

ZERO SHOT AUDIO TO AUDIO EMOTION TRANSFER WITH SPEAKER DISENTANGLEMENT

Soumya Dutta and Sriram Ganapathy

Learning and Extraction of Acoustic Patterns (LEAP) lab, Indian Institute of Science, Bangalore.
E-mail - {soumyadutta,sriramg}@iisc.ac.in

ABSTRACT

The problem of audio-to-audio (A2A) style transfer involves replacing the style features of the source audio with those from the target audio while preserving the content related attributes of the source audio. In this paper, we propose an efficient approach, termed as **Zero-shot Emotion Style Transfer (ZEST)**, that allows the transfer of emotional content present in the given source audio with the one embedded in the target audio while retaining the speaker and speech content from the source. The proposed system builds upon decomposing speech into semantic tokens, speaker representations and emotion embeddings. Using these factors, we propose a framework to reconstruct the pitch contour of the given speech signal and train a decoder that reconstructs the speech signal. The model is trained using a self-supervision based reconstruction loss. During conversion, the emotion embedding is alone derived from the target audio, while rest of the factors are derived from the source audio. In our experiments, we show that, even without using parallel training data or labels from the source or target audio, we illustrate zero shot emotion transfer capabilities of the proposed ZEST model using objective and subjective quality evaluations.

Index Terms: Speech Emotion modeling, Style Transfer, Disentangled Representation Learning.

1. INTRODUCTION

Artificial emotional intelligence [1] encompasses methods that enable machines to understand and interact with human expressions of emotions. The style transfer approach to manipulating emotion, given a source and target data sample, is the task of converting emotion in the source sample to match the emotional style of the target sample while retaining rest of the attributes of the source. While the task has shown promising results in image domain [2], the applications in audio domain is more challenging [3, 4]. In this paper, we explore the task of emotion style transfer in speech data.

Voice conversion of speech primarily explored converting the speaker identity of a voice [5]. However, speech also contains information about the underlying emotional trait of the speaker in varying levels [6]. The initial frameworks using Gaussian mixture model (GMM) [7], hidden Markov model [8] and deep learning [9] based conversion approaches have recently been advanced with generative adversarial networks (GAN) [10] and sequence-to-sequence auto-encoding models [11].

In many of the prior emotion conversion approaches, the emotion targets are treated as discrete labels. However, emotion is a fine-grained attribute which has varying levels of granularities [6]. Forcing the emotion attribute to a small number of discrete target labels may not allow the models to capture the wide range of diverse

and heterogeneous sentiments elicited in human speech. Hence, we argue that the most natural form of emotion conversion is to transfer the emotion expressed in a target audio to the source audio, a.k.a A2A emotion style transfer. This motivation is also echoed in a recent work on A2A style transfer [12]. In spite of these efforts, audio-to-audio (A2A) style transfer in zero shot setting (unseen speakers and emotion classes) is challenging.

On a separate front, representation learning of speech has shown remarkable progress in the recent years. The wav2vec [13] models have been improved with masked language modeling (MLM) objectives (for example, HuBERT [14]) in self-supervised learning (SSL) settings. The derivation of speaker representations have mostly been pursued with a supervised model [15]. Further, reconstructing speech from factored representations of speaker, content and pitch contours [16] has shown that models like Tacotron [17], AudioLM [18] and HiFi-GAN [19] allow good quality speech generation.

In this paper, we propose a framework called, zero shot emotion transfer - ZEST, which leverages the advances made in representation learning and speech reconstruction. The proposed framework decomposes the given audio into semantic tokens (using HuBERT model [14]), speaker representations (x-vectors [15]), and emotion embeddings (derived from a pre-trained emotion classifier). Inspired by speaker disentangling proposed for speech synthesis [20], we also perform a single step of speaker and emotion disentanglement in the embeddings. Since pitch (F0) is also a component that embeds content, speaker and emotion, we investigate a cross-attention based model for predicting the F0 contour of a given utterance. Using the three representations (speech, speaker and emotion) along with the predicted F0 contours, the proposed ZEST framework utilizes the HiFi-GAN [19] decoder model for reconstructing the speech. During emotion conversion, the proposed ZEST approach does not use text or parallel training and simply imports the emotion embedding from the target audio for style transfer.

