ClapFM-EVC: High-Fidelity and Flexible Emotional Voice Conversion with Dual Control from Natural Language and Speech

Anonymous submission to Interspeech 2025

Abstract

Despite great advances, achieving high-fidelity emotional voice 2 conversion (EVC) with flexible and interpretable control re-3 mains challenging. This paper introduces ClapFM-EVC, a 4 novel EVC framework capable of generating high-quality con-5 verted speech driven by natural language prompts or refer-6 ence speech with adjustable emotion intensity. We first pro-7 8 pose EVC-CLAP, an emotional contrastive language-audio pretraining model, guided by natural language prompts and cat-9 egorical labels, to extract and align fine-grained emotional el-10 ements across speech and text modalities. Then, a FuEn-11 coder with an adaptive intensity gate is presented to seamless 12 fuse emotional features with Phonetic PosteriorGrams from a 13 pre-trained ASR model. To further improve emotion expres-14 siveness and speech naturalness, we propose a flow matching 15 model conditioned on these captured features to reconstruct 16 Mel-spectrogram of source speech. Subjective and objective 17 evaluations validate the effectiveness of ClapFM-EVC. 18 Index Terms: emotional voice conversion, natural language 19 prompt, contrastive language-audio pretraining, conditional 20

flow matching

21

22

1. Introduction

Emotional voice conversion (EVC) aims to convert the emotional state of source speech to a target category while preserving original content and speaker identity [1]. Recently, EVC
has garnered great attention within the speech processing realms
and holds great potential for many practical applications such as
voice assistant, audiobook production, and dubbing [2, 3, 4].

In general, the key challenges for EVC lie in the accurate 29 and efficient extraction, decoupling, and use of various speech 30 31 attributes, such as the emotion [5, 6], content [7, 8], and timbre [9, 10] information from speech signals. With rapid progress 32 in deep learning, existing EVC methods are generally based on 33 generative adversarial networks (GANs) [11, 12, 13] and au-34 toencoder models [14, 15, 16]. StarGAN [11] used a cycle-35 consistent and class-conditional GAN to achieve EVC. [16] pro-36 posed AINN, an attention-based interactive disentangling net-37 work with a two-stage pipeline for fine-grained EVC. However, 38 the converted speech of these systems lacks emotional diversity, 39 which is crucial for realistic speech synthesis [17]. To this end, 40 several studies [16, 17, 18] shifted towards incorporating inten-41 sity control modules into EVC framework to allow more pre-42 cise manipulation of emotional expression. Emovox [18] dis-43 entangled speaker style and controlled emotional intensity by 44 encoding emotion in a continuous space. EINet [17] predicted 45 emotional class and intensity via an emotion evaluator and in-46 tensity mapper, incorporating controllable emotional intensity 47

48 to enhance naturalness and diversity of emotion conversion.

Despite impressive advances, these approaches still face 49 challenges. First, current EVC systems based on GANs and au-50 toencoders, while promising, have great potential for improve-51 ments in emotional diversity, naturalness, and speech quality 52 [17, 19]. Second, current methods typically rely on reference 53 speech or categorical text labels as conditions to control a lim-54 ited set of emotional expressions. Nevertheless, this paradigm 55 not only imposes constraints on the user experience, but restricts 56 the diversity of emotional expressions, while falling short in in-57 tuitiveness and interpretability of conveyed emotions. 58

