

# SELF-SUPERVISED REPRESENTATIONS FOR SINGING VOICE CONVERSION

Tejas Jayashankar<sup>1,\*</sup>, Jilong Wu<sup>2</sup>, Leda Sari<sup>2</sup>, David Kant<sup>2</sup>,  
Vimal Manohar<sup>2</sup>, Qing He<sup>2</sup>

<sup>1</sup>Massachusetts Institute of Technology <sup>2</sup>Meta AI

## ABSTRACT

A singing voice conversion model converts a song in the voice of an arbitrary source singer to the voice of a target singer. Recently, methods that leverage self-supervised audio representations such as HuBERT and Wav2Vec 2.0 have helped further the state-of-the-art. Though these methods produce more natural and melodic singing outputs, they often rely on confusion and disentanglement losses to render the self-supervised representations speaker and pitch-invariant. In this paper, we circumvent disentanglement training and propose a new model that leverages ASR fine-tuned self-supervised representations as inputs to a HiFi-GAN neural vocoder for singing voice conversion. We experiment with different  $f_0$  encoding schemes and show that an  $f_0$  harmonic generation module that uses a parallel bank of transposed convolutions (PBTC) alongside ASR fine-tuned Wav2Vec 2.0 features results in the best singing voice conversion quality. Additionally, the model is capable of making a spoken voice sing. We also show that a simple  $f_0$  shifting scheme during inference helps retain singer identity and bolsters the performance of our singing voice conversion model. Our results are backed up by extensive MOS studies that compare different ablations and baselines.

**Index Terms**— Singing voice conversion, self-supervised representations, neural vocoder,  $f_0$  encoder

## 1. INTRODUCTION

The goal of a singing voice conversion task is to convert a song in the voice of a source singer, say A, to the voice of a target singer, say B. The conversion model should retain the singer-invariant content of the song such as the linguistic content and only change the singer dependent characteristics such as the  $f_0$  and the singer's timbre.

Early work on singing voice conversion used HMM based architectures [1, 2, 3] for modeling the latent linguistic and melodic content in singing voices. Subsequent work based on parametric statistical models tackled the problem of parallel singing voice conversion, wherein paired singing samples in the voice of both the source singer and target singer are available [4, 5]. Learning parallel singing voice conversion models is difficult as it is expensive to collect parallel samples from multiple different singers.

To circumvent the data collection difficulties associated with parallel singing voice conversion models, recent work has focused on non-parallel singing voice conversion. In addition to this, the success of neural network architectures, specifically neural vocoders such as WaveNet, [6], WaveRNN [7] and HiFi-GAN [8], has once again made the research area of voice conversion a hot topic. The method of Unsupervised Singing Voice Conversion [9] employs an

autoencoder architecture with a WaveNet decoder to convert between a fixed set of singers. The model is trained with a domain confusion loss to extract singer-invariant features from the encoder. The decoder is conditioned on the target singer identity to synthesize a song in the desired voice. PitchNet [10] builds on this model by employing an additional adversarial pitch confusion term to extract pitch-invariant and singer-invariant features from the encoder. The WaveNet decoder is conditioned on additional source  $f_0$  information which can help improve the quality of the converted singing voice and retain important melodic information.

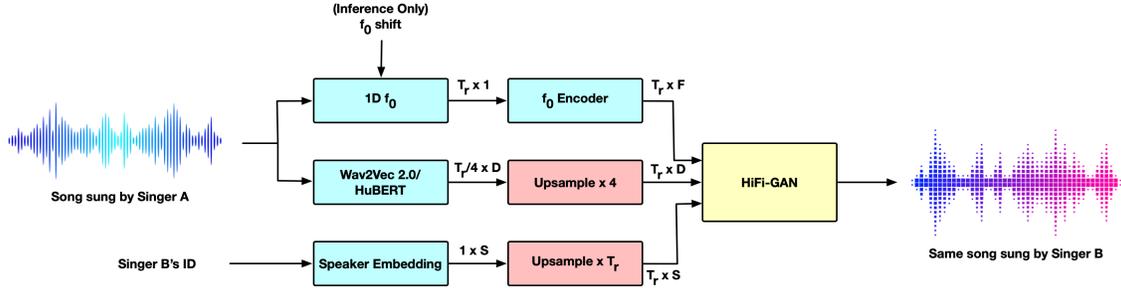
Rather than relying on a confusion loss, the method of Unsupervised Cross-Domain Singing Voice Conversion [11] trains a GAN based architecture with a WaveNet generator and a convolutional discriminator. The generator synthesizes singing audio given  $f_0$  information, a speaker embedding and intermediate ASR features, which are assumed to be disentangled. During inference the source speaker embedding can be swapped out with a target speaker embedding to perform voice conversion.

