

EAD-VC: Enhancing Speech Auto-Disentanglement for Voice Conversion with IFUB Estimator and Joint Text-Guided Consistent Learning

Ziqi Liang^{1,2‡}, Jianzong Wang^{1‡}, Xulong Zhang^{1✉}, Yong Zhang¹, Ning Cheng¹, Jing Xiao¹

¹*Ping An Technology (Shenzhen) Co., Ltd.*

²*University of Science and Technology of China*

Abstract—Using unsupervised learning to disentangle speech into content, rhythm, pitch, and timbre for voice conversion has become a hot research topic. Existing works generally take into account disentangling speech components through human-crafted bottleneck features which can not achieve sufficient information disentangling, while pitch and rhythm may still be mixed together. There is a risk of information overlap in the disentangling process which results in less speech naturalness. To overcome such limits, we propose a two-stage model to disentangle speech representations in a self-supervised manner without a human-crafted bottleneck design, which uses the Mutual Information (MI) with the designed upper bound estimator (IFUB) to separate overlapping information between speech components. Moreover, we design a Joint Text-Guided Consistent (TGC) module to guide the extraction of speech content and eliminate timbre leakage issues. Experiments show that our model can achieve a better performance than the baseline, regarding disentanglement effectiveness, speech naturalness, and similarity. Audio samples can be found at <https://largeaudiomodel.com/eadvc>.

Index Terms—voice conversion, speech disentanglement, self-supervised learning, mutual information

I. INTRODUCTION

Voice Conversion (VC) involves transforming the vocal characteristics of an source speaker into those of a target speaker. This is achieved by altering the speech para-linguistic aspects (e.g. speaker identity, prosody), without compromising the original linguistic information. The maturity of VC has brought benefits to various industries [1]–[3].

With the advancement of deep learning, various voice conversion solutions have been proposed. Some approaches utilize auxiliary models such as ASR or TTS models to achieve VC [4]–[6]. In [7]–[13], researchers use GAN and VAE to produce speech resembling the target speakers, but training GAN-based models is typically challenging. Recently, a lot of VC systems have enabled speech representation of speaker-dependent and independent information [14]–[17]. These works decompose speech into speaker and content representations, ensuring not only distribution matching like GANs but also easy training as easily as VAEs. [17] proposed a method for disentangling content and speaker information from speech. They used a vector quantization-based method to eliminate speaker information in the content information, and then add the unseen speaker information in the decoding stage, which greatly

improves the model generalization for unknown speakers. [18] proposed to unify other components besides timbre and content into prosodic features, and proposed a VC model that uses F0 as a condition. But these methods can only decompose speech into speaker, content and prosody that achieve coarse-grained disentanglement, while the information of pitch and rhythm are still mixed together. SpeechFlow [19] uses multiple autoencoders to disentangle speech into four components of pitch, rhythm, timbre, and content by introducing three well-designed information Bottlenecks. SpeechSplit2.0 [20] relies on the architecture provided by SpeechFlow. By employing additional signal processing techniques, the speech can be disentangled without the need for laborious bottleneck tuning.

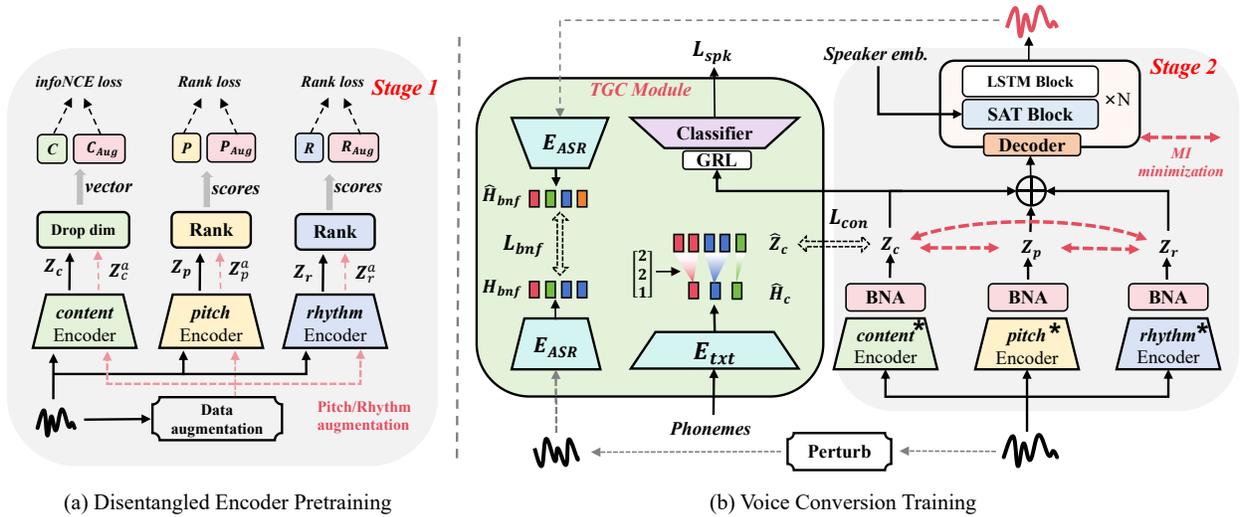
However, the model should be capable of independently distinguishing speech, thereby eliminating the need for manual extraction of bottleneck features. It not only conserves time but also minimizes the potential for bias and subjectivity that can stem from the manual selection of features. Liu et al. [21] realize that disentangling speech representation to the content, pitch, and rhythm by comparing the speech and its augmented version through ranking, which does not require bottleneck fine-tuning. However, the decoupling method based on data augmentation also has the same problem of insufficient decoupling, causing content, pitch, and rhythm to still be entangled. Moreover, without the guidance of linguistic information, the embedding extracted by the content encoder may be mixed with speaker information, causing timbre leakage and content inconsistencies.