To mitigate the aforementioned issues, this paper presents 59 ClapFM-EVC, an innovative any-to-one EVC framework that 60 enables flexible and intuitive control of emotion conversion in 61 a user-friendly manner. To elaborate, we first propose EVC-62 CLAP (Contrastive Language Audio Pretraining), which is 63 guided by both natural language prompts and emotional cate-64 gorical labels, so as to extract and align the emotion features 65 across speech-text modalities. Additionally, we introduce an 66 end-to-end voice conversion (VC) model, termed AdaFM-VC, 67 composed of pre-trained ASR model, FuEncoder and condi-68 tional flow matching (CFM) model. Using FuEncoder with an 69 adaptive intensity gate (AIG), AdaFM-VC is able to integrate 70 the captured emotional representations with Phonetic Posteri-71 orgrams (PPGs) from a pre-trained ASR model HybridFormer 72 [20], while allowing flexible control over the emotional inten-73 sity of the converted waveform. To further enhance emotional 74 expressiveness, naturalness and speech quality, we incorporate 75 a CFM-based decoder [21, 22] that samples the output of the 76 FuEncoder from random Gaussian noise and reconstructs the 77 Mel-spectrogram of the source speech. During inference, EVC-78 CLAP can generate target emotional embeddings based on the 79 given natural language prompt, and then AdaFM-VC lever-80 ages the target emotion vectors, source PPGs, and predefined 81 emotional intensity to reconstruct the target Mel-spectrogram, 82 which is ultimately converted into the target speech by a pre-83 trained vocoder [23]. Extensive experiments and ablation stud-84 ies demonstrate that our ClapFM-EVC significantly outper-85 forms several existing EVC approaches in terms of emotional 86 expressiveness, speech naturalness, and speech quality. 87

2. METHODOLOGY

2.1. System Overview

As illustrated in Fig. 1, ClapFM-EVC can be characterized as a conditional latent model, where the proposed EVC-CLAP, FuEncoder, CFM-based decoder, as well as pretrained ASR [20] and vocoder [23] models serve as its core components.

Similarly to [24, 25], we first train EVC-CLAP using a symmetric Kullback-Leibler divergence based contrastive loss (symKL-loss) along with soft labels derived from natural lan-

88 89

90

91

92

93



Figure 1: Overall training architecture of the proposed ClapFM-EVC framework.

guage prompts and their corresponding categorical emotion la-97 bels. In this way, EVC-CLAP can effectively extract and align 98 emotional representations across audio and text modalities, 99 while enabling ClapFM-EVC to capture fine-grained emotional 100 information conveyed by natural language prompts. Then, we 101 train the AdaFM-VC using the obtained emotional elements 102 and content representations extracted by EVC-CLAP and pre-103 trained HybridFormer, respectively. The FuEncoder within 104 AdaFM-VC facilitates the seamless integration of emotional 105 and content characteristics, with its AIG module explicitly con-106 trolling the intensity of emotional conversion. Concurrently, the 107 CFM model in AdaFM-VC samples the outputs of the FuEn-108 109 coder from random Gaussian noise, and, conditioned on the target emotional vector produced by EVC-CLAP, it generates 110 the Mel-spectrogram features of target speech. Finally, the 111 generated Mel-spectrogram features are fed into a pre-trained 112 113 vocoder to synthesize the converted speech.

During inference, it is worth noting that our ClapFM-EVC 114 framework provides three modes for obtaining the target emo-115 tional embeddings: (1) directly based on the provided reference 116 speech; (2) directly based on the given natural language emo-117 tional prompt; and (3) EVC-CLAP retrieves relevant data from 118 a pre-constructed high-quality reference speech corpus using 119 specified natural language emotional prompt, subsequently ex-120 tracting the target emotion elements from the retrieved speech. 121

122 2.2. Soft-Labels-Guided EVC-CLAP

Overall, the purpose of Emo-CLAP training is to minimize the
 distance between data pairs within the same class, while simul taneously maximizing the distance between pairs of data from
 different categories.

Assume that the input data pair is $\{X_i^a, X_i^y, X_i^p\}$, where X_i^a is the source speech, X_i^y and X_i^p denote its corresponding emotional label and natural language prompt, $i \in [0, N]$ and N is the batch size. Our EVC-CLAP first adopts a pre-trained HuBERT¹ [26] based audio encoder and a pre-trained XLM-RoBERTa² [27] based text encoder to compress X_i^a and X_i^p

into two latent variables
$$Z_a \in \mathbb{R}^{N \times D}$$
 and $Z_p \in \mathbb{R}^{N \times D}$, where D

equals 512, representing the hidden state dimension. Following this, we compute their corresponding similarity matrices S^a_{pred} and S^p_{pred} as: 136

$$S_{pred}^{a} = \varepsilon_{a} \times (Z_{a} \cdot Z_{p}^{-1})$$

$$S_{pred}^{p} = \varepsilon_{t} \times (Z_{p} \cdot Z_{a}^{-T})$$
(1)

where ε_a and ε_t are two learnable hyper-parameters, with their values empirically initialized to 2.3. Subsequently, we employ symKL-loss to train Emo-CLAP with the guidance of the soft labels $M_{GT}^s \in \mathbb{R}^{N \times N}$ derived from X_i^y and X_i^p .