The experiments are performed on emotion speech dataset (ESD) [21]. We also explore a zero shot setting, where an unseen emotion from a different dataset is used as the reference audio. Further, a setting where the source speech is derived from an unseen speaker is also investigated. We perform several objective and subjective quality evaluations and compare with benchmark methods to highlight the style transfer capability of the proposed framework. The key contributions from this work can be summarized as follows,

- Proposing a novel framework for predicting the pitch contour of a given audio file using the semantic tokens from HuBERT, speaker embeddings and emotion embeddings.
- Enabling speaker-emotion disentanglement using adversarial training.
- Illustrating zero shot emotion transfer capabilities from unseen emotion categories, novel speakers and content.

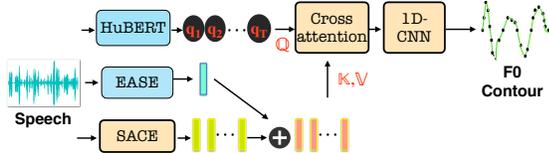


Fig. 1: Training of the F0 contour predictor. EASE - Emotion Agnostic Speaker Encoder. SACE - Speaker Adversarial Classifier of Emotions. The blue blocks are kept frozen during training.

2. RELATED PRIOR WORKS

A2A EST Using World Vocoder: One of the earliest attempts for EST involved using the world vocoder, as proposed by Gao et al. [22]. This work used the statistics of F0 and spectral components from the target speaker before reconstruction using the decoder. Our work uses recent advances in speaker, emotion and content embeddings for emotion style transfer. Further, we aim to transfer the emotion style from the reference speech to the source speech signal rather than just modifying the emotion category.

Expressive text-to-speech synthesis: The work by Li et al. [20] explored using speaker disentanglement for generating emotional speech from text. Similarly, Emovox proposed by Zhou et al. [6] used phonetic transcription for emotional voice conversion. However, our work explores EST without using any linguistic or phonetic transcriptions of the source or target speech.

Non-parallel and unseen emotion conversion: Recent work by Chen et al. [12] explored using attention models for performing EST. However, this work forced the source and target speech to be from the same speaker, limiting the utility of the EST applications.

3. METHOD

3.1. Content encoder

The content encoder used in the proposed framework is the HuBERT SSL model [14]. The HuBERT model gives continuous valued vector representations for each speech segment, which is subsequently converted into discrete tokens with a k-means clustering.

3.2. Emotion Agnostic Speaker Encoder (EASE)

The speaker embeddings for each audio file are extracted using an enhanced channel attention-time delay neural network (ECAPA-TDNN) [15] model. This model is pre-trained on 2794 hours and 7363 speakers from the VoxCeleb dataset [23], for the task of speaker classification. The model involves an utterance level pooling of the frame-level embeddings, called x-vectors. The x-vectors have been shown to encode emotion information [24, 25]. In order to suppress the emotion information, inspired by the disentanglement approach proposed in Li et al. [20], we add two fully connected layers to the x-vector model and further train the model with an emotion adversarial loss [26]. We refer to these vectors as the Emotion Agnostic Speaker Encoder (EASE) vectors. The loss function is given by

$$\mathcal{L}_{tot-spkr} = \mathcal{L}_{ce}^{spkr} - \lambda_{adv} \mathcal{L}_{ce}^{emo} \quad (1)$$

3.3. Speaker Adversarial Classifier of Emotions (SACE)

A speaker adversarial classifier of emotions (SACE) classifier is designed based on the wav2vec2.0 representations [13], similar to the one proposed by Pepino et al. [27]. The wav2vec model, pre-trained

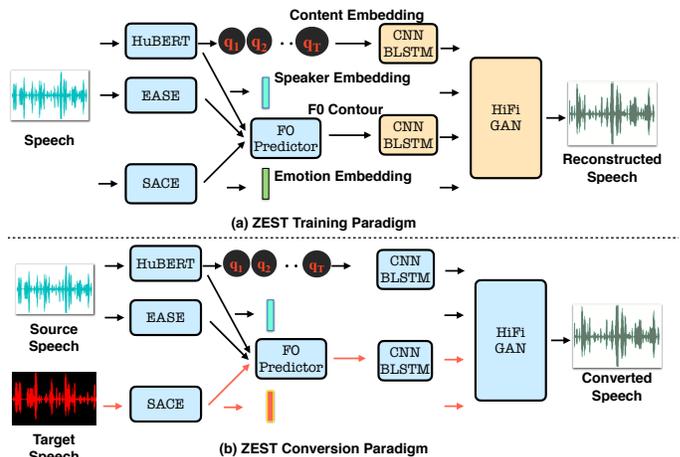


Fig. 2: (a) During training, ZEST is learned to reconstruct the speech signal. (b) During emotion conversion, the components that are derived from target speech are coded in orange color. The yellow blocks are the learnable parts of the model using an auto-encoding objective, while the blue blocks indicate frozen components.

