

Disentangling Style and Speaker Attributes for TTS Style Transfer

Xiaochun An, *IEEE Student Member*

Frank K. Soong, *IEEE Fellow*

Lei Xie*, *IEEE Senior Member*

Abstract—End-to-end neural TTS has shown improved performance in speech style transfer. However, the improvement is still limited by the available training data in both target styles and speakers. Additionally, degenerated performance is observed when the trained TTS tries to transfer the speech to a target style from a new speaker with an unknown, arbitrary style. In this paper, we propose a new approach to seen and unseen style transfer training on disjoint, multi-style datasets, i.e., datasets of different styles are recorded, one individual style by one speaker in multiple utterances. An inverse autoregressive flow (IAF) technique is first introduced to improve the variational inference for learning an expressive style representation. A speaker encoder network is then developed for learning a discriminative speaker embedding, which is jointly trained with the rest neural TTS modules. The proposed approach of seen and unseen style transfer is effectively trained with six specifically-designed objectives: reconstruction loss, adversarial loss, style distortion loss, cycle consistency loss, style classification loss, and speaker classification loss. Experiments demonstrate, both objectively and subjectively, the effectiveness of the proposed approach for seen and unseen style transfer tasks. The performance of our approach is superior to and more robust than those of four other reference systems of prior art.

Index Terms—Neural TTS, style transfer, disjoint datasets, variational inference, style and speaker attributes

I. INTRODUCTION

END-TO-END neural TTS models, such as Tacotron 2 [1], can produce high quality speech with naturalness close to that of human speakers [2]–[5]. The neural TTS models usually consist of an encoder-decoder neural network [6], [7] which is trained to map a given text sequence to a sequence of speech frames. Extensions of these models have shown that the *speech styles* (e.g., speaker identity, speaking style, emotion and prosody) can be modelled and controlled in an effective way [8]–[11]. Many TTS application scenarios, such as audiobook narration, news broadcasting, and conversational assistants, demand single-speaker, multi-style speech synthesis, i.e., a TTS to speak simultaneously in multiple styles. However, the corresponding performance in this area is inadequate due to the lack of single-speaker, multi-style speech data in general.

Currently most neural TTS models [12]–[16] are trained with an expressive, single-style corpus. Acquisition of a large set of single-speaker speech data with multiple styles, which is useful for training a good neural expressive TTS, is usually expensive and time consuming. Alternatively, a more effective solution is to perform *speech style transfer* [17], which allows a speaker to learn the desired style from the speech data in the same style but recorded by other speakers and preserves the target speaker’s timbre. The neural TTS models with a reference encoder [18], [19], global style tokens (GST) [20], and a variational autoencoder (VAE) [21] have become popular for controlling and transferring *speech styles*. Theoretically, these models can model any complex styles in a continuous latent space, hence one can control and transfer style by manipulating the latent variables or variational inference from a reference audio. However, these approaches tend to have too much entangled information to make the style rendering robust and interpretable, and independent control of specific speech characteristics (e.g., speaker identity and speaking style) is not clear and direct. In style transfer, one needs to transfer all styles whether desired or not, which may not fit the contexts thus hurts generalization. When conducting style control, one can hardly find the direct relationships between the target styles and the parameters in the embedding vector of the style representations to facilitate a direct control strategy.

To address the above issues, Bian *et al.* [22] introduce a multi-reference encoder to GST [20] to model multiple styles simultaneously and adopts intercross training to extract and separate different classes of speech styles. Occasionally, the model shows a successful style transfer, but the intercross training does not guarantee each possible combination of style classes is seen during training, leading to a missed opportunity to learn disentangled representations of styles and sub-optimal results on disjoint, multi-style datasets. In [23], the authors address the challenges of multi-reference style transfer on disjoint datasets by using an adversarial cycle consistency training scheme. Different from intercross training, their training tries all possible combinations of style classes via paired and unpaired triplets. Thus results in the disentanglement of multiple style dimensions and classes, and enables the style transfer to be more faithful than other existing methods.