$$M_{GT}^s = \alpha_e M_{GT}^y + (1 - \alpha_e) M_{GT}^p \tag{2}$$

where α_e is a hyper-parameter to adjust M_{GT}^y and M_{GT}^p , empiri-141 cally set to 0.2 in our case. In detail, if the categorical emotional 142 labels or natural language prompt labels of different data pairs 143 within the same batch are identical, their corresponding ground 144 truth is assigned a value of 1; otherwise, it is set to 0. To en-145 sure the consistency of the label distributions across the batch, 146 the class similarity matrices M_{GT}^E and M_{GT}^P are normalized such 147 that the sum of each row equals 1, effectively capturing the rel-148 ative similarity between data pairs. Therefore, the training loss 149 of EVC-CLAP can be formulated as: 150

$$L_{symKL} = \frac{1}{4} \left(KL \left(S_{pred}^{a} || M_{GT}^{s} \right) + KL \left(\tilde{M}_{GT}^{s} || S_{pred}^{a} \right) + KL \left(S_{pred}^{p} || M_{GT}^{s} \right) + KL \left(\tilde{M}_{GT}^{s} || S_{pred}^{p} \right) \right)$$
(3)

$$\tilde{M}_{GT}^s = (1 - \alpha) \cdot M_{GT}^s + \frac{\alpha}{N} \tag{4}$$

153

154

$$KL(S||M) = \sum_{i,j} S(i,j) \log \frac{S(i,j)}{M(i,j)}$$
(5)

where α is a hyper-parameter, empirically set to 1×10^{-8} .

¹https://huggingface.co/TencentGameMate/chinese-hubert-large ²https://huggingface.co/FacebookAI/xlm-roberta-base

2.3. AdaFM-VC 155

2.3.1. FuEncoder with AIG 156

As a pivotal intermediate component within ClapFM-EVC, 157 FuEncoder aims to seamlessly integrate content features ex-158 tracted by HybridFormer with emotional embeddings derived 159 from EVC-CLAP, while offering flexible control over the emo-160 tion intensity through the adaptive intensity gate, namely AIG. 161 Detailed, FuEncoder comprises a preprocessing network 162 (PreNet), a positional encoding module, an AIG module, an 163 adaptive fusion module, and a linear mapping layer. First, 164 PreNet is presented to compress the source content features 165 Z_c to a latent space, preventing overfitting through a dropout 166 mechanism. Next, a positional encoding module is advocated 167 to employ sinusoidal positional encoding to extract the posi-168 tional characteristics of Z_c and performs element-wise addition 169 with Z_c to ensure that FuEncoder learns its sequential and struc-170 tural information. Afterwards, we propose an AIG module to 171 172 multiply a learnable hyperparameter by the EVC-CLAP's emotional features to flexibly adjust the emotional intensity. As the 173 core of FuEncoder, the adaptive fusion module consists of mul-174 tiple fusion blocks, each of which contains a multi-head self-175 attention layer, two emotion adaptive layer norm layers [28], 176 and a position-wise feed-forward network layer, enabling effi-177 cient fusion of content and emotional information, thus generat-178 ing rich embedding representations that contain both linguistic 179 and emotional characteristics. The fused features are ultimately 180 mapped to the specific dimensions $f \in \mathbb{R}^{B \times T \times D}$ through a fully 181 connected layer. 182

