# Towards Better Disentanglement in Non-Autoregressive Zero-Shot Expressive Voice Conversion

Seymanur Akti<sup>1</sup>, Tuan Nam Nguyen<sup>1</sup>, Alexander Waibel<sup>1,2</sup>

<sup>1</sup>Interactive Systems Lab, Karlsruhe Institute for Technology, Germany <sup>2</sup>Carnegie Mellon University, USA

seymanur.akti@kit.edu, tuan.nguyen@kit.edu, alexander.waibel@kit.edu

## **Abstract**

Expressive voice conversion aims to transfer both speaker identity and expressive attributes from a target speech to a given source speech. In this work, we improve over a self-supervised, non-autoregressive framework with a conditional variational autoencoder, focusing on reducing source timbre leakage and improving linguistic-acoustic disentanglement for better style transfer. To minimize style leakage, we use multilingual discrete speech units for content representation and reinforce embeddings with augmentation-based similarity loss and mix-style layer normalization. To enhance expressivity transfer, we incorporate local F0 information via cross-attention and extract style embeddings enriched with global pitch and energy features. Experiments show our model outperforms baselines in emotion and speaker similarity, demonstrating superior style adaptation and reduced source style leakage.

Index Terms: speech synthesis, expressive voice conversion

## 1. Introduction

Voice conversion (VC) approaches aim to transform a source audio by transferring speaker characteristics from a target audio while preserving the source content. Conventional VC models perform well in replicating speaker identity but struggle when the target speech is highly expressive. Expressive voice conversion (EVC) expands on this by capturing both speaker and expressive cues, such as emotions, intensity, and pitch, during synthesis. This allows to not only generate read speech but also replicate the emotional nuances of the target speaker. EVC can be applied in dialogue systems such as [1, 2] to improve human-robot interactions or speech translation pipelines [3, 4, 5], ensuring that the emotions and expressions from the source language are retained in the target language speech [6, 7].

EVC was tackled as a supervised sequence-to-sequence task [8], however, due to the challenge of creating parallel emotional speech corpora, many methods explored non-parallel synthesis [9, 10]. A common strategy is disentangling linguistic and acoustic information in a self-supervised manner, then recombining the source's linguistic features with the target's acoustic attributes [11]. However, a persistent challenge is source timbre leakage, where residual speaker timbre from the source speech degrade style transfer. To address this, some studies introduce an information bottleneck in the linguistic encoder to suppress residual acoustic information [12, 13].

In this work, we enhance the information bottleneck to improve disentanglement and reduce source speaker leakage. We adopt a non-parallel, self-supervised speech generation model based on a conditional variational autoencoder, inspired by VITS [14] and FreeVC [15]. Our system uses self-supervised speech representations from mHuBERT-147 [16] as input,

leveraging its discrete speech units to eliminate non-linguistic information more effectively than continuous representations through quantization [17]. Additionally, its multilingual speech units enable cross-lingual EVC, making it particularly valuable for speech translation pipelines. To our knowledge, this is the first use of mHuBERT-147 in a VC setup, combining the benefits of both discrete speech units and multilinguality. To further improve disentanglement, we introduce a perturbation-based similarity loss to minimize variation in the content embedding distribution and integrate mixed-layer normalization from [18] to enhance the linguistic-acoustic disentanglement. We represent both speaker identity and emotional cues in a single global style embedding, using ECAPA-TDNN [19] complemented by global pitch and energy information. Additionally, we add a local F0 encoder for a better prosodic similarity to the target speech. We provide speech samples on the demo page. 1.

# 2. Related Work

EVC is a speech-to-speech task, where early approaches often relied on auto-regressive models trained with parallel speech data [8, 20]. These models typically required text supervision. To reduce reliance on transcripts, [11] proposed a textless approach that extracts discrete speech units directly from audio and learns a translation network between emotion classes. Given that obtaining original parallel data is challenging, some studies explored using synthesized parallel data for accent conversion [21, 22].