on 300 hours (543 speakers) of switchboard corpus [28], is used for extracting features from the raw speech signal [29]. The convolutional feature extractors are kept frozen while the transformer layers along with two position wise feed forward layers are trained for the task of emotion recognition on the Emotional Speech Dataset (ESD) [21]. The model is trained with speaker adversarial loss (the emotion classifier equivalent of Eq. 1). The representations averaged over the entire utterance are used as the emotion embedding.

3.4. Pitch Contour Predictor

The framework for the pitch (F0) predictor is shown in Figure 1. The HuBERT tokens for the speech signal are converted to an sequence of vectors by means of an embedding layer. This sequence, denoted by $\mathcal{Q} = \{q_1, \dots, q_T\}_{i=1}^T$, is used as the query sequence for cross-attention. The frame-level SACE embeddings are added with speaker embedding (EASE) to form the key-value pair for the cross attention module [30]. This is followed by a 1D-CNN network to predict the F0 contour. The target pitch contour is the one derived using the YAAPT algorithm [31]. We use the \mathcal{L}_1 loss between the predicted and target F0 contour.

3.5. Speech reconstruction

The speech reconstruction framework is shown in Fig. 2(a). For reconstructing the speech signal, the HuBERT tokens, SACE embedding, EASE vector and the predicted F0 contour are used. The HuBERT tokens are converted to a sequence of real-valued vectors with a learnable embedding layer. In order to add contextual information during the speech reconstruction phase, the tokens and $F0_{pred}$ are passed through two separate networks consisting of CNN layers and bidirectional long-short term memory (BLSTM) layers. Finally, all the components are passed through a HiFi-GAN vocoder [19] to reconstruct the speech signal. More details of the HiFi GAN model are provided in Polyak et al. [16].

3.6. Emotion conversion

The ZEST framework for emotion conversion is shown in Figure 2(b). The HuBERT tokens and the speaker vector are extracted

from the source speech while the emotion embeddings are derived from the target speech. The emotion embedding sequence, being extracted from the reference speech, may differ in length from the source speech. However, as the query sequence in the cross-attention (Figure 1) is driven by the HuBERT tokens of the source signal, the F0 contour generated during conversion will match the length of the source signal. The HuBERT tokens, speaker vector, predicted F0 contour and the emotion embedding are then used to generate the converted speech through the pre-trained HiFi-GAN model. The conversion phase does not involve any model training steps.

4. EXPERIMENTS AND RESULTS

4.1. Datasets and Pretraining

The HuBERT model is the pre-trained 12-layer base model described in Hsu et al. [14]. This model is pre-trained on 960 hours of Librispeech dataset. The speaker encoder is initialized with the ECAPA-TDNN model [15] pre-trained on the VoxCeleb dataset. Further, the emotion agnostic training of the EASE model is performed on the Emotional Speech Database (ESD) [21]. We also train the F0 predictor as well as the HiFiGAN based reconstruction module on this dataset. We use the gender balanced English partition of this dataset in our experiments with 10 speakers. The training, validation and test splits are used as suggested in the dataset where, 300 utterances from each of the 10 speakers, for the five labelled emotions - neutral, angry, happy, sad and surprise, are used for training. Thus, the training data consists of 15000 utterances, while 2500 (50 utterances per speaker per emotion) unseen utterances are used for validation and testing.

We further evaluate the ZEST method on unseen emotions and speakers. For unseen emotion targets, we use 2542 utterances from the CREMA-D dataset [32] belonging to two new emotion classes (fear and disgust). For unseen source speakers, we investigate the use of the TIMIT database [33]. We choose 100 utterances distributed across 10 speakers selected at random from the TIMIT dataset for these experiments. Note that, the TIMIT and the CREMA-D datasets are only used in the conversion setting (Figure 2(b)). These evaluations reflect the zero-shot capability of the model.