2.3.2. Conditional Flow Matching-based Decoder 183

To further enhance emotion expressiveness, speech natural-184 ness, and quality, we incorporate an optimal transport (OT)-185 based CFM model to reconstruct the target Mel-spectrogram 186 $x_1 = p_1(x)$ from a standard Gaussian noise $x_0 = p_0(x) =$ 187 $\mathcal{N}(x; 0, I)$. To elaborate, conditioned on the captured EVC-188 CLAP's emotional embeddings, an OT flow $\psi_t : [0,1] \times \mathbb{R}^d \to$ 189 \mathbb{R}^d is adopted to train our CFM-based decoder, which consists 190 191 of 6 CFM blocks with timestep fusion. Each CFM block contains a ResNet [29] module, a multi-head self-attention [30] 192 module and a FiLM [31] layer. By utilizing an ordinary differ-193 ential equation to model a learnable and time-dependent vector 194 field $v_t: [0,1] \times \mathbb{R}^d \to \mathbb{R}^d$, the flow can approximate the opti-195 mal transport path from $p_0(x)$ to the target distribution $p_1(x)$: 196

$$\frac{d}{dt}\psi_t(x) = v_t(\psi_t(x), t) \tag{6}$$

where $\psi_0(x) = x$ and $t \in [0, 1]$. Besides, drawing inspiration 197 from previous works [32], which suggest adopting straighter 198 trajectories, we simplify the OT flow formula as follows: 199

$$\psi_{t,z}(x) = \mu_t(z) + \sigma_t(z)x \tag{7}$$

where $\mu_t(z) = tz$, $\sigma_t(z) = (1 - (1 - \sigma_{\min})t)$, and z repre-200 sents the random conditioned input. σ_{\min} denotes the minimum 201 standard deviation of the white noise introduced to perturb in-202 dividual samples, with its value empirically set to 0.0001. 203

Consequently, the training loss for AdaFM-VC is defined 204 205 as:

$$\mathcal{L} = \mathbb{E}_{t,p(x_0),q(x_1)} \| (x_1 - (1 - \sigma) x_0) - v_t(\psi_{t,x_1}(x_0) | h) \|^2$$
(8)

where $x_0 \sim p(x_0)$, $x_1 \sim q(x_1)$, $t \sim U[0, 1]$, $q(x_1)$ denotes the true 206 yet potentially non-Gaussian distribution of the data, h refers to 207 the conditional emotion embeddings extracted by EVC-CLAP. 208

3. EXPERIMENTS

3.1. Experimental Setups

3.1.1. Datasets

Since no open-source EVC corpus with comprehensive emo-212 tional natural language prompts is currently available, we lever-213 age an internally developed expressive single-speaker Mandarin 214 corpus for training the proposed ClapFM-EVC system. This 215 corpus encompasses 20 hours of speech data sampled at 24 kHz. 216 From this, we specifically selected 12,000 utterances represent-217 ing 7 original categorical emotion classes (neutral, happy, sad, 218 angry, fear, surprise, disgust). To ensure high-quality annota-219 tions, we enlisted 15 professional annotators to provide natural language prompts for the selected waveforms. 221

3.1.2. Implementation Details

In all experiments, the proposed EVC-CLAP and AdaFM-VC 223 models are trained end-to-end using 8 NVIDIA RTX 3090 224 GPUs. For training the EVC-CLAP model, the Adam opti-225 mizer is employed with an initial learning rate of 1×10^{-5} and 226 a batch size of 16. All models are trained for 40 epochs within 227 the PyTorch framework. Afterwards, the AdaFM-VC approach 228 is trained using the AdamW optimizer over 500,000 iterations 229 on the same GPU setup, with an initial learning rate of 2×10^{-4} 230 and a batch size of 32. During inference, the Mel-spectrogram 231 of the target waveform is sampled using 25 Euler steps within 232 the CFM-based decoder, with a guidance scale set to 1.0. 233

3.1.3. Evaluation Metric

To assess speech quality and emotion similarity of ClapFM-235 EVC, we perform subjective and objective evaluations. Regard-236 ing speech quality, a variety of objective metrics are used, in-237 cluding Mel-cepstral distortion (MCD), root mean squared error 238 (RMSE), character error rate (CER), and predicted MOS (UT-239 MOS). For emotion similarity, we employ a pretrained speech 240 emotion recognition model³ [33] to compute the cosine similar-241 ity of emotion embeddings between the converted and reference 242 speech, referred to as the emotion embedding cosine similarity 243 (EECS). The CER and UTMOS are calculated using pretrained 244 CTC-based ASR⁴ and MOS prediction⁵ approaches. Besides, 245 the subjective Mean Opinion Score (MOS) with a 95% confi-246 dence interval is used to measure naturalness (nMOS) and emo-247 tion similarity (eMOS). In practice, we invite 12 professional 248 raters to participate in the evaluation. The scoring scale ranges 249 from 1 to 5, with increments of 1, where higher scores indicate 250 better performance. The audio samples are available online⁶. 251

