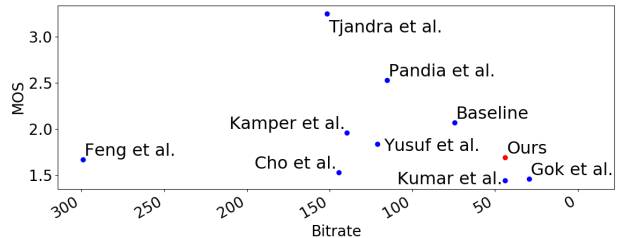


Figure 3: Partial results of the ZeroSpeech Challenge.

crete and differentiable nature of MBV, the proposed method can be used in other unsupervised end-to-end clustering or classification tasks, where other approaches may fail to generalize.

#### 4.4. The zero resource speech challenge competition

We compete with other teams in the ZeroSpeech 2019 Challenge [19] at a global scale. We use the proposed method with a dimension of 6 to achieve extremely low bitrate, and were able to encode a whole language with less than 64 distinct units, human evaluation also suggests that the produced speech are still acceptable (Table 3). On the Surprise dataset [20][21] leaderboard, the proposed method is 2<sup>nd</sup> place in terms of low bitrate, while achieving higher Mean Opinion Score (MOS) and lower CER than the 1<sup>st</sup> place team, as shown in Figure 3. Although the proposed approach does not achieve high MOS because it is an inevitable trade-off with extremely low bitrate, in Table 3 we have shown that with a larger encoding dimension we were able to generate outstanding voice converted speech.

## 5. Conclusions

We proposed to use multilabel-binary vectors to represent the content of human speech, as its discrete nature offers a strong extraction of speaker-independent representation. We show that these discrete units naturally possess the ability of disentangling speech content and style, which makes them extremely suitable for voice conversion tasks. Also, we show that these discrete units indeed produce better style disentanglement than ordinary settings, and finally we were able to improve voice conversion results through the addition of residual augmented signals.