4.2. Implementation

The HuBERT representations for the encoder are taken from 9-th layer of the pre-trained model and clustered with a k-means algorithm ($K = 100$) [16]. The EASE model is trained for a total of 10 epochs with λ_{adv}^{emo} set to 10 and batch size set to 32. The SACE setup is trained with a speaker adversarial loss ($\lambda = 1000$) along with the pitch predictor (Section 3.4) with a batch size of 24 and a learning rate of $1e - 4$, while the speech reconstruction model is trained with a batch size of 32 and a learning rate of $2e - 4$. The pitch predictor is trained for 50 epochs and the HiFi-GAN model is trained for 100K steps. The values of λ_{adv}^{emo} and λ are set based on the validation set performance of EASE and the pitch predictor respectively. Our code, converted audio samples and further ablations are available at this link¹.

We use the VAWGAN model proposed by Zhou et al. [34] as the baseline setup for bench-marking the proposed ZEST framework. As described in the paper, we supply the target emotion classes instead of the reference audio file for conversion. We also explore the

reconstruction model proposed in Polyak et al. [16] for benchmarking. Here, the F0 contour is derived from the target audio, while the speech and speaker contents are derived from the source audio.

4.3. Evaluation settings

- **Same-Speaker-Same-Text (SSST):** The source and reference speech are from the same speaker and with the same textual content. Further, the source speech is always in neutral emotion, while the reference speech can be from any other emotional category. This test set has a total of 1200 samples (30 neutral samples from the 10 speakers to be converted to 4 other emotion categories)
- **Same-Speaker-Different-Text (SSDT):** The source and reference speech are from the same speaker but with different textual content. This test set is created by randomly choosing 10 neutral utterances (one per speaker) as the source speech signals, and the remaining 29 test set utterances for conversion into each of the four target emotions. This leads to a total of 1160 test utterances.
- **Different-Speaker-Same-Text (DSST):** The source and reference speech are derived from different speakers using the same textual content. Each source speech signal (in neutral emotion) has 36 reference signals (the other 9 speakers having the same utterance in each of the 4 emotions). As each speaker has 30 test utterances, this test set has a total of 10800 utterances.
- **Different-Speaker-Different-Text (DSDT):** This is the most generic setting, where the source and reference speech do not share speaker or textual content. In this case, 10 neutral utterances are chosen as the source speech signal set (one per speaker). For each of these utterances, the remaining 29 utterances are selected as the reference. This results in a test set of 10400 samples.
- **Unseen Target Emotions (UTE):** The 2542 utterances from CREMA-D are considered as target speech signals for each of the 10 randomly chosen neutral source utterances from ESD dataset, resulting in a total of 25420 evaluation utterances.
- **Unseen Source Speakers (USS):** In this setting, the 100 utterances from TIMIT are chosen to be the source speech signals and 8 randomly chosen emotional utterances from ESD dataset (2 per emotion class barring neutral) are used as the reference speech signals. Each of the 8 utterances from ESD is from a different speaker. This results in a total of 800 evaluation utterances.

4.3.1. Objective evaluation

We use three objective metrics on the converted speech signals.

- **Emotion conversion accuracy:** We evaluate the target emotion class accuracy on the converted audio using the trained emotion classifier (Sec. 3.3)
- **Textual content preservation:** The converted speech is passed through an automatic speech recognition (ASR) system. This ASR system used the HuBERT encoder as a pre-training module [14] and was trained on the supervised 960h audio in the Librispeech [35] dataset using the CTC loss. The character error rate (CER) is computed for the converted audio with respect to the ground truth transcription of the source speech signal.
- **Speaker preservation:** We train a speaker classification model using the source speech and the ECAPA-TDNN embeddings to classify among the 10 training speakers. A two-layer feed-forward model is employed as backend for this task. The speaker recognition accuracy of the converted audio, with the source speaker as target, is used as measure of the speaker preservation property.

¹<https://github.com/iiscleap/ZEST>

Table 1: Objective evaluation results (%). Here, Emo. Acc. - Emotional Accuracy, CER - Character Error Rate, Spk. Acc. - Speaker accuracy. The different test settings are described in Sec. 4.3.