3.2. Main Results

3.2.1. EVC by Reference Speech

To evaluate the performance of ClapFM-EVC, we compare it 255 with several existing EVC methods, i.e., StarGAN-EVC⁷ [11], Seq2seq-EVC⁸ [34], and MixEmo⁹ [35]. Since these baselines 257

⁴https://huggingface.co/facebook/hubert-large-ls960-ft

⁵https://github.com/tarepan/SpeechMOS

⁶https://anonymous.4open.science/w/clapfm-evc-ACF1

⁷https://github.com/glam-imperial/EmotionalConversionStarGAN

⁸https://github.com/KunZhou9646/seq2seq-EVC

⁹https://github.com/KunZhou9646/Mixed_Emotions

209 210

211

220

222

234

256

252

253

254

³https://github.com/ddlBoJack/emotion2vec

 Table 1: Overall comparison results of the speech quality and emotion similarity between our proposed ClapFM-EVC system and other

 SOTA baseline methods using reference speech. nMOS and eMOS are presented with 95% confidence intervals.

Model	MCD (\downarrow)	RMSE (\downarrow)	$CER(\downarrow)$	UTMOS (†)	nMOS (†)	EECS (†)	eMOS (†)
StarGAN-EVC [11]	8.85	19.48	13.07	1.45	2.09 ± 0.12	0.49	1.97 ± 0.09
Seq2seq-EVC [34]	6.93	15.79	10.56	1.81	2.52 ± 0.11	0.54	2.23 ± 0.11
MixEmo [35]	6.28	13.84	8.93	2.09	2.98 ± 0.07	0.65	2.58 ± 0.13
ClapFM-EVC	5.83	10.91	6.76	3.68	$\textbf{4.09} \pm \textbf{0.09}$	0.82	$\textbf{3.85} \pm \textbf{0.06}$

employ the reference waveform to facilitate EVC, we first examine their performance using reference speech.

As evidenced in Table 1, our ClapFM-EVC consistently 260 attains state-of-the-art performance in both speech quality 261 and emotion similarity. With regard to emotion similarity, 262 ClapFM-EVC exhibits significant advancements over baseline 263 approaches, achieving notable relative improvements of at least 264 26.2% in EECS and 53.1% in eMOS, respectively. This in-265 266 dicates that our proposed ClapFM-EVC framework has a remarkable capability to precisely capture and effectively trans-267 fer the target emotional characteristics during emotional voice 268 conversion. Regarding speech quality, the experimental results 269 reveal that ClapFM-EVC exhibits superior results across multi-270 ple objective metrics. Detailed, ClapFM-EVC shows enhanced 271 performance by achieving the lowest values in MCD, RMSE, 272 and CER metrics, respectively. Moreover, subjective evalua-273 tion confirms that ClapFM-EVC gains the highest scores in both 274 nMOS and UTMOS, with relative improvements of 37.2% and 275 49.2%, respectively, over the best-performing baseline method. 276 These results underscore its exceptional capability in maintain-277 ing superior perceptual quality. 278

279 3.2.2. EVC by Natural Language Prompt

²⁸⁰ To compare the performance of ClapFM-EVC when using ref-

281 erence speech (Reference) versus natural language prompts

(Prompt), we further conduct an ABX preference test.



Figure 2: The ABX preference test results compare the Reference with Prompt.