To obtain the converted speech without human-crafted bottleneck tuning, we turn to the decoupling method based on self-supervised learning and focus on eliminating information overlap between different speech components. Thus, we propose the two-stage VC model for Enhanced speech Auto-Disentanglement (EAD-VC). The primary contributions of our work are as follows:

- Disentangling speech representations into content, pitch, rhythm, and timbre in a self-supervised manner without human-crafted bottleneck tuning.
- We design a new upper bound estimator IFUB of MI to enhance the decoupling between content, pitch, rhythm, and timbre. The bottleneck adapter (BNA) trained with MI with IFUB is designed to separate overlapping information between different speech components.

[‡] Equal Contributions

[✉] Corresponding author: Xulong Zhang (zhangxulong@ieee.org).



(a) Disentangled Encoder Pretraining

(b) Voice Conversion Training

Fig. 1: Framework of EAD-VC, which shows the two stages of our method: (I) Train the encoder based on the data and its augmented versions to disentangle speech as (a). (II) Freeze encoders to extract Z_c , Z_p , and Z_r in (b); Desired content embedding \hat{Z}_c from phonemes, which is used to guide the content encoder training. E_{ASR} is used to keep the content consistent after VC.

- We propose a joint text-guided consistent (TGC) module to solve timbre leakage in the content extraction and avoid content inconsistencies after conversion under the guidance of text transcriptions and ASR-Bottleneck (ASR-BNFs).

II. RELATED WORK

With the advancement of deep learning, many voice conversion has explored many techniques employing VAEs or GANs to facilitate the transfer of speaker information. VAE-VC [22] accomplishes voice conversion by generating speaker-independent content embeddings through its encoder. CDVAE-VC [23] uses two VAEs to reconstruct two different speech features of straight spectra and mel-cepstral coefficients (MCCs) respectively. ACE-VC [24] incorporates an auxiliary speaker classifier after the decoder and prevents the classifier from correctly classifying speaker information. Influenced by style transfer techniques in computer vision, [8] and [7] introduced CycleGAN and StarGAN respectively to implement voice conversion. Later, some methods utilize TTS or ASR models as auxiliary models to implement voice conversion [25]–[27]. Some researchers have harnessed the power of GANs and VAEs to produce voice that closely resemble the voice of a target speaker. [12], [28]–[30], but GAN-based voice conversion models are usually difficult to train. Recently, many VC systems have achieved the decoupling of speaker representation and speech-independent representation (content representation) [14], [31]–[35]. Compared with traditional voice conversion, it can achieve no need for the paired source speaker and target speaker data during training. High-quality source and target speech inputs significantly enhance the performance of the disentangled-based VC method, leading to substantial advancements in both the fidelity and the resemblance. Current mainstream VC models usually focus on speech representation

disentanglement, aiming to disentangle speaker information and content information as much as possible [14], [31], [36]. AutoVC [14] combines the ideas of GAN [29] and CVAE [24] to decouple content information and speaker information, and achieves One-Shot voice conversion. At the same time, there are also some works that pay more attention to the prosody of speech. [18] unifies components other than timbre and content into prosody features, and proposes a VC model using F0 as a condition, which adjusts the Auto-Encoder through constraining bottleneck features to achieve prosody decoupling.

If the model simply focus on timbre and content but ignores other components, it may result in less natural and expressive generated speech. Therefore, [37] additionally extracts the rhythm and pitch components from speech to achieve more fine-grained information decoupling. [19], [20], [38] use multiple encoders and signal processing techniques to disentangle speech into pitch, rhythm, content, and timbre by introducing three well-designed information bottlenecks. Based on SpeechFlow [19], our method explores the use of self-supervised learning to decouple content, pitch, rhythm, and timbre information without manually extracting pitch and rhythm information as guidance. At the same time, we enhance the information constraints between decoupled components and alleviate the information leakage problem existing in the above work.

III. METHODOLOGY

A. SSL-based speech disentanglement

As in [19], [20], we employ three encoders in our model. Rank loss and contrastive learning are applied to extract pitch, rhythm, and content representations of speech in a self-supervised manner. Pitch shift and time stretch are applied to modify the speech pitch and rhythm. As shown in Fig.1(a), speech data y and augmented data y_{Aug} were send to the disentangled encoders. We get Z_c , Z_p and Z_r from y with the

three disentangled encoders, and Z_c^a , Z_p^a and Z_r^a from y_{Aug} . Afterward, we apply the Rank layer to map Z_r and Z_p into two individual scores R and P , which is inspired by [21], [39], [40]. Hyperparameter $\gamma \in (0, 1)$ is used to indicate the data augmentation intensity. $\gamma < 0.5$ means negative augmentations such as decreasing pitch or rhythm, while $\gamma > 0.5$ means positive augmentations such as increasing pitch or rhythm. $\gamma = 0.5$ indicates no augmentation is applied. To ensure that the encoders produce disentangled representations by recognizing this augmentation intensity, we first apply a sigmoid function on scores of pitch and rhythm:

$$S_r = \frac{1}{1 + e^{-(R - RAug)}} \quad S_p = \frac{1}{1 + e^{-(P - PAug)}} \quad (1)$$

then we get the rank loss based on rhythm scores S_r and pitch scores S_p :

$$\mathcal{L}_r = -\gamma^r \log(S_r) - (1 - \gamma^r) \log(1 - S_r) \quad (2)$$

$$\mathcal{L}_p = -\gamma^p \log(S_p) - (1 - \gamma^p) \log(1 - S_p) \quad (3)$$

Concurrently, the content is initially condensed into a vector C for subsequent processing. Furthermore, to ensure that the content encoder solely produces representations pertinent to the content, we implement a contrastive learning loss on both the original content C and its augmented counterpart C_{Aug} :