The use of phonetic posteriors (PPGs), namely the penultimate layer output of an ASR model, has gained large traction in the voice conversion community. The model in [12] uses PPG outputs to model the linguistic content of the source singing sample and the state-of-the-art HiFi-GAN vocoder for synthesis. They build upon [11] by using HuBERT [13] features to model the melodic content. These self-supervised features are learned embeddings that encapsulate linguistic and melodic information within windows of the input audio. They are used as additional melodic conditioning to the model on top of the source  $f_0$  for improving the singing voice quality. The HuBERT features are fine-tuned to be pitch-invariant and speaker-invariant to avoid mismatch between the input  $f_0$  and the pitch encoded within HuBERT features.

In this paper we further investigate the benefits of pre-trained self-supervised features for singing voice conversion/synthesis. Our main contribution is a new singing voice conversion model that uses a parallel bank of transposed convolutions (PBTC) for encoding  $f_0$  and ASR fine-tuned self-supervised features for modeling the singer invariant features. Furthermore, to retain target singer similarity, we revisit a simple  $f_0$  shifting technique from voice conversion literature and demonstrate how it can be used with our model to boost target singer identity similarity after conversion. We perform extensive MOS studies to evaluate different architectural designs for singing voice conversion.

Our paper is organized into five sections. In Section 2, we explain our proposed model architecture and introduce concepts related to  $f_0$  encoding, self-supervised representations and neural vocoders. In Section 3, we provide implementation details and report results from the MOS studies we performed to evaluate our models. Finally, we wrap up with a conclusion, some remarks on future work and acknowledgements.

\* Work performed during internship at Meta AI



**Fig. 1:** The overall architecture of our proposed singing voice conversion system. The HiFi-GAN decodes a new singing sample in the voice of Singer B given a singer embedding corresponding to this singer and content (pitch + self-supervised) features extracted from the original song in Singer A’s voice. Note that the features are extracted at different frequencies and need to be upsampled appropriately before concatenation. During inference,  $f_0$  scale and shift is performed to match the source  $f_0$  with the target singer’s  $f_0$  distribution.

## 2. METHOD

The overall architecture of our proposed model is shown in Figure 1. The model is composed of an  $f_0$  feature module, a singer embedding module, a self-supervised feature extractor and a HiFi-GAN decoder. In this section we provide more details about each of these components and the overall working of the system.

### 2.1. System Overview

Our model operates by extracting two different types of features from the source audio: 1) the source  $f_0$  and 2) self-supervised feature embeddings of dimension  $D$ . The  $f_0$  is further processed by embedding it into a continuous latent space of dimension  $F$  using an  $f_0$  feature encoder. Additionally, a target singer ID is provided to the model to indicate in whose voice the song must be sung. The singer ID is fed to an embedding model, which in practice is a learned look up table (LUT) of size  $N \times S$  where the dimension of each embedding is  $S$  and  $N$  is the number of singers. The three different features are then upsampled and concatenated together resulting in a feature with dimensions  $T_r \times (F + D + S)$  where  $T_r$  is the total number of frames in the input source audio. The HiFi-GAN decoder then synthesizes a singing voice at 24 kHz from the concatenated features.

During training, the singer ID is unchanged and the task is to decode the source audio from the concatenated features. During inference, we perform voice conversion by swapping the source singer ID with a different target singer ID. We additionally shift the source  $f_0$  to the target  $f_0$ ’s domain by performing a simple scale and shift operation (see Section 2.5).

### 2.2. Self-supervised Features: Wav2Vec 2.0 and HuBERT

Our model uses self-supervised features to encode the speaker-invariant characteristics of a singing voice such as the linguistic/speech content and pitch independent melodic content (e.g., timing information such as the duration of phonemes).