	Method	SSST	SSDT	DSST	DSDT	UTE	USS
CER ↓	VAWGAN [34]	7.9	6.6	7.9	6.6	-	-
	Polyak [16]	7.0	5.7	7.0	5.7	6.1	14.6
	ZEST-no-adv.	7.9	5.6	7.6	5.7	6.2	14.9
	ZEST-no-F0-pred.	8.6	6.3	8.5	6.3	7.1	14.9
	ZEST	7.2	5.6	7.2	5.4	5.9	14.8
Emo. Acc. ↑	VAWGAN [34]	41.4	44.5	41.4	44.5	-	-
	Polyak [16]	36.4	33.3	28.3	25.9	-	26.6
	ZEST-no-adv.	56.1	42.4	46.8	38.6	-	36.4
	ZEST-no-F0-pred.	55.8	39.9	45.3	36.6	-	28.4
	ZEST	69.0	63.8	60.0	55.6	-	72.5
Spk. Acc. ↑	VAWGAN [34]	62.4	76.0	62.4	76.0	-	-
	Polyak [16]	99.6	100	99.7	100	99.1	-
	ZEST-no-adv.	98.3	100	97.5	99.8	98.7	-
	ZEST-no-F0-pred.	96.8	99.4	93.4	97.9	97.6	-
	ZEST	99.4	99.8	97.8	99.7	98.6	-

The results for these objective tests are shown in Table 1. The following are the insights drawn from the objective evaluations.

- The emotion transfer accuracy for ZEST is seen to be significantly improved over the baseline systems compared here. Even in challenging conditions, like different speakers and different text content present in the target audio, and for unseen source speakers, the ZEST is seen to perform emotion transfer effectively.
- The CER results show that both the baseline system (Polyak et al. [16]) and the ZEST provide similar ASR performance for all the conditions. Even on the unseen target emotions (UTE), the ZEST is seen to preserve the CER. On the USS condition, both models show a degradation in the CER, potentially due to the mismatch in the training/test domain of HiFi-GAN model.
- The speaker accuracy of the VAWGAN is found to be inferior (SSST setting). The baseline model of Polyak et al. [16] allows near perfect speaker classification on all settings, while the proposed ZEST also matches this performance on all settings except the DSST. However, the model based on Polyak et al. [16] is not seen to be effective in transfer of emotions.
- We show two ablations of ZEST in Table 1, where either the emotion adversarial loss in the EASE module or the pitch predictor module is absent. The adversarial training allows disentangled representations of speaker and emotion, which is shown to improve the objective quality results. Further, the utility of EASE is seen from the improvement in the emotion accuracy across all the 6 test settings. The pitch prediction module also aids the system to achieve a better ASR performance.

4.4. Subjective tests

We conduct listening tests in order to judge the effectiveness of our method. We used Prolific² tool to setup the subjective listening tests. We recruited 20 participants to perform the subjective evaluation. We choose 52 recordings, with 8 recordings from each of the 4 test settings (SSST, SSDT, DSST, DSDT) and 10 recordings each from UTE, and USS settings. The recordings were presented in a random order. The participants were also provided with training examples to illustrate the objective of the test.

All the participants in the survey were asked to give their opinion score on the audio files (range of 1-5) based on three criteria

²<https://www.prolific.co>

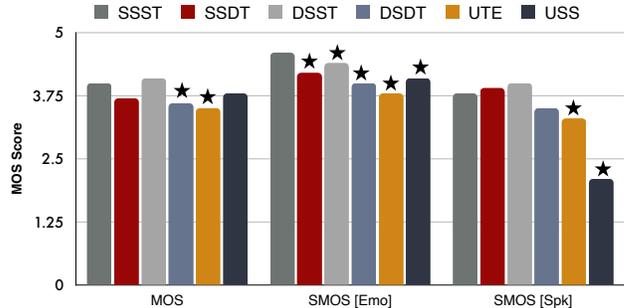


Fig. 3: Subjective evaluation on the different test settings. Abbreviations used: MOS- Mean Opinion Score, SMOS - Similarity Mean Opinion Score. The definition of the different test settings is given in Sec. 4.3. ★ indicates that the difference in scores from the SSST test setting is statistically significant ($p < 0.05$)

- i) Emotion transfer between the converted and the reference signal, ii) Quality of the converted speech and iii) Speaker similarity between the converted and the source signal. The subjective evaluation results (in terms of mean opinion score (MOS)) are reported in Figure 3. As seen in these results, the emotion transfer and reconstruction speech quality MOS values are the best for SSST condition. This is expected as the transfer of emotion from a reference speech of the same speaker with same textual content is the easiest test setting among all. However, the MOS results for the other challenging conditions also compare well with the SSST setting (except the speaker MOS on the unseen source speaker (USS) condition).