Although the style transfer performance has been improved in [23], it is still limited to a style transfer from a speaker seen in training, but inadequate to transfer to a target style from a new speaker with an unknown, arbitrary style. In addition,

Xiaochun An and Lei Xie are with the Audio, Speech and Language Processing Group (ASLP@NPU), School of Computer Science, Northwestern Polytechnical University, Xi’an, China (email: xiaochunan@npu-aslp.org, lxie@npu.edu.cn)

Frank K. Soong is with the Speech Group, Microsoft Research Asia (MSRA), Beijing, China (email: frankkps@microsoft.com)

* Corresponding author

collecting training samples needed of new styles is always challenging and labor-intensive. Transferring style from one dataset to another (i.e., disjoint, multi-style datasets) is an appealing feature for a TTS system. Unseen style transfer on disjoint, multi-style datasets needs to be further improved.

In this paper, we propose an encoder-decoder neural network to improve *performance of seen and unseen style transfer* on disjoint, multi-style datasets. Our preliminary work has been presented in [24], in which how to distinguish different style types and capture the characteristics of individual speakers are not explicitly considered. In the current work, our preliminary work is extended with explicit constraints of different style types and different speakers to learn more discriminative style and speaker representations, which demonstrates improved discrimination of different styles and speakers. The modules for learning the style and speaker representations are sharpened in the current work, and the improvement has been confirmed in the Section V-E and Section V-F.

Our first contribution is to adopt an inverse autoregressive flow (IAF) to improve variational inference and learn discriminative style representations. Previous style embeddings [14], [21] are obtained by adopting the VAE [25] network. On the basis of the mean-field approximation, VAE assumes the independence of utterances and models them with a corresponding isotropic, Gaussian latent space. Despite the tractability of computation, the dimension-independent Gaussian distribution is not expressive enough, which has been investigated in various work [26], and utterances with the same style are not dependent but connected by sharing one global style space. Hence, we introduce an IAF in our proposed network to perform the variational inference. In the IAF model, posterior distributions are formulated by a series of cascaded invertible transformations to map a simple initial density to an arbitrarily complex, flexible distribution with tractable Jacobians [27]. Thus a flexible approximate distribution can be used for discriminative style embeddings. Note that IAF has been recently applied to speech processing tasks, e.g., fast and high-fidelity WaveNet based speech synthesis, Oord *et al.* [28] propose to use probability density distillation as a bridge between trained WaveNet (teacher model) and IAF (student model). Esling *et al.* [29] propose a universal audio synthesizer built with normalizing flows [30] to learn the latent space representation for semantic control of a synthesizer by interpolating latent variables. In this study, IAF is used as normalizing flows to perform the variational inference for learning discriminative and expressive style embeddings.

Our second contribution is to develop speaker encoder network for joint training of our proposed network to learn discriminative speaker representations. In recent years, researchers usually adopt a speaker recognition model [31] or a speaker verification model [32] to learn the speaker representations. Although they can extract the speaker embeddings, such speaker extractor models always need to be pretrained by using an independent dataset in advance. If the training data is inadequate, as shown in [33], it leads to poor speaker representations. Therefore, we introduce a

well-designed speaker encoder in our proposed network to train jointly with the rest network to learn discriminative speaker representations from any speakers, even if the speaker are not seen in the training data.

Our third contribution is to use six specifically-designed losses in network training. The style transfer accuracy and speaker preservation are both considered in our proposed approach of seen and unseen speech style transfer. The network is trained to learn more discriminative style and speaker representations in a disentangled manner, by optimizing the tradeoff between style transfer accuracy and speaker identity preservation through the six specifically-designed objectives: reconstruction loss, adversarial loss, style distortion loss, cycle consistency loss, style classification loss, and speaker classification loss. The reconstruction loss is used to measure the distortions in both source and target reconstructions; the adversarial loss to “fool” a well-trained discriminator; the style distortion loss to constrain the style representation of a source utterance to be closer to the target style representation; the cycle consistency loss to ensure that the transferred utterance can preserve the speaker identity of the source utterance; the style classification loss and the speaker classification loss to further improve their style and speaker representations and to make the category of each style and each speaker more distinguishable. In this way, we can transfer the speech to a target style from a new speaker with an unknown, arbitrary style, which does not need to be seen in training.