Meanwhile, in EVC, many recent methods have adopted self-supervised speech reconstruction techniques to remove the need for aligned data. One of the earliest studies in this area, [9], disentangles emotion-invariant and emotion-variant features and uses an autoencoder to synthesize speech from these features. [23] employs StarGAN for the EVC task while [10] introduces variational autoencoders with non-autoregressive speech synthesis for EVC. [24] applies a VAE for language-agnostic EVC with limited data. [25] leverages self-supervised speech representations and a k-nearest neighbors model for feature retrieval. [26] uses discrete speech units as linguistic inputs and a prosody encoder for acoustic inputs, synthesizing speech with an auto-regressive decoder. The most similar works to ours are conditional VAE based methods adapted for the style conversion task where [12] uses a hierarchical adaptive generator for generating the waveform, [27] adds the style consistency loss for a better style transfer, [28] uses a similar architecture to [12] and adapts it for jointly trained TTS and cross-lingual EVC and [29] introduces prosody extraction and fusing methods for improving VITS for EVC task.

https://seymanurakti.github.io/evc/

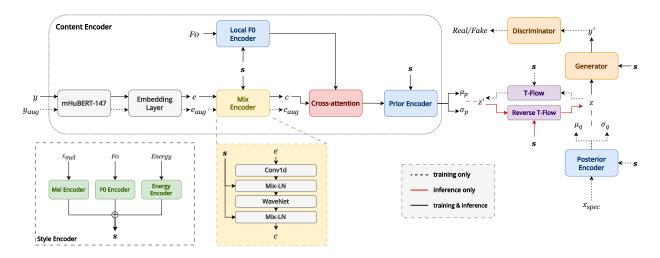


Figure 1: Overall architecture of the proposed system.

# 3. Methodology

We follow the architecture of FreeVC [15] and adapt it for EVC with significant modifications as overall architecture illustrated in Fig. 1. The aim is to learn two distributions: one for linguistic features  $p=N(\mu_p,\sigma_p)$  and one for spectrogram features  $q=N(\mu_q,\sigma_q)$ . During training, the normalizing flow maps the posterior distribution q to the prior distribution p, while the decoder generates speech from samples of q. The content encoder aims to capture only linguistic information, with acoustic features injected via style embeddings through conditional layers. During inference, the reverse flow generates style-injected representations from linguistic features, allowing speech synthesis that preserves the source content while adopting the target style.

#### 3.1. Content Encoder

First, mHuBERT-147 [16] units are extracted from waveform y and quantized. Then, unit embeddings e are obtained through an embedding layer and are processed through the Mix Encoder, which helps to produce style-agnostic content embeddings by applying mixed-layer normalization [18] with scale and bias vectors derived from the mixed style embeddings as:

$$\gamma_{mix}(s) = \lambda \gamma(s) + (1 - \lambda)\gamma(\tilde{s}) \tag{1}$$

$$\beta_{mix}(s) = \lambda \beta(s) + (1 - \lambda)\beta(\widetilde{s}) \tag{2}$$

$$MixLN(e, s) = \gamma_{mix}(s) \times LN(e) + \beta_{mix}(s)$$
 (3)

where  $\widetilde{s}$  is batch-wise shuffled style embeddings and  $\lambda$  is a parameter from Beta distribution. This introduces random style information from other samples in the batch, injecting some noise, reducing the dependency of content embeddings on their original style embeddings and making them more style-agnostic. Our experiments show that Mix-LN also improves content-independence in style embeddings implicitly making them time-invariant. This is particularly useful for mitigating the train-inference mismatch, where training samples share the same content for source and target speech, whereas this alignment does not hold during inference.