Taking inspiration from recent voice conversion models [14], we experiment with the recent and popular Wav2Vec 2.0 [15] and HuBERT [13] embeddings. Both Wav2Vec 2.0 and HuBERT consist of a CNN feature encoder that encodes the input audio frames and a BERT-like [16] transformer trained with a masked predictive task that enforces that the learned encodings contain strong acoustic and linguistic content given the past and future context.

Wav2Vec 2.0 is trained with a contrastive loss to correctly predict the quantized CNN encoder output at time step  $t$  from a set of

distractors. On the other hand, HuBERT training forces the model to correctly predict the quantized latent from a random layer  $k$  of the transformer at time step  $t$  amongst a set of distractors taken at different time steps from the same layer — hence the name **Hidden unit BERT**. In both settings, the continuous input at time step  $t$  to the transformer is masked out.

We also experiment with ASR fine-tuned Wav2Vec 2.0 features. The pre-trained Wav2Vec 2.0 features are fine-tuned on an ASR prediction task by fusing a softmax layer on top of the BERT transformer with a CTC [17] loss. Some recent voice conversion models have relied on ASR outputs such as phonetic posteriors (PPGs) for modeling the linguistic content of singing voices. We postulate that the ASR fine-tuning task helps preserve linguistic content in the embeddings akin to PPGs. These embeddings also contains important singer-invariant melody information that would otherwise need to be extracted using a special encoder or via disentanglement training/confusion losses as in [9, 10, 12].

### 2.3. $f_0$ Feature Encoder

We initially experimented with 1D  $f_0$  features, i.e.,  $F = 1$ . However, the HiFi-GAN synthesized singing outputs sounded monotonic and off-pitch. In most cases the desired prosody was not synthesized. Subsequently, we experimented with two  $f_0$  feature encoders to learn higher dimensional  $f_0$  representations as described next.

#### 2.3.1. $Q$ -LUT Embedding Module

The continuous 1D  $f_0$  sequence is first mapped to the range  $[0, 1]$  using mean and variance normalization. The entries are then uniformly quantized to  $L$  bins. A look up table (LUT) with  $L$  entries of dimension  $F$  is learned end-to-end with the HiFi-GAN decoder to learn a rich pitch embedding for each possible quantized value of the 1D  $f_0$ . An illustration of the architecture is shown in Figure 2a.

#### 2.3.2. PBTC Harmonic Generation Module

The parallel bank of transposed convolutions (PBTC) was introduced in [18] for generating harmonic information from 1D  $f_0$  sequences. The continuous 1D  $f_0$  sequence is again globally normalized and quantized to  $L$  bins. The quantized sequence of length  $T_r$  is then one-hot projected to an embedding of shape  $T_r \times L$ . This embedding is parallelly processed with  $K$  different 1D transposed convolutions with linearly increasing dilation rates and  $F$  filters. Since the input embedding is sparse (due to the one-hot encoding), the transposed convolutions learn appropriate filters to apply at the

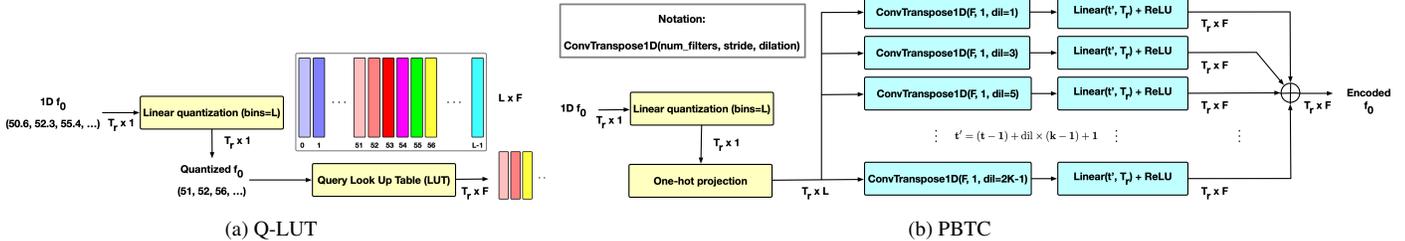


Fig. 2: Architectures of the Q-LUT and PBTC  $f_0$  encoding models that we experiment with in our system.

non-zero positions to generate rich harmonic information at different rates. The output of each transposed convolution is linearly projected to the same duration  $T_r$  and then summed together resulting in a new feature of shape  $T_r \times F$ . An illustration of the architecture is shown in Figure 2b. As will be described in Section 3.4, MOS studies indicate that the PBTC architecture results in higher quality singing audio in comparison to the Q-LUT method.