$$\mathcal{L}_{InfoNCE} = -\beta \cdot \log \frac{\text{sim}(C, C_{Aug})}{\text{sim}(C, C_{Aug}) + \sum_{x_{Neg}} \text{sim}(C, C_{Neg})} \quad (4)$$

where $\text{sim}(\cdot, \cdot)$ is the exponential dot product, with a temperature t . β is a decay coefficient, the initial value is 1.0. Ultimately, we initiate the pre-training process for the trio of disentangled encoders utilizing the loss:

$$\min_{E_c(\cdot), E_p(\cdot), E_r(\cdot)} \mathcal{L}_{enc} = L_r + L_p + L_{InfoNCE} \quad (5)$$

B. Mutual information with IFUB estimator

Since full fine-tuning could distort pre-trained features and lead to worse performance in the presence of large distribution shifts, we freeze the pre-trained encoder in Fig.1 (a) to preserve model decoupling capabilities. We design a trainable bottleneck adaptor (BNA), and add it to each frozen encoder, which was trained to achieve further disentangling. Moreover, we designed a new mutual information upper bound estimator IFUB combined with InfoNCE, which inspired by vCLUB [41]. We utilize BNA trained with IFUB of MI to eliminate information overlap between pitch, rhythm, and content and enhance the decoupling of speech components. During training, the gradient only updates BNA layers while retaining the pre-trained parameters from frozen encoders. We remove disentangled information overlap and enhance disentanglement by MI minimization.

While $Q_\theta(Y|X)$ in vCLUB [41], [42] can be represented by any neural network, it is common in practice to parameterize $Q_\theta(Y|X)$ using a Gaussian family. The potential reason for avoiding the use of MLP to parameterize $Q_\theta(Y|X)$ is that it is difficult to converge in CLUB, so it makes sense to design a function that can parameterize $Q_\theta(Y|X)$ with any

neural network and converge easily. Therefore, we integrate vCLUB with InfoNCE to design a new upper bound estimator IFUB. This integration involves initially using the trained $f(x_i, y_i)$ from InfoNCE to substitute $\log p(Y|X)$ in vCLUB for computing the value of $\hat{\mathcal{I}}(X, Y)$, and subsequently minimizing it to reduce MI. We depict the process of minimizing Mutual Information using IFUB in Algorithm 1.

Algorithm 1 MI minimization with IFUB estimator

Training:

for each iteration do

Sample $\{x_i, y_i\}_{i=1}^N$ from $E_{c,p,r}(x, y)$.

Update Critic $f_{c,p,r}(\cdot)$ by maximizing $\mathcal{L}(x_i, y_i)$:

$$\mathcal{L}(x_i, y_i) = \frac{1}{N} \sum_{i=1}^N \log \frac{\exp(f_{c,p,r}(x_i, y_i))}{\frac{1}{N} \sum_{j=1}^N \exp(f_{c,p,r}(x_i, y_j))}$$

for $i = 0$; $i < N$; $i++$ **do**

$$UB_{c,p,r}^i = f_{c,p,r}(x_i, y_i) - \frac{1}{N} \sum_{j=1}^N f_{c,p,r}(x_i, y_j)$$

end for

Update \mathbf{E}_c and \mathbf{E}_p by :

$$\hat{\mathcal{I}}(Z_c, Z_p) = \frac{1}{N} \sum_{i=1}^N UB_{c,p}^i$$

Update \mathbf{E}_p and \mathbf{E}_r by :

$$\hat{\mathcal{I}}(Z_p, Z_r) = \frac{1}{N} \sum_{i=1}^N UB_{p,r}^i$$

Update \mathbf{E}_p and \mathbf{E}_r by :

$$\hat{\mathcal{I}}(Z_c, Z_r) = \frac{1}{N} \sum_{i=1}^N UB_{c,r}^i$$

end for

During the training phase, the neural network $f_{c,p,r}(\cdot)$ and $E_{c,p,r}(\cdot)$ are trained alternately, and we parameterize $f_{c,p,r}(\cdot)$ using a fully connected layer. By minimizing \mathcal{L}_{MI} , we can reduce the correlation among various speaker-irrelevant speech representations (content, pitch, and rhythm) and realize enhanced speech auto-disentanglement. The estimated MI loss is as:

$$\min_{Ada(\cdot)} \mathcal{L}_{MI} = \hat{\mathcal{I}}(Z_c, Z_p) + \hat{\mathcal{I}}(Z_p, Z_r) + \hat{\mathcal{I}}(Z_c, Z_r) \quad (6)$$

C. Joint text-guided consistent learning

We propose a TGC module to solve the problem of timbre leakage in the content encoder, which contains four submodules: (I) Text2Content module with length regulator. (II) Shared Speech2Content module trained with CTC loss. (III) Adversarial speaker classifier. (IV) Timbre fusion.

Text2Content module: To guide the content encoder in generating speaker-independent content embedding, we utilize a text encoder that produces text embedding \hat{Z}_c from phonemes. Then we calculate content consistent loss:

$$Z_c = E_{Ada}(E_c(x_i^{sp})) \quad (7)$$

$$\hat{Z}_c = F_{dur}(E_{txt}(x_i^{phn})) \quad (8)$$

$$\min_{E_c(\cdot)Ada(\cdot)} \mathcal{L}_{con} = \frac{1}{N} \sum_{i=1}^N (\|Z_c - \hat{Z}_c\|_1) \quad (9)$$

where x_i^{sp} represents mel-spectrogram and x_i^{phn} represents phonemes. x_i^{sp} passes through the content encoder E_c and BNA layer E_{Ada} in turn to get content embedding Z_c . x_i^{phn} passes through the text encoder E_{txt} and duration predictor

F_{dur} to get text2content embedding \hat{Z}_c , which is obtained by phoneme alignment of \hat{H}_c according to the duration. The aim of L_{con} is to encourage the content encoder to produce speaker-independent content embedding from source speech with the guidance of text embedding from phonemes.