Subjective and objective experiments show that our approach to seen and unseen speech style transfer can improve 1) speech naturalness, 2) style similarity, and 3) speaker similarity, as compared with the four reference systems of prior art. Specially, on the unseen style transfer task, the reference systems, in most cases, fail to transfer an unseen style to a target style, and are not effective in preserving the speaker’s timbre, resulting in lower similarity scores.

To summarize our contributions, this paper proposes a novel approach to seen and unseen speech style transfer that can significantly improve performance of seen and unseen style transfer on disjoint, multi-style datasets. Specifically, the network is trained to minimize six specifically-designed losses to ensure the style representation of a source utterance is closer to the target style representation after the transfer and the transferred utterance can preserve the speaker identity of the source utterance, even when an utterance to be transferred is from a new speaker with an arbitrary, unknown style. In addition, the proposed scheme can be used as a data augmentation method to generate a single-speaker, multi-style speech data, which is useful for various speech applications, such as multi-style TTS and voice conversion.

The rest of this paper is organized as follows. Section II reviews the previous studies on single-reference and multi-reference speech style transfer. Section III presents our proposed approach to seen and unseen speech style transfer on disjoint, multi-style datasets. Section IV and Section V introduce the experimental setup and results, respectively. Section VI concludes this paper and discusses our future

work.

II. RELATED WORKS

Currently, there are mainly two major approaches, supervised and unsupervised to speech style transfer in TTS. The supervised approach, which takes manual style labels as additional TTS model input, has been shown effective in multi-speaker TTS [34]. However, this approach can hardly deal with more complex styles like different speaking styles, varying emotion levels and prosodic changes, etc., because there are no clean objective measures to annotate these styles. The unsupervised approach, which combines neural TTS models and a single-reference or multi-reference encoder, is more popular than the supervised approach. In this section, we will briefly review the related work on single-reference and multi-reference speech style transfer.

A. Single-reference Speech Style Transfer

Skerry-Ryan *et al.* extend the Tacotron architecture [2] by adding a reference encoder module [18], [35] that compresses the style of a variable-length audio signal into a fixed-length vector, as the reference embedding. Here, the reference encoder is to learn the embedding space of style from the speech data directly in the training process. The learned embedding, when used as a condition in synthesis, can generate speech signals with a style similar to that of the reference signal, even when the reference and target speakers are different.

Another popular single-reference speech style transfer is the GST model [20], which augments the reference encoder by utilizing a multi-head attention [36] based style token layer to extract rich style information in the training data. The extracted information is then used to control the synthesis, such as varying speed and speaking style. Similarly, it can be used for style transfer, replicating the speaking style of a single audio clip across a long-form text corpus.

Deep generative models, such as VAE [25] and GAN [37], [38], are powerful architectures which can learn complicated distributions in an unsupervised manner. Particularly, VAE, which can explicitly model the corresponding latent variables, has become one of the most popular schemes and achieved significant advancement in text generation [39], image generation [40], [41] and speech generation [12], [42]. Zhang *et al.* [21] introduce the VAE in the neural TTS model, to learn the latent representations of speaking style in an unsupervised manner. The learned style representations are then fed into a TTS network to control the style of the synthesized speech.

Theoretically, the above approaches can model any complex styles in a continuous latent space, so that one can control the transferred style by manipulating the latent variables or conducting a variational inference from a reference audio. However, these methods model all speech styles into one single representation, which is not versatile enough to control specific speech attributes independently. When conducting style transfer, one then has to transfer all the embedded styles, desired or not, which may not fit well with the contexts and

may hurt its generalization. In addition, these approaches fail to generalize to a new domain which is unseen in training. For example, to create speech in different speaker identities and speaking style classes by using a single model, a dataset containing audio samples for each speaking style class and speaker identity is needed, and yet the model can still fail to transfer the speaking style from a new speaker with an arbitrary, unknown style.