#### 2.4. Synthesis with the HiFi-GAN Vocoder

We use a slightly modified version of the original HiFi-GAN [8] for synthesizing the singing voice from the input features. All features are upsampled to the same rate before being fed to the generator.

In our experiments the inputs to the generator are upsampled to a rate of 200 Hz. The generator  $G$  consists of four upsampling blocks with upsample rates of [3, 4, 5, 8]. Each upsampling block consists of three residual blocks with kernel sizes of [3, 7, 11] and dilation rates of [1, 3, 5] per residual block. The discriminator comprises two networks: a Multi-Period Discriminator (MPD) and a Multi-Scale Discriminator (MSD). The MPD consists of five sub-discriminators each operating on equally spaced samples. The five sub-discriminators use different spacings of [2, 3, 5, 7, 11]. The MSD uses three sub-discriminators operating at three different scales of the input signal: original scale,  $\times 2$  downsampled signal and  $\times 4$  downsampled signal.

#### 2.5. $f_0$ Shifting During Inference

During training of our model, the speaker embedding LUT is learned end-to-end with the  $f_0$  encoder and the HiFi-GAN vocoder, and qualitative experiments suggested that the learned singer embedding internally encodes the moments of the corresponding singer’s  $f_0$  distribution and hence are not disentangled. However, during inference, the source singer embedding can only be swapped out for a target singer embedding under the assumption that the embedding is disentangled from the  $f_0$ . Attempts at disentangling the  $f_0$  and singer embedding features led to poor results. Instead, a simpler method that aligns the  $f_0$  from the source singer with the internal  $f_0$  information encoded in the target singer embedding worked well.

Assume that the  $f_0$  distribution for singer A is a normal distribution  $f_0^A \sim \mathcal{N}(\mu_A, \sigma_A)$  and that the  $f_0$  distribution for singer B is  $f_0^B \sim \mathcal{N}(\mu_B, \sigma_B)$ . Then given a source  $f_0$  sequence, say  $f_A$  from singer A, we shift it to match the  $f_0$  distribution of singer B during inference.

$$f_A^{shift} = \frac{\sigma_B}{\sigma_A} (f_A - \mu_A) + \mu_B. \quad (1)$$

### 3. EXPERIMENTS

This section details the dataset, implementation details and results from different Mean Opinion Scores (MOS) studies that ran.

#### 3.1. Dataset

We trained our model on a dataset consisting of speech and singing samples. We found that training only on singing samples resulted in poor synthesized audio quality. We used an internal speech dataset consisting of  $\sim 200$  hours of English clean speech recorded in a voice production studio by contracted professional voice talents. Our singing voice dataset consists of 10+ hours of samples taken from the NUS-48E dataset [19], the Children’s Song Dataset (CSD) [20] and the AmericanSong corpus licensed from SpeechOcean<sup>1</sup>. We use a 70-30 train-test split for training and evaluating our models. All speech and singing samples are sampled at 24 kHz. Our dataset consists of total of 17 different singers and over 20 different speakers.

#### 3.2. Implementation Details

We use the CREPE [21] pitch tracker to extract 1D  $f_0$ . The extraction range is set to [50, 800] Hz, where the upper bound was chosen to accommodate high pitched singing voices. We use large versions of the pre-trained Wav2Vec 2.0 models found on the original authors’ website<sup>2</sup>. Our best performing model uses the ASR fine-tuned version of the Wav2Vec 2.0 model which was fine-tuned on 960 hours of LibriSpeech [22] data. Similarly for HuBERT we experiment with the extra large version of the pre-trained model<sup>3</sup> that was trained on 60k hours of the Libri-Light dataset [23]. Both Wav2Vec 2.0 and HuBERT extract features at a rate of 50 Hz from 16 kHz audio. Thus, we first downsample the audio and then upsample the extracted features to the same extraction frequency as the  $f_0$  features (200 Hz).