As shown in Fig. 2, the first test aims to evaluate the emo-283 tional similarity between the converted speech driven by Ref-284 erence and Prompt. 47 participants were asked to rate speech 285 samples generated by Prompt, with Reference serving as the 286 benchmark, on a scale from -1 to 1. Here, -1 indicates that 287 the converted speech driven by Reference shows better emo-288 tion similarity, and 0 indicates no preference. The results 289 revealed that 57.4% of participants selected "no preference," 290 while 19.1% favored the "Prompt," suggesting that ClapFM-291 EVC can effectively control the emotional expression of con-292 verted speech via Prompt. In addition, we assess the quality 293 of the converted speech relative to ground truth samples. Par-294 ticipants were required to choose the converted speech sample 295 that is closer to ground truth in speech quality. As shown in the 296 figure, the preference rates for speech driven by Reference and 297 Prompt are 25.5% and 23.4%, showcasing that ClapFM-EVC is 298

able to achieve high-quality EVC driven by Prompt.

3.3. Ablation Study

299

320

To evaluate the contributions and validity of each component of 301

the proposed system, we conduct ablation studies. All results 302 are summarized in Table 2. 303

Table 2: Ablation results of the proposed ClapFM-EVC driven by natural language prompts. 'w/o emo label' denotes removing emotional categorical labels when training EVC-CLAP, 'w/o symKL' represents replacing symKL-loss with KL-Loss, 'w/o AIG' denotes removing the AIG module of AdaFM-VC.

Model	UTMOS (†)	nMOS (†)	EECS (†)	eMOS (†)
ClapFM-EVC	3.63	$\textbf{4.01} \pm \textbf{0.06}$	0.79	$\textbf{3.72} \pm \textbf{0.08}$
w/o emo label	3.61	3.96 ± 0.11	0.66	3.01 ± 0.07
w/o symKL	3.57	3.89 ± 0.05	0.71	3.28 ± 0.08
w/o AIG	3.25	3.62 ± 0.12	0.74	3.52 ± 0.05

From the table above, we can easily reach the following 304 conclusions: (1) In the absence of categorical emotional la-305 bels, the EECS and eMOS scores exhibit a significant decline, 306 while the speech quality metrics remain largely unaffected. This 307 demonstrates the effectiveness and rationality of the proposed 308 soft-label-guided training strategy in our EVC-CLAP. (2) When 309 training EVC-CLAP with KL-Loss, the EECS and eMOS val-310 ues showed a relative performance drop of 10.1% and 11.8%, 311 indicating that the proposed symKL-Loss can effectively en-312 hance the emotion representation capability of EVC-CLAP. (3) 313 The removal of the AIG module leads to a notable deteriora-314 tion in speech quality and a slight performance reduction in 315 emotional similarity. This underscores the critical role of the 316 proposed AIG module in adaptively integrating content and 317 emotional characteristics, thereby improving the overall perfor-318 mance of our proposed ClapFM-EVC system. 319

4. CONCLUSIONS

In this study, we propose ClapFM-EVC, an innovative and ef-321 fective high-fidelity any-to-one EVC framework that features 322 flexible and interpretable emotion control along with adjustable 323 emotion intensity. Specifically, the proposed ClapFM-EVC 324 initially employs EVC-CLAP to extract and align emotional 325 elements across audio-text modalities. To enhance the emo-326 tional representational capacity, we utilize symKL-Loss to train 327 the proposed EVC-CLAP model, guided by soft labels derived 328 from the natural language prompts and their corresponding cat-329 egorical emotion labels. Subsequently, we introduce AdaFM-330 VC based on the pre-trained ASR model, the proposed FuEn-331 coder and CFM-based decoder to achieve high-fidelity and flex-332 ible emotion voice conversion. Extensive experiments indicate 333 that our ClapFM-EVC is able to generate converted speech with 334 precise emotion control and high speech quality driven by nat-335 ural language prompts. 336 337