## 6. References

- [1] Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio *et al.*, “Tacotron: Towards end-to-end speech synthesis,” *arXiv preprint arXiv:1703.10135*, 2017.
- [2] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Skerry-Ryan *et al.*, “Natural tts synthesis by conditioning wavenet on mel spectrogram predictions,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 4779–4783.
- [3] W. Ping, K. Peng, A. Gibiansky, S. O. Arik, A. Kannan, S. Narang, J. Raiman, and J. Miller, “Deep voice 3: Scaling text-to-speech with convolutional sequence learning,” *arXiv preprint arXiv:1710.07654*, 2017.
- [4] S.-F. Huang, Y.-C. Chen, H. yi Lee, and L.-S. Lee, “Improved audio embeddings by adjacency-based clustering with applications in spoken term detection,” *CoRR*, vol. abs/1811.02775, 2018.
- [5] W. He, W. Wang, and K. Livescu, “Multi-view recurrent neural acoustic word embeddings,” *arXiv preprint arXiv:1611.04496*, 2016.
- [6] S. Settle and K. Livescu, “Discriminative acoustic word embeddings: Recurrent neural network-based approaches,” *2016 IEEE Spoken Language Technology Workshop (SLT)*, pp. 503–510, 2016.
- [7] Y.-A. Chung, C.-C. Wu, C.-H. Shen, H. yi Lee, and L.-S. Lee, “Audio word2vec: Unsupervised learning of audio segment representations using sequence-to-sequence autoencoder,” in *INTER-SPEECH*, 2016.
- [8] K. Levin, K. Henry, A. Jansen, and K. Livescu, “Fixed-dimensional acoustic embeddings of variable-length segments in low-resource settings,” in *2013 IEEE Workshop on Automatic Speech Recognition and Understanding*. IEEE, 2013, pp. 410–415.
- [9] W.-N. Hsu, S. Zhang, and J. R. Glass, “Unsupervised learning of disentangled and interpretable representations from sequential data,” in *NIPS*, 2017.
- [10] A. Jansen, M. Plakal, R. Pandya, D. P. W. Ellis, S. Hershey, J. Liu, R. C. Moore, and R. A. Saurous, “Unsupervised learning of semantic audio representations,” *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 126–130, 2018.
- [11] J.-c. Chou, C.-c. Yeh, H.-y. Lee, and L.-s. Lee, “Multi-target voice conversion without parallel data by adversarially learning disentangled audio representations,” *arXiv preprint arXiv:1804.02812*, 2018.
- [12] 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 *2016 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA)*. IEEE, 2016, pp. 1–6.
- [13] Y. Gao, R. Singh, and B. Raj, “Voice impersonation using generative adversarial networks,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 2506–2510.
- [14] L. Badino, C. Canevari, L. Fadiga, and G. Metta, “An auto-encoder based approach to unsupervised learning of subword units,” in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2014, pp. 7634–7638.
- [15] C.-T. Chung, C.-Y. Tsai, H.-H. Lu, C.-H. Liu, H. yi Lee, and L.-S. Lee, “An iterative deep learning framework for unsupervised discovery of speech features and linguistic units with applications on spoken term detection,” *2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*, pp. 245–251, 2015.
- [16] J. Chorowski, R. J. Weiss, S. Bengio, and A. v. d. Oord, “Un-supervised speech representation learning using wavenet autoencoders,” *arXiv preprint arXiv:1901.08810*, 2019.
- [17] A. van den Oord, O. Vinyals *et al.*, “Neural discrete representation learning,” in *Advances in Neural Information Processing Systems*, 2017, pp. 6306–6315.
- [18] M. Arjovsky, S. Chintala, and L. Bottou, “Wasserstein gan,” *arXiv preprint arXiv:1701.07875*, 2017.
- [19] E. Dunbar, R. Algayres, J. Karadayi, M. Bernard, J. Benjumea, X.-N. Cao, L. Miskic, C. Dugrain, L. Ondel, A. W. Black, L. Besacier, S. Sakti, and E. Dupoux, “The Zero Resource Speech Challenge 2019: TTS without T,” in *20th Annual Conference of the International Speech Communication Association (INTERSPEECH 2019): Crossroads of Speech and Language*, 2019, submitted. [Online]. Available: <https://zerospeech.com/2019/>
- [20] S. Sakti, R. Maia, S. Sakai, T. Shimizu, and S. Nakamura, “Development of hmm-based Indonesian speech synthesis,” *2008 Proc. Oriental COCOSA*, pp. 215–220, November 2008.
- [21] S. Sakti, E. Kelana, H. Riza, S. Sakai, K. Markov, and S. Nakamura, “Development of Indonesian large vocabulary continuous speech recognition system within a-star project,” *2008 Proc. Technologies and Corpora for Asia-Pacific Speech Translation (TCAST)*, pp. 19–24, January 2008.
- [22] E. Jang, S. Gu, and B. Poole, “Categorical reparameterization with gumbel-softmax,” *arXiv preprint arXiv:1611.01144*, 2016.
- [23] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, “Image-to-image translation with conditional adversarial networks,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [24] 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 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2014, pp. 2672–2680.
- [25] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, “Improved training of Wasserstein GANs,” in *Advances in Neural Information Processing Systems*, 2017, pp. 5767–5777.
- [26] A. Odena, C. Olah, and J. Shlens, “Conditional image synthesis with auxiliary classifier GANs,” in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR.org, 2017, pp. 2642–2651.
- [27] Z. Zhang, Y. Song, and H. Qi, “Decoupled learning for conditional adversarial networks,” in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2018, pp. 700–708.
- [28] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, “Image-to-image translation with conditional adversarial networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.
- [29] W. Shi, J. Caballero, F. Huszar, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, “Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1874–1883.
- [30] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [31] T. Schatz, V. Peddinti, F. Bach, A. Jansen, H. Hermansky, and E. Dupoux, “Evaluating speech features with the Minimal-Pair ABX task: Analysis of the classical MFC/PLP pipeline,” in *INTERSPEECH 2013 : 14th Annual Conference of the International Speech Communication Association*, Lyon, France, Aug. 2013, pp. 1–5.
- [32] L. Ondel, L. Burget, and J. Černocký, “Variational inference for acoustic unit discovery,” *Procedia Computer Science*, vol. 81, pp. 80–86, 2016.

- [33] Z. Wu, O. Watts, and S. King, "Merlin: An open source neural network speech synthesis system," in *9th ISCA Speech Synthesis Workshop (2016)*, Sep. 2016, pp. 218–223.