Concurrently, the Local F0 Encoder extracts frame-based F0 embeddings from the audio's F0 contours which share the same sampling rate (320) as the content embeddings c. F0 embeddings are then fused with the content embeddings via

multi-head cross-attention, with content embeddings as query and F0 embeddings as key and value. This more effectively captures the pitch flow of the target compared to re-normalizing the source F0 and fusing via summation, as done in [18, 30]. It also allows the target F0 to be directly used as the pitch input during inference—regardless of any length differences. Finally, the prior encoder generates the distribution p(z | c) from the F0-enriched content embeddings.

Additionally, we introduce a perturbation-based similarity loss to improve the quality of content embeddings. We apply augmentation via Parselmouth<sup>2</sup>, modifying the original audio by reducing its pitch range to create less expressive and more uniform speech, and applying pitch shifting to augment speaker identity. The similarity loss then ensures that content embeddings from the original and augmented samples remain close, reducing their dependency on non-linguistic variations. The similarity loss function is formally defined in Eq. 4.

$$L_{sim} = (1 - cos(e, e_{aug})) + (1 - cos(c, c_{aug}))$$
 (4)

#### 3.2. Posterior Encoder

Posterior Encoder takes the linear spectrogram  $x_{spec}$  as the input and generates the posterior distribution  $q(z|x_{spec})$ , sharing the same architecture with the Prior Encoder with convolution layers and projection layer for learning distribution parameters. In order to make the posterior and prior distributions closer, The KL divergence loss is calculated as given in Eq. 5.

$$L_{kl} = KL(q(z|x_{spec})||p(z|c))$$
(5)

# 3.3. Normalizing Flow

The normalizing flow layer learns the mapping from the posterior distribution to the prior distribution during training. We use a Transformer-based normalizing flow, following the approach in [30], which has demonstrated superior performance compared to convolution-only normalizing flows due to its ability to capture longer time dependencies [31]. During inference, the content representation of the source speech is mapped to the posterior distribution using reverse normalizing flow to be decoded by Generator.

<sup>&</sup>lt;sup>2</sup>https://github.com/YannickJadoul/Parselmouth

## 3.4. Waveform Synthesizer

For speech generation, we utilized the HiFi-GAN vocoder [32]. The Generator generates speech waveforms through a series of upsampling layers, while the Discriminator aims to distinguish between real and generated waveforms. The synthesizer employs three loss functions similar to conventional generative adversarial networks, as defined in Eq. 6-8.

$$L_{adv}(G) = \mathbb{E}_z[(D(G(z)) - 1)^2]$$
 (6)

$$L_{adv}(D) = \mathbb{E}_{(y,z)}[(D(y) - 1)^2 + (D(G(z)))^2]$$
 (7)

$$L_{fm}(G) = \mathbb{E}_{(y,z)} \left[ \sum_{l=1}^{T} \frac{1}{N_l} \|D^l(y) - D^l(G(z))\|_1 \right]$$
 (8)

Also, a reconstruction loss is calculated between the original and generated mel spectrograms as shown in Eq. 9.

$$L_{rec} = \|x_{mel} - \hat{x}_{mel}\|_1 \tag{9}$$

#### 3.5. Style Encoder

We use the ECAPA-TDNN model [19] for style encoding. Unlike style encoders relying solely on mel spectrograms, we also incorporate F0 and energy contours. We extract 512-dimensional embeddings from mel spectrograms, F0, and energy contours, then fuse them using a trainable weighted summation, as defined in Eq. 10.

$$s = \lambda_{mel}.e_{mel} + \lambda_{f0}.e_{f0} + \lambda_{energy}.e_{energy}$$
 (10)

Style embeddings are then used to condition all encoders, injecting style information throughout the process. The overall loss function is given in Eq. 11.

$$L = L_{adv}(G) + L_{adv}(D) + L_{fm}(G) + L_{kl} + L_{rec} + L_{sim}$$
(11)

# 4. Experiments and Results

For training, we used a combination of LibriTTS-100 [33], ESD [34] (English only), subset of GigaSpeech [35], and Expresso [36]. All datasets are English and total duration is around 228 hours with more than 920 speakers. We used 2 NVIDIA A600 GPUs for training with batch size of 64 for 1M steps. For the ablation study, we trained the models for 300k steps.