We also measure a statistical significance (unpaired t-test) between the scores obtained for the SSST test setting with each of the other 5 conditions. The emotion MOS for SSST is statistically significantly higher than all the other conditions. However on the speech quality, all the conditions, except UTE and DSDT, generate statistically similar results compared to SSST setting. For the speaker MOS scores, the unseen emotion targets and unseen speaker sources generate statistically different results compared to SSST. Part of the reason for this behavior may be attributed to the limited (10 speakers with same speech content) training employed for the reconstruction model (HiFi-GAN). In future, we plan to train the reconstruction model on a larger resource of neutral speech along with the emotional speech to improve the generalization to novel speakers. Further, leveraging larger emotional speech datasets for training the emotion embedding extractor and the F0 predictor may improve the quality of the emotion transfer to unseen emotions.

5. SUMMARY

We have presented an approach for zero shot emotion style transfer (ZEST) for audio-to-audio emotion conversion. The proposed ZEST method leverages pre-trained representations of speech content, speaker embeddings and emotion embeddings. Further, a pitch predictor model is designed in a self-supervised setting in order to learn the mapping from the representations to the pitch contour. A reconstruction module based on the HiFi-GAN allows the re-composition of the factored representations to generate the speech signal. For emotion conversion, we only derive the emotion embeddings from the target speech and perform the reconstruction. Various objective quality evaluation experiments with different levels of source/target mis-match, transfer to unseen emotions and from unseen source speakers elicit the zero shot emotion transfer capability of the proposed model. The subjective listening tests also validate the benefits of the proposed setting.