We found that our  $f_0$  encoders work well with  $L = 400$  bins and an embedding size of  $F = 256$ . Our PBTC architecture worked well with  $K = 10$  parallel layers. We extract Wav2Vec 2.0 features from layer 12 of the transformer with dimensionality  $D = 1024$  and we extract HuBERT features from the penultimate layer of the transformer with dimensionality  $D = 1024$  as well. Our speaker embedding module is a simple look up table with dimensionality  $S = 128$ . We train our models for 5 million steps on 8 NVIDIA A100 GPUs. We use a learning rate of 0.0002 with hyperparameters  $\lambda_{recon} = 40, \lambda_{fm} = 1$ .

#### 3.3. Experimental Setup

We evaluated our singing voice conversion model by experimenting with different  $f_0$  encoding schemes and *pre-trained* self-supervised representations. Additionally, we implemented a baseline based on disentanglement training to learn pitch and speaker-invariant self-supervised representations. Since no publicly available baselines were available we implemented everything ourselves. We found that a disentanglement task based on mutual information minimization

<sup>1</sup><https://en.speechocean.com/datacenter/details/3069.html>

<sup>2</sup><https://github.com/facebookresearch/fairseq/tree/main/examples/wav2vec>

<sup>3</sup><https://github.com/facebookresearch/fairseq/tree/main/examples/hubert>

[24] lead to training convergence, but in general we found disentanglement training to be very unstable.

### 3.4. Results

We performed four different MOS studies to evaluate the various systems that we implemented<sup>4</sup>. If a system uses ASR fine-tuned self supervised features, we will mention it explicitly. We synthesized singing samples in 4 speaker voices from our internal speech dataset, with 10 random singing segments from the dev sets. For each MOS study, we asked 300 raters to rate the quality/speaker similarity of 10 randomly chosen samples on a scale of 1-5. All studies were performed using MTurk and unreliable raters were filtered out before computing the scores. We tested the systems along two major axis — 1) *audio quality*, i.e., how natural and human-like the audio clips sounded and 2) *target singer/speaker similarity*, i.e., how similar did a voice converted audio clip sound to a reference sample from the same singer/speaker.

#### 3.4.1. MOS Study 1: Naturalness across $f_0$ encoding schemes and baselines.

We fix the self-supervised representation to HuBERT. We generate converted voice samples using our different  $f_0$  encoding schemes. A baseline based on disentanglement training and one trained only on singing voice samples is also included. We also experiment with  $f_0$  shifting during inference. As shown in Table 1, an  $f_0$  encoding method that uses the PBTC architecture has the highest audio quality. We also see that training only on singing samples and with additional disentanglement losses leads to a drop in performance.

#### 3.4.2. MOS Study 2: Target singer similarity with fixed self-supervised features.

We fix the self-supervised feature to HuBERT and provide the raters with a reference ground truth audio sample of the target singer and a converted sample from our model. As shown in Table 2,  $f_0$  shifting during inference significantly improves the speaker similarity.

#### 3.4.3. MOS Studies 3 & 4: Varying the self-supervised features along with $f_0$ shifting during inference.

We now vary the self-supervised representations being used and perform an  $f_0$  shift to retain as much target speaker identity as possible. As shown in Table 3, using ASR fine-tuned Wav2Vec 2.0 features along with a PBTC encoder results in excellent audio quality. It significantly outperforms baselines that use vanilla Wav2Vec 2.0 features without ASR fine-tuning.

We also test for target singer identity similarity as shown in Table 4. The results show similar trends with the architectures using ASR fine-tuned features retaining the target singer identity well. Surprisingly, the architecture that uses HuBERT + PBTC also achieves higher target singer similarity, probably suggesting that the HuBERT embedding encodes less speaker dependent information in it.