# 4.1. Evaluation

For evaluation, we used test sets of ESD [34], Expresso [36], and LibriTTS [33]. For ESD, the source and target samples are from same speaker with same content and different emotions. In other two datasets, source and target had different speakers and content. For objective evaluation, we use several metrics as follows:

- WER: We use Whisper-Large-3 [37] with text normalization.
- SECS: Speaker embedding cosine similarity computed between synthesized and target audio using Resemblyzer<sup>3</sup>.
- **EECS:** Emotion embedding cosine similarity between synthesized and target audio using Emotion2Vec+ [38].
- ECA: Emotion classification accuracy calculated on the synthesized samples using Emotion2Vec+.
- EER: Equal error rate computed with a speaker verification model [19], with synthesized sample as query, source speech as negative, and target as positive. A lower EER signals better target speaker matching and less source speaker leakage.

We use three subjective metrics: naturalness MOS (nMOS) for speech quality, speaker MOS (sMOS) for speaker similarity, and emotion MOS (eMOS) for emotion accuracy, all rated on a 1–5 scale. We use ESD test set samples for ESD-only models and RAVDESS [39] for zero-shot EVC for human evaluation. 15 users participated in evaluation and users were given 10-15 samples per model across different emotions.

#### 4.2. Style Transfer Results

We compare our approach with three models with VAE architecture and a style encoder for copying expressive style information by integrating emotional datasets in the training. For a fair comparison, we report results in two settings based on the training data: (1) ESD-only setting, comparing against X-E-Speech [28] (trained on ESD) and Consistency-VC [27] (trained on ESD and VCTK); and (2) multi-dataset setting, comparing against Hierspeech++ [30], which is trained on a larger dataset.

Table 1 shows our model surpasses Hierspeech++ in emotion copying across all test sets, achieving higher SECS for both seen (Expresso) and unseen (LibriTTS) speakers, while lower EER on LibriTTS suggests better source speaker identity removal. In the ESD-only comparison, we outperform X-E-Speech and Consistency-VC in converting emotional speech to neutral, proving effective in eliminating source style. It also achieves higher ECA than X-E-Speech in 'neutral to others' and 'overall' but is slightly surpassed by Consistency-VC. In the zero-shot setting with unseen speakers from Expresso, our model shows slightly better emotion and speaker transfer. However, its WER is higher in both setups, likely due to discrete speech units reducing speaker information but also eliminating some linguistic cues, leading to pronunciation artifacts [40, 17].

The subjective metric results at Table 3, support the objective evaluation, as we surpass others in emotion transfer capability for both setups. For speaker transfer, we perform better for zero-shot, but others perform better for ESD version. Given the nMOS scores of our model exceeding 3, it is possible to claim that the audio remains reasonably natural and intelligible, despite the higher WER compared to other methods.

## 4.3. Ablation Study

For assessing the effect of each contribution on EVC, we conducted an ablation study, as shown in Table 2. Results show that replacing F0 injection via summation with cross-attention significantly enhances style transfer. The addition of Mix-LN improves both emotion and speaker similarity, and as reflected in higher ECA (O  $\rightarrow$  N) and lower EER scores, Mix-LN effectively reduces source style leakage. Furthermore, the lower difference between WERs of ESD (source and target has the same content) and LibriTTS (source and target has different content) suggests that Mix-LN mitigates the train-inference mismatch caused by content differences between the source and target speech, possibly by reducing content leakage in the style embeddings. To investigate this, we measured the intra-speaker cosine similarity of style embeddings from the style encoders of both the proposed model and the model without Mix-LN. Ideally, these embeddings should be content-agnostic, meaning similarity should remain high for the same speaker regardless of the content. For unseen speakers in LibriTTS, the average cosine similarity increased from 0.719 to 0.748 with Mix-LN, demonstrating improved disentanglement in style embeddings.