6. REFERENCES

- [1] Dagmar Schuller and Björn W Schuller, “The age of artificial emotional intelligence,” *Computer*, vol. 51, no. 9, pp. 38–46, 2018.
- [2] Yongcheng Jing et al., “Neural style transfer: A review,” *IEEE transactions on visualization and computer graphics*, vol. 26, no. 11, pp. 3365–3385, 2019.
- [3] Kaizhi Qian et al., “Autovc: Zero-shot voice style transfer with only autoencoder loss,” in *ICML*. PMLR, 2019, pp. 5210–5219.
- [4] Shrutina Agarwal, Sriram Ganapathy, and Naoya Takahashi, “Leveraging symmetrical convolutional transformer networks for speech to singing voice style transfer,” *arXiv preprint arXiv:2208.12410*, 2022.
- [5] Berrak Sisman et al., “An overview of voice conversion and its challenges: From statistical modeling to deep learning,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 132–157, 2020.
- [6] Kun Zhou et al., “Emotion intensity and its control for emotional voice conversion,” *IEEE Transactions on Affective Computing*, vol. 14, no. 1, pp. 31–48, 2022.
- [7] Ryo Aihara et al., “GMM-based emotional voice conversion using spectrum and prosody features,” *American Journal of Signal Processing*, vol. 2, no. 5, pp. 134–138, 2012.
- [8] Zeynep Inanoglu and Steve Young, “Data-driven emotion conversion in spoken english,” *Speech Communication*, vol. 51, no. 3, pp. 268–283, 2009.
- [9] Huaiping Ming et al., “Deep Bidirectional LSTM Modeling of Timbre and Prosody for Emotional Voice Conversion,” in *INTERSPEECH 2016*, 2016, pp. 2453–2457.
- [10] Georgios Rizos et al., “StarGAN for emotional speech conversion: Validated by data augmentation of end-to-end emotion recognition,” in *ICASSP*. IEEE, 2020, pp. 3502–3506.
- [11] Ravi Shankar et al., “Multi-speaker emotion conversion via latent variable regularization and a chained encoder-decoder-predictor network,” *arXiv preprint arXiv:2007.12937*, 2020.
- [12] Yun Chen et al., “Attention-based Interactive Disentangling Network for Instance-level Emotional Voice Conversion,” in *INTERSPEECH 2023*, 2023, pp. 2068–2072.
- [13] Alexei Baevski et al., “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *NeurIPS*, vol. 33, pp. 12449–12460, 2020.
- [14] Wei-Ning Hsu et al., “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.
- [15] Brecht Desplanques, Jenthe Thienpondt, and Kris Demuynck, “Ecapa-t2nn: Emphasized channel attention, propagation and aggregation in t2nn based speaker verification,” *arXiv preprint arXiv:2005.07143*, 2020.
- [16] Adam Polyak et al., “Speech resynthesis from discrete disentangled self-supervised representations,” *arXiv preprint arXiv:2104.00355*, 2021.
- [17] Jonathan Shen et al., “Natural tts synthesis by conditioning wavenet on mel spectrogram predictions,” in *ICASSP*. IEEE, 2018, pp. 4779–4783.
- [18] Zalán Borsos et al., “Audiolm: a language modeling approach to audio generation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023.
- [19] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae, “Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis,” *NeurIPS*, vol. 33, pp. 17022–17033, 2020.
- [20] Tao Li et al., “Cross-speaker emotion disentangling and transfer for end-to-end speech synthesis,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 1448–1460, 2022.
- [21] Kun Zhou et al., “Seen and unseen emotional style transfer for voice conversion with a new emotional speech dataset,” in *ICASSP*. IEEE, 2021, pp. 920–924.
- [22] Jian Gao et al., “Nonparallel emotional speech conversion,” *arXiv preprint arXiv:1811.01174*, 2018.
- [23] Arsha Nagrani et al., “Voxceleb: Large-scale speaker verification in the wild,” *Computer Speech & Language*, vol. 60, pp. 101027, 2020.
- [24] Raghavendra Pappagari et al., “x-vectors meet emotions: A study on dependencies between emotion and speaker recognition,” in *ICASSP*. IEEE, 2020, pp. 7169–7173.
- [25] Zein Shaheen et al., “Exploiting Emotion Information in Speaker Embeddings for Expressive Text-to-Speech,” in *INTERSPEECH 2023*, 2023, pp. 2038–2042.
- [26] Yaroslav Ganin et al., “Domain-adversarial training of neural networks,” *The journal of machine learning research*, vol. 17, no. 1, pp. 2096–2030, 2016.
- [27] Leonardo Pepino, Pablo Riera, and Luciana Ferrer, “Emotion recognition from speech using wav2vec 2.0 embeddings,” *arXiv preprint arXiv:2104.03502*, 2021.
- [28] John J Godfrey, Edward C Holliman, and Jane McDaniel, “Switchboard: Telephone speech corpus for research and development,” in *ICASSP*. IEEE Computer Society, 1992, vol. 1, pp. 517–520.
- [29] Soumya Dutta and Sriram Ganapathy, “HCAM–Hierarchical Cross Attention Model for Multi-modal Emotion Recognition,” *arXiv preprint arXiv:2304.06910*, 2023.
- [30] Ashish Vaswani et al., “Attention is all you need,” *NeurIPS*, vol. 30, 2017.
- [31] Kavita Kasi and Stephen A Zahorian, “Yet another algorithm for pitch tracking,” in *ICASSP*. IEEE, 2002, vol. 1, pp. 1–361.
- [32] Houwei Cao et al., “CREMA-D: Crowd-sourced emotional multimodal actors dataset,” *IEEE transactions on affective computing*, vol. 5, no. 4, pp. 377–390, 2014.
- [33] John S Garofolo et al., “DARPA TIMIT acoustic-phonetic continuous speech corpus CD-ROM. NIST speech disc 1-1.1,” *NASA STI/Recon technical report n*, vol. 93, pp. 27403, 1993.
- [34] Kun Zhou, Berrak Sisman, and Haizhou Li, “Vaw-gan for disentanglement and recombination of emotional elements in speech,” in *IEEE SLT*. IEEE, 2021, pp. 415–422.
- [35] Vassil Panayotov et al., “Librispeech: an asr corpus based on public domain audio books,” in *ICASSP*. IEEE, 2015, pp. 5206–5210.