5. References

- [1] K. Zhou, B. Sisman, C. Busso, and H. Li, "Mixed emotion mod-338 elling for emotional voice conversion," computer, vol. 6, p. 7, 339 2022 340
- [2] Y. Zheng, R. Zhang, M. Huang, and X. Mao, "A pre-training 341 based personalized dialogue generation model with persona-342 343 sparse data," in Proceedings of the AAAI Conference on Artificial Intelligence, vol. 34, no. 05, 2020, pp. 9693-9700. 344
- [3] S. Chen, Y. Feng, L. He, T. He, W. He, Y. Hu, B. Lin, Y. Lin, 345 Y. Pan, P. Tan *et al.*, "Takin: A cohort of superior quality zero-shot speech generation models," *arXiv preprint arXiv:2409.12139*, 346 347 2024. 348
- [4] Z. Zhang, L. Li, G. Cong, H. Yin, Y. Gao, C. Yan, A. v. d. Hengel, 349 and Y. Qi, "From speaker to dubber: movie dubbing with prosody 350 and duration consistency learning," in Proceedings of the 32nd 351 ACM International Conference on Multimedia, 2024, pp. 7523-352 353 7532.
- [5] Y. Pan, Y. Yang, Y. Huang, J. Yao, J. Yin, Y. Hu, H. Lu, L. Ma, 354 355 and J. Zhao, "Msac: Multiple speech attribute control method for reliable speech emotion recognition," (No Title), 2023. 356
- 357 [6] Y. Pan, Y. Yang, H. Lu, L. Ma, and J. Zhao, "Gmp-atl: Genderaugmented multi-scale pseudo-label enhanced adaptive trans-358 fer learning for speech emotion recognition via hubert," arXiv 359 preprint arXiv:2405.02151, 2024. 360
- [7] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, 361 362 S. Wang, Z. Zhang, Y. Wu et al., "Conformer: Convolutionaugmented transformer for speech recognition," arXiv preprint 363 arXiv:2005.08100, 2020. 364
- [8] Y. Yang, Y. Pan, J. Yin, and H. Lu, "Lmec: Learnable multi-365 plicative absolute position embedding based conformer for speech 366 367 recognition," arXiv preprint arXiv:2212.02099, 2022.
- H. Wang, S. Zheng, Y. Chen, L. Cheng, and Q. Chen, "Cam++: A [9] 368 fast and efficient network for speaker verification using context-369 370 aware masking," arXiv preprint arXiv:2303.00332, 2023.
- [10] J. Yao, Y. Yang, Y. Lei, Z. Ning, Y. Hu, Y. Pan, J. Yin, H. Zhou, 371 H. Lu, and L. Xie, "Promptvc: Flexible stylistic voice conversion 372 373 in latent space driven by natural language prompts," in ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech 374 and Signal Processing (ICASSP). IEEE, 2024, pp. 10571-375 10575. 376
- [11] G. Rizos, A. Baird, M. Elliott, and B. Schuller, "Stargan for emo-377 tional speech conversion: Validated by data augmentation of end-378 to-end emotion recognition," in ICASSP 2020-2020 IEEE Inter-379 380 national Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2020, pp. 3502-3506. 381
- [12] K. Zhou, B. Sisman, and H. Li, "Transforming spectrum and 382 prosody for emotional voice conversion with non-parallel train-383 ing data," arXiv preprint arXiv:2002.00198, 2020. 384
- R. Shankar, J. Sager, and A. Venkataraman, "Non-parallel emo-385 [13] tion conversion using a deep-generative hybrid network and an 386 adversarial pair discriminator," arXiv preprint arXiv:2007.12932, 387 2020. 388
- [14] T.-H. Kim, S. Cho, S. Choi, S. Park, and S.-Y. Lee, "Emotional 389 voice conversion using multitask learning with text-to-speech," in 390 ICASSP 2020-2020 IEEE International Conference on Acoustics, 391 Speech and Signal Processing (ICASSP). IEEE, 2020, pp. 7774-392 7778 393
- [15] W. Lu, X. Zhao, N. Guo, Y. Li, J. Wei, J. Tao, and J. Dang, 394 395 "One-shot emotional voice conversion based on feature separation," Speech Communication, vol. 143, pp. 1-9, 2022. 396
- [16] Y. Chen, L. Yang, Q. Chen, J.-H. Lai, and X. Xie, "Attention-397 based interactive disentangling network for instance-level emo-398 tional voice conversion," arXiv preprint arXiv:2312.17508, 2023. 399
- [17] T. Qi, S. Wang, C. Lu, Y. Zhao, Y. Zong, and W. Zheng, "Towards 400 realistic emotional voice conversion using controllable emotional 401 intensity," arXiv preprint arXiv:2407.14800, 2024. 402

- [18] K. Zhou, B. Sisman, R. Rana, B. W. Schuller, and H. Li, "Emotion 403 intensity and its control for emotional voice conversion." IEEE 404 Transactions on Affective Computing, vol. 14, no. 1, pp. 31-48, 405 2022 406
- [19] H.-H. Chou, Y.-S. Lin, C.-C. Sung, Y. Tsao, and C.-C. Lee, "Toward any-to-any emotion voice conversion using disentangled diffusion framework," arXiv preprint arXiv:2409.03636, 2024.