Addition of  $L_{sim}$  improves unseen speaker conversion scores on LibriTTS, indicating reduced speaker identity leakage. The higher overall ECA score further suggests that

<sup>3</sup>https://github.com/resemble-ai/Resemblyzer

Table 1: Comparative results.  $N \rightarrow O$  indicates conversion from neutral to the other emotions. (overall) is for cross-emotion conversion.

Models		ESD			Expresso		LibriTTS		
	WER ↓	ECA (N $\rightarrow$ O) $\uparrow$	ECA (O $\rightarrow$ N) $\uparrow$	ECA (overall) ↑	SECS ↑	EECS ↑	WER↓	SECS ↑	EER ↓
Consistency-VC	4.71%	78.3%	72.6%	77.0%	66.6%	79.4%	-	-	_
X-E-Speech	5.69%	69.6%	74.1%	68.9%	67.1%	78.7%	-	-	-
Ours (ESD only)	10.44%	76.3%	78.8%	76.5%	67.9%	81.2%	-	-	-
Hierspeech++	5.01%	35.7%	45.0%	37.0%	73.1%	82.0%	3.48%	80.3%	14.6%
Ours	7.98%	81.2%	73.8%	78.9%	81.2%	85.3%	8.84%	83.2%	7.4%

Table 2: Ablation results.  $N \rightarrow O$  indicates conversion from neutral to other emotions. (overall) indicates cross-emotion conversion.

Models	ESD			Expresso		LibriTTS		
1/104015	WER↓	ECA (N $\rightarrow$ O) $\uparrow$	ECA (overall) ↑	SECS ↑	EECS ↑	WER↓	SECS ↑	EER ↓
w/o F0 cross-attention	8.82%	68.3%	64.0%	79.6%	83.6%	7.42%	82.0%	11.1%
w/o Mix-LN	10.45%	72.9%	74.6%	78.8%	84.5%	19.83%	79.4%	13.9%
w/o $L_{sim}$	9.90%	<b>77.1</b> %	73.0%	80.1%	85.2%	10.35%	81.9%	9.7%
w/o global F0 and energy	9.38%	76.8%	75.4%	80.3%	84.4%	9.58%	80.4%	10.3%
w/ MMS features	5.01%	62.2%	60.5%	79.1%	83.4%	5.21%	80.3%	16.5%
Proposed	9.25%	77.1%	76.5%	80.6%	84.9%	9.51%	82.3%	9.3%

Table 3: Subjective evaluation results.

	nMOS	eMOS	sMOS
GT	$4.40 \pm 0.15$	$4.17 \pm 0.18$	$4.60 \pm 0.12$
Consistency-VC	$4.25 \pm 0.15$	$3.88 \pm 0.19$	$4.73 \pm 0.08$
X-E-Speech	$4.12 \pm 0.17$	$3.88 \pm 0.16$	$4.63 \pm 0.10$
Ours (ESD only)	$3.88 \pm 0.17$	$4.00 \pm 0.16$	$4.50 \pm 0.13$
GT	$4.50 \pm 0.12$	$4.27 \pm 0.16$	$4.66 \pm 0.12$
Hierspeech++	$2.68 \pm 0.17$	$2.94 \pm 0.16$	$2.97 \pm 0.19$
Ours (w/ F0 sum)	$3.41 \pm 0.17$	$3.51 \pm 0.16$	$3.74 \pm 0.17$
Ours (w/ MMS)	$3.41 \pm 0.17$	$3.46 \pm 0.17$	$3.58 \pm 0.18$
Ours (proposed)	$3.37 \pm 0.19$	$3.72 \pm 0.17$	$3.65 \pm 0.18$

Table 4: Cross-lingual expressive voice conversion results.