Speech2Content module: We use the shared ASR encoder E_{asr} [43] trained by CTC Loss to extract the ASR bottleneck features (ASR-BNFs) from converted audio and source audio with formant perturbing. ASR-BNFs is the bottleneck feature extracted from the penultimate layer of the ASR model trained based on CTC loss. It only contains linguistic information. We first perform a formant perturb DA_{per} on the x_i^{SP} and eliminate the timbre information to obtain \tilde{x}_i^{SP} . Both the perturbed audio \tilde{x}_i^{SP} and the synthesized audio \hat{x}_i^{SP} pass through the shared encoder E_{asr} to extract ASR-BNFs H_{bnf} and \hat{H}_{bnf} . Then we calculate the L1 Loss to optimize the decoder and BNA layers.

$$\tilde{x}_i^{SP} = DA_{per}(x_i^{SP}) \quad (10)$$

$$H_{bnf} = E_{asr}(\tilde{x}_i^{SP}) \quad \hat{H}_{bnf} = E_{asr}(\hat{x}_i^{SP}) \quad (11)$$

$$\min_{D(\cdot), Ada(\cdot)} \mathcal{L}_{bnf} = \frac{1}{N} \sum_{i=1}^N (\|\hat{H}_{bnf} - H_{bnf}\|_1) \quad (12)$$

The content information in utterance should be the same after VC, only the speaker-dependent information has changed. Finally, we get relatively pure content information and have good disentangling performance with ASR-BNFs' constraint.

Adversarial speaker classifier: We integrate a speaker classifier with a Gradient Reversal Layer (GRL) to remove the timbral information from the embedded content representation. The expectation is that the content encoder will be trained to reduce the incorporation of speaker-specific details. The gradient is first reversed by the GRL to aim for the diminishment of speaker-specific attributes within the content embedding, prior to its backward propagation towards the content encoder. The adversarial loss is as follows:

$$\min_{E_c(\cdot), Ada(\cdot)} \mathcal{L}_{adv} = \sum_{n=1}^N \mathbb{I}(id_{spk} == n) \log P_{spk}^n \quad (13)$$

where $\mathbb{I}(\cdot)$ is the indicator function, P_{spk}^n represents the probability that the speaker classifier's output, which indicates the probability of being classified as spk_n .

Timbre fusion: Inspired by [44], we design a decoder with Speaker-Attention (SAT) block, which learns the mapping between the timbre features F_{spk} and the fusion features Z by slightly modifying the cross-attention mechanism. The SAT is designed as follows:

$$Z_{con} = \mathbf{Con}(Z_c, Z_p, Z_r) \\ Z = \mathbf{Con}(Softmax((W_q Z_{con})(W_k F_{spk})^T) W_v F_{spk}, Z_{con}) \quad (14)$$

where $\mathbf{Con}(\cdot)$ means concatenation. Since it is difficult to accurately estimate speaker representations for unseen speakers, using inaccurate speaker representations as input to the decoder will lead to a mismatch between training and inference. The SAT block was designed to improve the generalization to

unseen speakers compared to concatenation [19]. Finally, the training loss of voice conversion is:

$$\mathcal{L} = \mathcal{L}_{recon} + \alpha_1 \mathcal{L}_{MI} + \alpha_2 \mathcal{L}_{con} + \alpha_3 \mathcal{L}_{adv} + \alpha_4 \mathcal{L}_{bnf} \quad (15)$$

IV. EXPERIMENTS

A. Experiment setup

We use the VCTK corpus [45], which are randomly split into 100, 3 and 6 speakers as training, validation and testing sets respectively. Each speaker has about 400 sentences, and the audio is downsampled to 16kHz. To perform objective and subjective tests, we randomly select 50 conversion pairs to generate synthesized samples from both models. 20 listeners participated in subjective evaluation and scored the naturalness and similarity of the test audio.

In stage 1, we pre-train our content, pitch, and rhythm encoders for 50k iterations, with the learning rate set to 5e-5. Additionally, the temperature t in Eq.4 is set to 0.1. Due to the presence of numerous variables, we employ cosine annealing [46] to adjust the learning rate and govern the gradient update step, mitigating the risk of falling to a local optimal solution. At the same time, we use learning rate warm-up [47] to slow down the overfitting phenomenon of the mini-batch model in the initial stage and maintain the stability of the distribution.

In stage 2, The BNA and decoder are trained using a learning rate of 1e-4 for 500k iterations and we use a pre-trained WaveNet as vocoder.

B. Evaluation

We compare our proposed method with other systems, including:

- SpeechFlow [19], which can decompose speech into pitch, rhythm, content, and timbre by introducing three carefully designed information bottlenecks;
- VQMIVC [42], which use vector quantization to generate content embedding and employ MI as the correlation metric to decompose speech;
- Liu et al. [21], which is the first method to automatically decouple speech into different components through data enhancement and rank modules without hand-crafted features.
- EAD-VC: The VC method we proposed using SAT to fuse timbre; EAD-VC(Con) is a method of concatenating timbre in the channel dimension which is the same as the previous work [19].