407

408

409

410

411

412

413

414

418

419

420

421

422

423

424

428

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

- [20] Y. Yang, Y. Pan, J. Yin, J. Han, L. Ma, and H. Lu, "Hybridformer: Improving squeezeformer with hybrid attention and nsr mechanism," in ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2023, pp. 1-5.
- [21] J. Yao, Y. Yan, Y. Pan, Z. Ning, J. Ye, H. Zhou, and L. Xie, "Sta-415 blevc: Style controllable zero-shot voice conversion with condi-416 tional flow matching," arXiv preprint arXiv:2412.04724, 2024. 417
- [22] Y. Pan, Y. Yang, J. Yao, J. Ye, H. Zhou, L. Ma, and J. Zhao, "Ctefm-vc: Zero-shot voice conversion based on content-aware timbre ensemble modeling and flow matching," arXiv preprint arXiv:2411.02026, 2024.
- [23] S.-g. Lee, W. Ping, B. Ginsburg, B. Catanzaro, and S. Yoon, "Bigvgan: A universal neural vocoder with large-scale training," arXiv preprint arXiv:2206.04658, 2022.
- B. Elizalde, S. Deshmukh, M. Al Ismail, and H. Wang, "Clap [24] 425 learning audio concepts from natural language supervision," in 426 ICASSP 2023-2023 IÈEE International Conference on Acoustics, 427 Speech and Signal Processing (ICASSP). IEEE, 2023, pp. 1–5.
- [25] Y. Pan, Y. Hu, Y. Yang, W. Fei, J. Yao, H. Lu, L. Ma, and J. Zhao, 429 "Gemo-clap: Gender-attribute-enhanced contrastive language-430 audio pretraining for accurate speech emotion recognition," in 431 ICASSP 2024-2024 IEEE International Conference on Acous-432 tics, Speech and Signal Processing (ICASSP). IEEE, 2024, pp. 433 10 021-10 025. 434
- [26] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdi-435 nov, and A. Mohamed, "Hubert: Self-supervised speech represen-436 tation learning by masked prediction of hidden units," IEEE/ACM 437 transactions on audio, speech, and language processing, vol. 29, 438 pp. 3451-3460, 2021.
- [27] A. Conneau, "Unsupervised cross-lingual representation learning at scale," arXiv preprint arXiv:1911.02116, 2019.
- D. Min, D. B. Lee, E. Yang, and S. J. Hwang, "Meta-stylespeech: [28] Multi-speaker adaptive text-to-speech generation," in Interna-PMLR, 2021, pp. tional Conference on Machine Learning. 7748-7759.
- [29] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770-778.
- [30] A. Vaswani, "Attention is all you need," Advances in Neural Information Processing Systems, 2017.
- [31] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville, "Film: Visual reasoning with a general conditioning layer," in Proceedings of the AAAI conference on artificial intelligence, vol. 32, no. 1, 2018.
- [32] Y. Yang, Y. Pan, J. Yao, X. Zhang, J. Ye, H. Zhou, L. Xie, L. Ma, and J. Zhao, "Takin-vc: Zero-shot voice conversion via jointly hybrid content and memory-augmented context-aware timbre modeling," arXiv preprint arXiv:2410.01350, 2024.
- [33] 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.
- [34] K. Zhou, B. Sisman, and H. Li, "Limited data emotional voice conversion leveraging text-to-speech: Two-stage sequence-tosequence training," arXiv preprint arXiv:2103.16809, 2021.
- [35] K. Zhou, B. Sisman, R. Rana, B. W. Schuller, and H. Li, "Speech 465 synthesis with mixed emotions," IEEE Transactions on Affective 466 Computing, vol. 14, no. 4, pp. 3120-3134, 2022. 467