## 4. CONCLUSION

In this paper we propose a singing voice conversion architecture that uses ASR fine-tuned Wav2Vec 2.0 features along with a specialized  $f_0$  encoder. Experiments show that an  $f_0$  encoder based on a parallel bank of transposed convolutions (PBTC) leads to the best audio

<sup>4</sup>Audio samples: <https://tkj516.github.io/Self-Supervised-Singing-Voice-Conversion/>

System	MOS
Ground Truth	4.22 ± 0.04
Q-LUT (Sing only)	3.44 ± 0.07
Q-LUT	3.69 ± 0.04
Q-LUT (Disentanglement)	3.37 ± 0.07
<b>PBTC</b>	<b>3.76 ± 0.06</b>
Q-LUT + $f_0$ shift	3.66 ± 0.06
PBTC + $f_0$ shift	3.68 ± 0.06

**Table 1:** MOS Study 1: Comparison of audio quality with different singing voice conversion models while using HuBERT self-supervised features.

System	MOS
Q-LUT	3.70 ± 0.04
Q-LUT (Disentanglement)	3.57 ± 0.04
<b>Q-LUT + <math>f_0</math> shift</b>	<b>3.79 ± 0.04</b>
<b>PBTC + <math>f_0</math> shift</b>	<b>3.80 ± 0.04</b>

**Table 2:** MOS Study 2: Comparison of target singer/speaker similarity with different singing voice conversion models while using HuBERT self-supervised features.

System	MOS
Wav2Vec2 + $f_0$ shift (Q-LUT)	3.67 ± 0.05
Wav2Vec2 + $f_0$ shift (PBTC)	3.77 ± 0.05
HuBERT + $f_0$ shift (Q-LUT)	3.81 ± 0.04
HuBERT + $f_0$ shift (PBTC)	3.81 ± 0.04
<b>Wav2Vec2-ASR + <math>f_0</math> shift (Q-LUT)</b>	<b>3.83 ± 0.06</b>
<b>Wav2Vec2-ASR + <math>f_0</math> shift (PBTC)</b>	<b>3.91 ± 0.05</b>

**Table 3:** MOS Study 3: Comparison of audio quality while varying the self-supervised feature and  $f_0$  encoder.

System	MOS
Wav2Vec2 + $f_0$ shift (Q-LUT)	3.79 ± 0.04
Wav2Vec2 + $f_0$ shift (PBTC)	3.75 ± 0.05
HuBERT + $f_0$ shift (Q-LUT)	3.82 ± 0.04
<b>HuBERT + <math>f_0</math> shift (PBTC)</b>	<b>3.85 ± 0.06</b>
Wav2Vec2-ASR + $f_0$ shift (Q-LUT)	3.84 ± 0.05
Wav2Vec2-ASR + $f_0$ shift (PBTC)	3.83 ± 0.06

**Table 4:** MOS Study 4: Comparison of target singer/speaker similarity while varying the self-supervised feature and  $f_0$  encoder.

quality and that when combined with  $f_0$  distribution matching during inference helps retain target singer identity. We conduct extensive MOS studies to test our architecture and detail how our model is implemented. Additionally, since our model was trained on both speech and singing data, it can be used to convert a song to the voice of target speaker for which no singing data is available.

For future work, we are looking to use disentangled ASR fine-tuned Wav2Vec 2.0 features for speech to singing voice conversion, with the main challenge being the unsupervised learning of an alignment between speech and melody features. We are also looking into additional architectures for singing voice synthesis with quantized self-supervised representations.

## 5. ACKNOWLEDGEMENTS

We would like to thank Yossi Adi for insightful discussions and feedback.