C. Subjective evaluation results

Listeners utilize subjective evaluation to gauge the naturalness of speech and the similarity of the speaker's voice in the converted speeches, which are produced by a variety of models. We use Mean Opinion Score (MOS) to describe the naturalness of various model outputs, and Speaker Similarity MOS (SMOS) to evaluate the likeness between the converted audio and the intended audio with 95% confidence intervals, including timbre and prosody. 20 listeners (10 males and 10

TABLE I: Subjective and Objective Evaluation results of different methods.

Methods	Many-to-Many VC					One-Shot VC				
	MOS \uparrow	SMOS \uparrow	MCD \downarrow	log F_0 PCC \uparrow	WER \downarrow	MOS \uparrow	SMOS \uparrow	MCD \downarrow	log F_0 PCC \uparrow	WER \downarrow
SpeechFlow [19]	3.56 \pm 0.08	3.04 \pm 0.11	6.37	0.686	23.5%	2.62 \pm 0.10	2.45 \pm 0.09	8.01	0.544	31.7%
VQMIVC [42]	3.70 \pm 0.13	3.61 \pm 0.09	5.46	0.829	16.9%	3.35 \pm 0.09	3.06 \pm 0.11	6.77	0.675	24.1%
Liu et al. [21]	3.62 \pm 0.09	3.20 \pm 0.16	7.25	0.670	25.6%	3.05 \pm 0.11	2.89 \pm 0.14	7.59	0.552	29.3%
EAD-VC(Con)	3.85\pm0.16	3.64\pm0.11	5.40	0.758	15.2%	3.47 \pm 0.16	3.05 \pm 0.17	6.45	0.606	22.8%
EAD-VC	3.83 \pm 0.13	3.63 \pm 0.12	5.21	0.793	14.6%	3.52\pm0.14	3.35\pm0.19	6.33	0.693	21.5%
w/o BNA& L_{MI}	3.66 \pm 0.13	3.38 \pm 0.17	5.96	0.749	16.4%	3.21 \pm 0.13	2.92 \pm 0.13	7.12	0.622	25.9%
w/o L_{con}	3.71 \pm 0.15	3.46 \pm 0.14	5.37	0.745	21.1%	3.36 \pm 0.09	2.98 \pm 0.18	6.96	0.675	28.6%
w/o L_{bnf}	3.80 \pm 0.11	3.58 \pm 0.12	5.23	0.762	18.4%	3.41 \pm 0.11	3.04 \pm 0.11	6.51	0.676	27.1%
w/o L_{spk}	3.75 \pm 0.09	3.51 \pm 0.11	5.31	0.737	17.6%	3.37 \pm 0.12	2.94 \pm 0.14	6.73	0.648	23.4%

females) are requested to assign scores within a range of 1 to 5 points respectively.

As depicted in Table I, EAD-VC(Con) enhances the similarity to target speakers and achieves higher speech naturalness than other systems in many-to-many VC. In one-shot VC scenario, the method of directly concatenating content, rhythm, pitch and timbre in the channel dimension will significantly reduce the naturalness and similarity of converted speech. The Timbre fusion method of timbre components concatenation is insufficient to simulate the timbres of various unseen speakers in real life. However, the performance drop of EAD-VC with SAT is smaller than that of other models including EAD-VC(Con). SMOS metrics in one-shot VC is improved by 0.3 on EAD-VC verifies that SAT can enhance the speaker similarity and the model generalization in the face of unseen speakers.

D. Objective evaluation results

As shown in Table I, we calculate mel-cepstrum distortion (MCD) and word error rate (WER) for converted audio. To evaluate pitch variations of the converted audio, the Pearson correlation coefficient (PCC) between F_0 of the source and the converted audio is calculated.

The EAD-VC and EAD-VC(Con) have the lowest WER among the various methods in preserving linguistic content. This shows that our method is very robust. At the same time, the lowest MCD among all methods demonstrates that our method improves the speaker similarity for less distortion between converted audio and target audio. By controlling the pitch variations, we can achieve high F_0 consistency in the audio, as demonstrated by the F_0 -PCC obtained from EAD-VC, which is higher than other methods except [42] in many-to-many VC. It proves that the converted audio generated by EAD-VC(Con) has a pitch contour that is more similar to the target audio. The F_0 -PCC of [42] is higher than EAD-VC(Con), we attribute this to the fact that the pitch embedding extracted from source speech in [42] is an external supervised label that is fed directly to the decoder without passing through the pitch encoder. However, in the face of unseen speakers in one-shot VC scenarios, [42] will degrade model performance more than EAD-VC which can disentangle pitch information from speech. It proves the advantages of SAT in One-shot VC compared to concatenation.

Furthermore, an additional evaluation involving a fake speech detection test is conducted in one-shot VC condition. Using

Resemblyzer toolkit¹, we randomly select 6 of the 10 real audios as groundtruth. The remaining four real audio samples along with the converted audio from diverse models will serve as the basis for evaluating timbral resemblance.

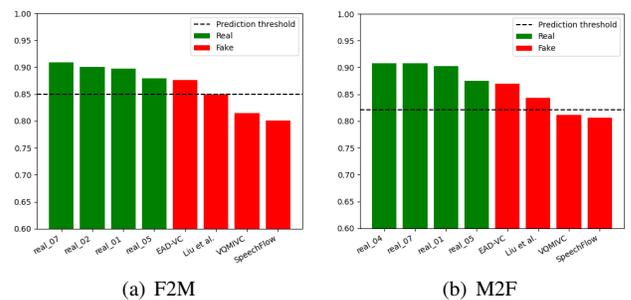


Fig. 2: Scores pertaining to VC are indicated. F denotes Female, and M denotes Male. Different models are represented on the x-axis, while prediction scores are represented on the y-axis.