	English to German			German to English			
	WER	SSIM	ECA	WER	SSIM	ECA	
Hierspeech++	4.51%	72.8%	61.7%	7.03%	65.5%	15.3%	
Consistency-VC	3.27%	64.1%	45.3%	13.54%	73.2%	25.1%	
X-E-Speech	4.77%	59.8%	51.9%	43.70%	76.1%	47.7%	
Ours	10.59%	74.8%	76.1%	30.84%	76.8%	50.7%	

 $L_{sim}$  enhances cross-emotional conversion by effectively removing source speech style. Comparing content features from MMS [41] and mHuBERT, using discrete speech units as content embeddings instead of continuous ones eliminates leaked acoustic information from the source but degrades quality, as reflected in higher WER scores. Lastly, incorporating global F0 and energy embeddings improves the style embeddings and result in better speaker and emotion transfer across all metrics.

#### 4.4. Cross-Lingual Style Transfer

Considering that the style encoder is content-agnostic, and linguistic features were extracted from a multi-lingual model, our

model can perform cross-lingual VC (XVC) even when trained solely on English data. In order to evaluate it, we used ESD (English) and EmoDB [42] (German) to perform neutral-to-emotional conversions across different languages. We choose German to demonstrate the performance on an unseen language for all models. We use XVC version for Consistency-VC which is trained on multilingual data including ESD. As shown in Table 4, our model more effectively preserves emotion and speaker identity, even for German speakers unseen during training. However, the WER scores reveal a significant drop in intelligibility when source language is not included in training, observed on German to English conversions for all models.

# 5. Conclusion

In this work, we proposed a novel zero-shot EVC framework that enhances linguistic and acoustic feature disentanglement to particularly reduce source style leakage. Our approach integrates F0 injection with cross-attention, Mix-LN, mHuBERT-147 units, perturbation-based similarity loss, and style embeddings enriched with F0 and energy contours in a novel framework. Experimental results demonstrate that proposed framework improves disentanglement, mitigates source style leakage more effectively than baselines, and achieves superior emotion transfer while preserving speaker similarity. Future work will focus on enhancing intelligibility and improving cross-lingual performance through multilingual training data.

## 6. Acknowledgements

The authors gratefully acknowledge support from the German Federal Ministry of Education and Research (BMBF) under grant 01EF1803B (RELATER), European Union's Horizon research and innovation programme under grant 101135798 (Meetween), and KIT Campus Transfer GmbH (KCT) staff in accordance with the collaboration with Carnegie-AI.