## 6. REFERENCES

- [1] Keijiro Saino, Heiga Zen, Yoshihiko Nankaku, Akinobu Lee, and Keiichi Tokuda, “An hmm-based singing voice synthesis system,” in *Ninth International Conference on Spoken Language Processing*, 2006.
- [2] Keiichiro Oura, Ayami Mase, Tomohiko Yamada, Satoru Muto, Yoshihiko Nankaku, and Keiichi Tokuda, “Recent development of the hmm-based singing voice synthesis system—sinsy,” in *Seventh ISCA Workshop on Speech Synthesis*, 2010.
- [3] Kazuhiro Nakamura, Keiichiro Oura, Yoshihiko Nankaku, and Keiichi Tokuda, “Hmm-based singing voice synthesis and its application to japanese and english,” in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2014, pp. 265–269.
- [4] Kazuhiro Kobayashi, Tomoki Toda, Graham Neubig, Sakriani Sakti, and Satoshi Nakamura, “Statistical singing voice conversion with direct waveform modification based on the spectrum differential,” in *Fifteenth Annual Conference of the International Speech Communication Association*. Citeseer, 2014.
- [5] Kazuhiro Kobayashi, Tomoki Toda, Graham Neubig, Sakriani Sakti, and Satoshi Nakamura, “Statistical singing voice conversion based on direct waveform modification with global variance,” in *Sixteenth Annual Conference of the International Speech Communication Association*. Citeseer, 2015.
- [6] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu, “Wavenet: A generative model for raw audio,” *arXiv preprint arXiv:1609.03499*, 2016.
- [7] Nal Kalchbrenner, Erich Elsen, Karen Simonyan, Seb Noury, Norman Casagrande, Edward Lockhart, Florian Stimberg, Aaron Oord, Sander Dieleman, and Koray Kavukcuoglu, “Efficient neural audio synthesis,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 2410–2419.
- [8] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae, “Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 17022–17033, 2020.
- [9] Eliya Nachmani and Lior Wolf, “Unsupervised singing voice conversion,” *arXiv preprint arXiv:1904.06590*, 2019.
- [10] Chengqi Deng, Chengzhu Yu, Heng Lu, Chao Weng, and Dong Yu, “Pitchnet: Unsupervised singing voice conversion with pitch adversarial network,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7749–7753.
- [11] Adam Polyak, Lior Wolf, Yossi Adi, and Yaniv Taigman, “Unsupervised cross-domain singing voice conversion,” *arXiv preprint arXiv:2008.02830*, 2020.
- [12] Chao Wang, Zhonghao Li, Benlai Tang, Xiang Yin, Yuan Wan, Yibiao Yu, and Zejun Ma, “Towards high-fidelity singing voice conversion with acoustic reference and contrastive predictive coding,” *arXiv preprint arXiv:2110.04754*, 2021.
- [13] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed, “Hubert: Self-supervised speech representation learning by masked prediction of hidden units,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [14] Adam Polyak, Yossi Adi, Jade Copet, Eugene Kharitonov, Kushal Lakhotia, Wei-Ning Hsu, Abdelrahman Mohamed, and Emmanuel Dupoux, “Speech resynthesis from discrete disentangled self-supervised representations,” *arXiv preprint arXiv:2104.00355*, 2021.
- [15] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 12449–12460, 2020.
- [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [17] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 369–376.
- [18] Jacob Webber, Olivier Perrotin, and Simon King, “Hider-finder-combiner: an adversarial architecture for general speech signal modification,” in *Interspeech 2020-21st Annual Conference of the International Speech Communication Association*, 2020.
- [19] Zhiyan Duan, Haotian Fang, Bo Li, Khe Chai Sim, and Ye Wang, “The nus sung and spoken lyrics corpus: A quantitative comparison of singing and speech,” in *2013 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference*. IEEE, 2013, pp. 1–9.
- [20] Soonbeom Choi, Wonil Kim, Saebul Park, Sangeon Yong, and Juhan Nam, “Children’s song dataset for singing voice research,” in *International Society for Music Information Retrieval Conference (ISMIR)*, 2020.
- [21] Jong Wook Kim, Justin Salamon, Peter Li, and Juan Pablo Bello, “Crepe: A convolutional representation for pitch estimation,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 161–165.
- [22] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, “Librispeech: An asr corpus based on public domain audio books,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 5206–5210.
- [23] J. Kahn, M. Riviere, W. Zheng, E. Kharitonov, Q. Xu, P.E. Mazare, J. Karadayi, V. Liptchinsky, R. Collobert, C. Fuegen, T. Likhomanenko, G. Synnaeve, A. Joulin, A. Mohamed, and E. Dupoux, “Libri-light: A benchmark for ASR with limited or no supervision,” in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. may 2020, IEEE.
- [24] Woo Hyun Kang, Jahangir Alam, and Abderrahim Fathan, “Robust speech representation learning via flow-based embedding regularization,” *arXiv preprint arXiv:2112.03454*, 2021.