As shown in Figure 2, the scores for real audios are indicated by the green clusters, whereas the red clusters correspond to the scores of the synthesized audios. The dashed line signifies the predictive threshold; any score surpassing this line is categorized as belonging to a real audio. Our proposed EAD-VC outperforms the other models on F2M and M2F VC by coming to highest scores over the dash-line among fake samples.

E. Generalization to unseen speaker

Fig.3 illustrates the different timbre embeddings visualized by the tSNE method. To prove the effectiveness of our method in enhancing speech representation auto-disentangling, we choose Liu et al. [21] as the baseline. The difference between the two is that EAD-VC uses timbre embeddings extracted from speaker encoder which jointly trained in the model and fused with other speech components, while the timbre embedding extracted from the pre-trained speaker encoder in [21] is directly concatenated with other speech components in the inference stage.

As shown in Fig.3(b), the better the clustering effect of embedding, the stronger the disentangling ability of VC model [15]. We achieve significant disentangling effects in one-shot VC. It demonstrates that EAD-VC can produce more clustered timbre representations on both VC and one-shot VC

¹Resemblyzer toolkit: <https://github.com/resemble-ai/Resemblyzer>

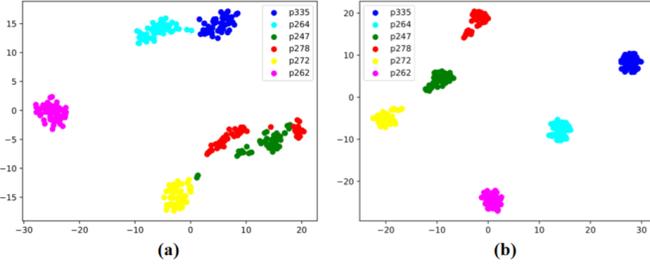


Fig. 3: Timbre embedding visualization on One-Shot VC. (a) Liu et al. [21]; (b) EAD-VC.

TABLE II: Ablation Study of P&R Decoupling Effectiveness.

Guidance	MOS \uparrow	SMOS \uparrow	MCD \downarrow
with F0	3.87 \pm 0.09	3.75 \pm 0.09	5.15
with Rhythm	3.67 \pm 0.11	3.41 \pm 0.13	6.11
with F0&Rhythm	3.72 \pm 0.08	3.52 \pm 0.11	5.76
Ours	3.83 \pm 0.13	3.63 \pm 0.12	5.21

and disentangle timbre information from speaker information. When we need to fine-tune the VC model on new unseen speakers, the encoding part only needs to update the parameters of the adaptor connected in series behind the disentangled encoder. Compared with other one-shot VC which require updating all model parameters, our method can reduce memory usage during training and improve computational efficiency.

F. Ablation study results

To validate the effect of our proposed method, we conduct ablation experiments. As shown in Table I, when the model is trained without the L_{MI} of bottleneck adaptor and (L_{con} , L_{bnf} , L_{spk}) of TGC module, both the quality and similarity scores drop when removing them. It still outperforms most of Liu et al., which demonstrates the efficiency of BNA&MI and TGC in improving speech naturalness and similarity.

When BNA& L_{MI} is removed, our method failed in almost all metrics in VC task, especially the MOS and SMOS. It means the L_{MI} is important in removing information overlap between decoupled speech components. What’s more, we observe that ASR performance (WER) and MCD drop significantly without using L_{con} , since the converted voice is compromised by undesired content information entangled with speaker representations. The existence of L_{con} can greatly reduce this timbre leakage problem. It also outperforms the model lacking L_{bnf} in terms of WER, which indicates the effectiveness of speech2content module in content inconsistency before and after voice conversion. Furthermore, it’s shown that the conversion quality of F0 is significantly degraded without the L_{spk} of the adversarial speaker classifier.

To assess the effectiveness of automatically disentangling pitch and rhythm through data augmentation and self-supervised learning, we employed the real pitch and rhythm extracted from speech as the label for supervised learning. We adopt this direct prediction approach to disentangle **Pitch** (F0) and **Rhythm** and compare it with our method in Table II.

When incorporating only real pitch as a guide during the training process (similar to VQMIVC), results show that

supervised training using real pitch (F0) as labels can improve naturalness and similarity to a certain extent. What’s more, the pre-train rhythm encoder [20] is used to extract rhythm embedding from speech as a guide for supervised learning. It proves that extracting rhythm information from speech is a challenging task to a certain extent as the disentangling effect drops significantly with rhythm guidance as shown in Table II. Our method can effectively disentangle rhythm from the latent space and does not require manual extraction of rhythm as a guide, which proves the feasibility of our method of disentangling pitch and rhythm.

G. Conversion rate

To demonstrate that our approach can enhance speech auto-disentanglement, we evaluate the conversion rate of the baseline [21] and EAD-VC. The conversion rate can serve as an indicator of the effectiveness of the disentangling process, which is followed by [19].

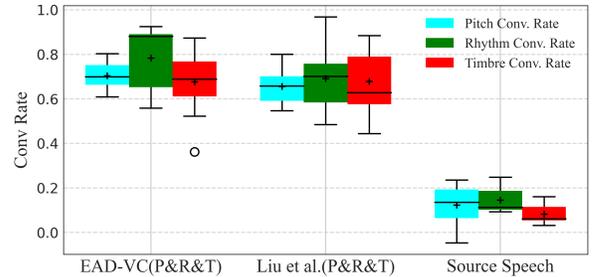


Fig. 4: Subjective conversion rate Evaluation. Each group encompasses three distinct subsections, which correspond to the conversion rates of pitch, rhythm, and timbre.