## 7. References

- [1] A. Waibel, H. Steusloff, R. Stiefelhagen *et al.*, "Chil: Computers in the human interaction loop," 2005.
- [2] M. Schmidt, J. Niehues, and A. Waibel, "Towards an open-domain social dialog system," *Dialogues with Social Robots: Enable*ments, Analyses, and Evaluation, pp. 271–278, 2017.
- [3] A. Waibel, M. Behr, D. Yaman, F. I. Eyiokur, T.-N. Nguyen, C. Mullov, M. A. Demirtas, A. Kantarci, S. Constantin, and H. K. Ekenel, "Face-dubbing++: Lip-synchronous, voice preserving translation of videos," in *ICASSP Workshops*, 2023.
- [4] I. S. Ahmad, A. Anastasopoulos, O. Bojar, C. Borg, M. Carpuat, R. Cattoni, M. Cettolo, W. Chen, Q. Dong, M. Federico et al., "Findings of the iwslt 2024 evaluation campaign," arXiv preprint arXiv:2411.05088, 2024.
- [5] A. Waibel and C. Fuegen, "Simultaneous translation of open domain lectures and speeches," Jan. 3 2012, uS Patent 8,090,570.
- [6] L. Barrault, Y.-A. Chung, M. C. Meglioli, D. Dale, N. Dong, M. Duppenthaler, P.-A. Duquenne, B. Ellis, H. Elsahar, J. Haaheim *et al.*, "Seamless: Multilingual expressive and streaming speech translation," *arXiv preprint arXiv:2312.05187*, 2023.
- [7] K. Song, Y. Ren, Y. Lei, C. Wang, K. Wei, L. Xie, X. Yin, and Z. Ma, "Styles2st: Zero-shot style transfer for direct speech-tospeech translation," *Interspeech*, 2023.
- [8] K. Zhou, B. Sisman, and H. Li, "Limited data emotional voice conversion leveraging text-to-speech: Two-stage sequence-tosequence training," *Interspeech*, 2021.
- [9] J. Gao, D. Chakraborty, H. Tembine, and O. Olaleye, "Nonparallel emotional speech conversion," *Interspeech*, 2019.
- [10] Y. Cao, Z. Liu, M. Chen, J. Ma, S. Wang, and J. Xiao, "Nonparallel emotional speech conversion using vae-gan." in *Interspeech*, 2020.
- [11] F. Kreuk, A. Polyak, J. Copet, E. Kharitonov, T.-A. Nguyen, M. Rivière, W.-N. Hsu, A. Mohamed, E. Dupoux, and Y. Adi, "Textless speech emotion conversion using discrete and decomposed representations," *EMNLP*, 2022.
- [12] S.-H. Lee, H.-Y. Choi, H.-S. Oh, and S.-W. Lee, "Hiervst: Hierarchical adaptive zero-shot voice style transfer," *Interspeech*, 2023.
- [13] Z. Ning, Q. Xie, P. Zhu, Z. Wang, L. Xue, J. Yao, L. Xie, and M. Bi, "Expressive-vc: Highly expressive voice conversion with attention fusion of bottleneck and perturbation features," in *ICASSP*, 2023.
- [14] J. Kim, J. Kong, and J. Son, "Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech," in *ICML*, 2021.
- [15] J. Li, W. Tu, and L. Xiao, "Freevc: Towards high-quality text-free one-shot voice conversion," in *ICASSP*, 2023.
- [16] M. Z. Boito, V. Iyer, N. Lagos, L. Besacier, and I. Calapodescu, "mhubert-147: A compact multilingual hubert model," *Interspeech*, 2024.
- [17] S. Akti, T. N. Nguyen, Y. Liu, and A. Waibel, "Voice privacyinvestigating voice conversion architecture with different bottleneck features," in *Proc. SPSC 2024*, 2024, pp. 44–49.
- [18] R. Huang, Y. Ren, J. Liu, C. Cui, and Z. Zhao, "Generspeech: Towards style transfer for generalizable out-of-domain text-tospeech," *NeurIPS*, 2022.
- [19] B. Desplanques, J. Thienpondt, and K. Demuynck, "Ecapa-tdnn: Emphasized channel attention, propagation and aggregation in tdnn based speaker verification," in *Interspeech*, 2020.
- [20] K. Zhou, B. Sisman, R. Rana, B. W. Schuller, and H. Li, "Emotion intensity and its control for emotional voice conversion," *IEEE Transactions on Affective Computing*, vol. 14, no. 1, 2022.
- [21] T.-N. Nguyen, N.-Q. Pham, and A. Waibel, "Accent conversion using pre-trained model and synthesized data from voice conversion." in *Interspeech*, 2022.