As shown in Fig.4, the conversion rate of our model exceeds the baseline in Pitch and Timbre. The findings indicate that while our disentangled encoders are structurally identical to the baseline, training them with BNA& L_{MI} and TGC modules enhances our model’s effectiveness in extracting disentangled pitch, rhythm, and timbre information.

V. CONCLUSIONS

In this paper, we introduce an innovative VC model with the designed IFUB estimator and joint text-guided consistent learning that can achieve SSL-based speech representation disentanglement. We use MI minimization with IFUB estimator to enhance the ability to speech auto-disentanglement into four components without the need for manual input of hand-crafted features and eliminate information overlap between components. Moreover, we use TGC module to solve timbre leakage and avoid content inconsistencies after VC. The experiments demonstrate that the EAD-VC which we have introduced, is capable of delivering superior disentanglement results and generate more authentic and natural speech.

VI. ACKNOWLEDGEMENT

Supported by the Key Research and Development Program of Guangdong Province (grant No. 2021B0101400003) and Corresponding author is Xulong Zhang (zhangxulong@ieee.org).

REFERENCES

- [1] H. Abdullah, W. Garcia, C. Peeters, P. Traynor, K. R. B. Butler, and J. Wilson, "Practical hidden voice attacks against speech and speaker recognition systems," in *NDSS*, 2019.
- [2] H. Abdullah, M. S. Rahman, W. Garcia, K. Warren, A. S. Yadav, T. Shrimpton, and P. Traynor, "Hear 'no evil', see 'kenansville'*: Efficient and transferable black-box attacks on speech recognition and voice identification systems," in *IEEE SP*, 2021, pp. 712–729.
- [3] S. Ahmed, I. Shumailov, N. Papernot, and K. Fawaz, "Towards more robust keyword spotting for voice assistants," in *USENIX Security Symposium*, 2022, pp. 2655–2672.
- [4] Z. Zhao, S. Ma, Y. Jia, J. Hou, L. Yang, and J. Wang, "Mix-guided VC: any-to-many voice conversion by combining ASR and TTS bottleneck features," in *ISCSLP*, 2022, pp. 96–100.
- [5] 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," *CoRR*, vol. abs/2211.04710, 2022.
- [6] X. Zhao, F. Liu, C. Song, Z. Wu, S. Kang, D. Tuo, and H. Meng, "Disentangling content and fine-grained prosody information via hybrid ASR bottleneck features for voice conversion," *CoRR*, vol. abs/2203.12813, 2022.
- [7] H. Kameoka, T. Kaneko, K. Tanaka, and N. Hojo, "Stargan-vc: non-parallel many-to-many voice conversion using star generative adversarial networks," in *IEEE SLT*, 2018, pp. 266–273.
- [8] F. Fang, J. Yamagishi, I. Echizen, and J. Lorenzo-Trueba, "High-quality nonparallel voice conversion based on cycle-consistent adversarial network," in *ICASSP*, 2018, pp. 5279–5283.
- [9] W. Huang, H. Luo, and H. Hwang, "Unsupervised representation disentanglement using cross domain features and adversarial learning in variational autoencoder based voice conversion," *IEEE Trans. Emerg. Top. Comput. Intell.*, vol. 4, no. 4, pp. 468–479.
- [10] Y. Saito, Y. Ijima, K. Nishida, and S. Takamichi, "Non-parallel voice conversion using variational autoencoders conditioned by phonetic posteriorgrams and d-vectors," in *ICASSP*, 2018, pp. 5274–5278.
- [11] H. Tang, X. Zhang, J. Wang, N. Cheng, Z. Zeng, E. Xiao, and J. Xiao, "Tgavc: Improving autoencoder voice conversion with text-guided and adversarial training," in *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2021, pp. 938–945.
- [12] X. Zhang, J. Wang, N. Cheng, E. Xiao, and J. Xiao, "Cyclegean: Cycle generative enhanced adversarial network for voice conversion," in *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2021, pp. 930–937.
- [13] H. Tang, X. Zhang, J. Wang, N. Cheng, and J. Xiao, "AVQVC: One-shot voice conversion by vector quantization with applying contrastive learning," in *ICASSP*, 2022, pp. 4613–4617.
- [14] K. Qian, Y. Zhang, S. Chang, X. Yang, and M. Hasegawa-Johnson, "AutoVC: Zero-shot voice style transfer with only autoencoder loss," in *ICML*, vol. 97, 2019, pp. 5210–5219.
- [15] J. Chou, C. Yeh, H. Lee, and L. Lee, "Multi-target voice conversion without parallel data by adversarially learning disentangled audio representations," in *Interspeech*, 2018.
- [16] J. Chou and H. Lee, "One-shot voice conversion by separating speaker and content representations with instance normalization," in *Interspeech*, 2019.
- [17] H. Tang, X. Zhang, J. Wang, N. Cheng, and J. Xiao, "Avqvc: One-shot voice conversion by vector quantization with applying contrastive learning," *ICASSP*, pp. 4613–4617, 2022.
- [18] K. Qian, Z. Jin, M. A. Hasegawa-Johnson, and G. J. Mysore, "F0-consistent many-to-many non-parallel voice conversion via conditional autoencoder," *ICASSP*, pp. 6284–6288, 2020.
- [19] K. Qian, Y. Zhang, S. Chang, and M. Hasegawa-Johnson, "Unsupervised speech decomposition via triple information bottleneck," in *ICML*, 2020.