- [22] T. N. Nguyen, S. Akti, N. Q. Pham, and A. Waibel, "Improving pronunciation and accent conversion through knowledge distillation and synthetic ground-truth from native tts," in *ICASSP*, 2025.
- [23] G. Rizos, A. Baird, M. Elliott, and B. Schuller, "Stargan for emotional speech conversion: Validated by data augmentation of end-to-end emotion recognition," in *ICASSP*, 2020.
- [24] B. Schnell, G. Huybrechts, B. Perz, T. Drugman, and J. Lorenzo-Trueba, "Emocat: Language-agnostic emotional voice conversion," arXiv preprint arXiv:2101.05695, 2021.
- [25] M. Baas, B. van Niekerk, and H. Kamper, "Voice conversion with just nearest neighbors," *Interspeech*, 2023.
- [26] L. Qu, T. Li, C. Weber, T. Pekarek-Rosin, F. Ren, and S. Wermter, "Disentangling prosody representations with unsupervised speech reconstruction," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023.
- [27] H. Guo, C. Liu, C. T. Ishi, and H. Ishiguro, "Using joint training speaker encoder with consistency loss to achieve cross-lingual voice conversion and expressive voice conversion," in ASRU, 2023.
- [28] —, "X-e-speech: Joint training framework of nonautoregressive cross-lingual emotional text-to-speech and voice conversion," in *Interspeech*, 2024.
- [29] J. Li and L. Zhang, "Zse-vits: A zero-shot expressive voice cloning method based on vits," *Electronics*, vol. 12, no. 4, 2023.
- [30] S.-H. Lee, H.-Y. Choi, S.-B. Kim, and S.-W. Lee, "Hierspeech++: Bridging the gap between semantic and acoustic representation of speech by hierarchical variational inference for zero-shot speech synthesis," arXiv preprint arXiv:2311.12454, 2023.
- [31] J. Kong, J. Park, B. Kim, J. Kim, D. Kong, and S. Kim, "Vits2: Improving quality and efficiency of single-stage text-to-speech with adversarial learning and architecture design," arXiv preprint arXiv:2307.16430, 2023.
- [32] J. Kong, J. Kim, and J. Bae, "Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis," Advances in neural information processing systems, vol. 33, 2020.
- [33] H. Zen, V. Dang, R. Clark, Y. Zhang, R. J. Weiss, Y. Jia, Z. Chen, and Y. Wu, "Libritts: A corpus derived from librispeech for textto-speech," *Interspeech*, 2019.
- [34] K. Zhou, B. Sisman, R. Liu, and H. Li, "Seen and unseen emotional style transfer for voice conversion with a new emotional speech dataset," in *ICASSP*, 2021.
- [35] G. Chen, S. Chai, G. Wang, J. Du, W.-Q. Zhang, C. Weng, D. Su, D. Povey, J. Trmal, J. Zhang et al., "Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio," arXiv preprint arXiv:2106.06909, 2021.
- [36] T. A. Nguyen, W.-N. Hsu, A. d'Avirro, B. Shi, I. Gat, M. Fazel-Zarani, T. Remez, J. Copet, G. Synnaeve, M. Hassid *et al.*, "Expresso: A benchmark and analysis of discrete expressive speech resynthesis," *arXiv preprint arXiv:2308.05725*, 2023.
- [37] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *ICML*, 2023.
- [38] Z. Ma, Z. Zheng, J. Ye, J. Li, Z. Gao, S. Zhang, and X. Chen, "Emotion2vec: Self-supervised pre-training for speech emotion representation," arXiv preprint arXiv:2312.15185, 2023.
- [39] S. R. Livingstone and F. A. Russo, "The ryerson audio-visual database of emotional speech and song (ravdess): A dynamic, multimodal set of facial and vocal expressions in north american english," *PloS one*, vol. 13, no. 5, 2018.
- [40] B. Van Niekerk, M.-A. Carbonneau, J. Zaïdi, M. Baas, H. Seuté, and H. Kamper, "A comparison of discrete and soft speech units for improved voice conversion," in *ICASSP*, 2022.
- [41] V. Pratap, A. Tjandra, B. Shi, P. Tomasello, A. Babu, S. Kundu, A. Elkahky, Z. Ni, A. Vyas, M. Fazel-Zarandi et al., "Scaling speech technology to 1,000+ languages," *Journal of Machine Learning Research*, vol. 25, no. 97, pp. 1–52, 2024.
- [42] F. Burkhardt, A. Paeschke, M. Rolfes, W. F. Sendlmeier et al., "A database of german emotional speech." in *Interspeech*, 2005.