- [20] C. H. Chan, K. Qian, Y. Zhang, and M. Hasegawa-Johnson, "Speech-split2.0: Unsupervised speech disentanglement for voice conversion without tuning autoencoder bottlenecks," in *ICASSP*, 2022.
- [21] Z. Liu, S. Wang, and N. Chen, "Automatic speech disentanglement for voice conversion using rank module and speech augmentation," in *Interspeech*, 2023.
- [22] C. Hsu, H. Hwang, Y. Wu, Y. Tsao, and H. Wang, "Voice conversion from non-parallel corpora using variational auto-encoder," in *APSIPA*, 2016, pp. 1–6.
- [23] W.-C. Huang, H.-T. Hwang, Y.-H. Peng, Y. Tsao, and H. Wang, "Voice conversion based on cross-domain features using variational auto encoders," *ISCSLP*, pp. 51–55, 2018.
- [24] H. Kameoka, T. Kaneko, K. Tanaka, and N. Hojo, "ACVAE-VC: non-parallel voice conversion with auxiliary classifier variational autoencoder," *IEEE ACM TASLP*, vol. 27, no. 9, pp. 1432–1443, 2019.
- [25] Z. Zhao, S. Ma, Y. Jia, J. Hou, L. Yang, and J. Wang, "Mix-guided vc: Any-to-many voice conversion by combining asr and tts bottleneck features," *ISCSLP*, pp. 96–100, 2022.
- [26] 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," *ICASSP*, pp. 1–5, 2023.
- [27] X. Zhao, F. Liu, C. Song, Z. Wu, S. Kang, D. Tuo, and H. M. Meng, "Disentangling content and fine-grained prosody information via hybrid asr bottleneck features for voice conversion," *ICASSP*, pp. 7022–7026, 2022.
- [28] H. Kameoka, T. Kaneko, K. Tanaka, and N. Hojo, "Stargan-vc: non-parallel many-to-many voice conversion using star generative adversarial networks," *IEEE SLT*, pp. 266–273, 2018.
- [29] F. Fang, J. Yamagishi, I. Echizen, and J. Lorenzo-Trueba, "High-quality nonparallel voice conversion based on cycle-consistent adversarial network," *ICASSP*, pp. 5279–5283, 2018.
- [30] X. Zhang, J. Wang, N. Cheng, and J. Xiao, "Voice conversion with denoising diffusion probabilistic gan models," in *19th International Conference on Advanced Data Mining and Applications*, 2023, pp. 154–167.
- [31] J. Chou and H. Lee, "One-shot voice conversion by separating speaker and content representations with instance normalization," in *Interspeech*, 2019, p. 664–668.
- [32] Y. Deng, H. Tang, X. Zhang, J. Wang, N. Cheng, and J. Xiao, "PMVC: data augmentation-based prosody modeling for expressive voice conversion," in *ACM MM*, 2023, pp. 184–192.
- [33] Y. Deng, J. Wang, X. Zhang, N. Cheng, and J. Xiao, "Learning expressive disentangled speech representations with soft speech units and adversarial style augmentation," in *IJCNN*, 2024.
- [34] Y. Deng, H. Tang, X. Zhang, N. Cheng, J. Xiao, and J. Wang, "Learning disentangled speech representations with contrastive learning and time-invariant retrieval," in *ICASSP*, 2024, pp. 7150–7154.
- [35] A. Polyak, Y. Adi, J. Copet, E. Kharitonov, K. Lakhota, W. Hsu, A. Mohamed, and E. Dupoux, "Speech resynthesis from discrete disentangled self-supervised representations," in *Interspeech*, 2021, pp. 3615–3619.
- [36] H. Lu, Z. Wu, D. Dai, R. Li, S. Kang, J. Jia, and H. M. Meng, "One-shot voice conversion with global speaker embeddings," in *Interspeech*, 2019.
- [37] W. Gan, B. Wen, Y. Yan, H. Chen, Z. Wang, H. Du, L. Xie, K. Guo, and H. Li, "Iqubbing: Prosody modeling based on discrete self-supervised speech representation for expressive voice conversion," *CoRR*, vol. abs/2201.00269, 2022.
- [38] S. Yang, M. Tantrawenith, H.-W. Zhuang, Z. Wu, A. Sun, J. Wang, N. Cheng, H. Tang, X. Zhao, J. Wang, and H. M. Meng, "Speech representation disentanglement with adversarial mutual information learning for one-shot voice conversion," in *Interspeech*, 2022.
- [39] S. Wang and D. Borth, "Zero-shot voice conversion via self-supervised prosody representation learning," in *IJCNN*, 2022, pp. 1–8.
- [40] Y. Souril, E. Noury, and E. Adeli, "Deep relative attributes," in *ACCV*, vol. 10115, 2016, pp. 118–133.
- [41] P. Cheng, W. Hao, S. Dai, J. Liu, Z. Gan, and L. Carin, "Club: A contrastive log-ratio upper bound of mutual information," in *ICML*, 2020.
- [42] D. Wang, L. Deng, Y. T. Yeung, X. Chen, X. Liu, and H. Meng, "VQMIVC: vector quantization and mutual information-based unsupervised speech representation disentanglement for one-shot voice conversion," in *Interspeech*, 2021, pp. 1344–1348.
- [43] Y. Zhang, H. Che, and X. Wang, "Non-parallel sequence-to-sequence voice conversion for arbitrary speakers," in *ISCSLP*, 2021, pp. 1–5.
- [44] D. Y. Park and K. H. Lee, "Arbitrary style transfer with style-attentional networks," in *CVPR*, 2019, pp. 5880–5888.
- [45] Y. J. M. K. Veaux Christophe, "Superseded - cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit," 2016.
- [46] I. Loshchilov and F. Hutter, "SGDR: stochastic gradient descent with warm restarts," in *ICLR*, 2017.
- [47] A. Gotmare, N. S. Keskar, C. Xiong, and R. Socher, "A closer look at deep learning heuristics: Learning rate restarts, warmup and distillation," in *ICLR*, 2019.