B. Multi-reference Speech Style Transfer

The single-reference speech style transfer has primarily focused on the transfer of a single-style reference audio sample. Those methods are inadequate for disjoint, multi-style datasets because of their lack of domain adaptation [43] capability.

Recently, Bian *et al.* [22] introduce a multi-reference encoder to GST [20] and adopt an intercross training scheme, to ensure that each sub-encoder of the multi-reference encoder disentangles and controls a specific style independently. The model shows successful style transfer in a multi-style scenario. However, its intercross training scheme does not guarantee each combination of style classes is seen in training, leading to a missed opportunity to learn disentangled representations of styles and sub-optimal results on disjoint datasets. To address the above problems, Whitehill *et al.* [23] propose an adversarial cycle consistency training scheme to ensure the use of information from all style classes. Different from intercross training, the scheme sweeps across all combinations of style classes via paired and unpaired triplets. This provides disentanglement of multiple style classes, enabling the model to transfer style in a more faithful manner than the existing methods.

However, similar to [22], the method of [23] still suffers a limitation that can only transfer the style seen in training, and is inadequate to transfer the speech to a target style from a new speaker with an unknown, arbitrary style, thus narrowing down applicable scenarios for neural TTS. In addition, recording training samples in a new style (e.g., poetry style) is challenging and labor-intensive, transferring style from one dataset to another (i. e., disjoint, multi-style datasets) is appealing for TTS systems.

More recently, the first author of this paper proposes an approach to perform seen and unseen style transfer on disjoint datasets [24]. Despite the realization of seen and unseen style transfer, the method cannot distinguish different style and speaker types well, especially for unseen style transfer, resulting in sub-optimal style transfer accuracy and speaker preservation. Hence, it is necessary to further enhance the performance of seen and unseen speech style transfer on disjoint, multi-style datasets.

III. PROPOSED APPROACH TO SEEN AND UNSEEN SPEECH STYLE TRANSFER

A. System Overview

Our approach to seen and unseen speech style transfer is built upon an encoder-decoder neural network [6], [7]. Fig. 1 illustrates our proposed framework for both

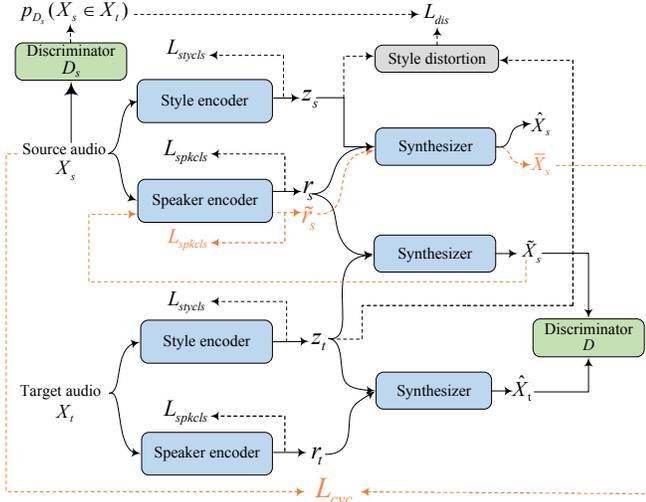


Fig. 1. The diagram of performing seen and unseen speech style transfer.

seen and unseen styles on disjoint, multi-style datasets. Let $X_s = \{x_s^{(1)}, \dots, x_s^{(n)}\}$ be the source utterances, and $X_t = \{x_t^{(1)}, \dots, x_t^{(m)}\}$, the target utterances, respectively. We assume that each speech utterance x can be decomposed into style representation, $z \in Z$, and speaker representation, $r \in R$. Each source utterance, $x_s^{(i)} \in X_s$, has its individual style, $z_s^{(i)}$, and the target utterance, $x_t^{(j)} \in X_t$, has the style, $z_t^{(j)}$. We use style encoder and speaker encoder, $E_z(x)$ and $E_r(x)$, to learn discriminative style representation z and speaker representation r of an utterance x , respectively: $z_s^{(i)} = E_z(x_s^{(i)})$, $r_s^{(i)} = E_r(x_s^{(i)})$, $z_t^{(j)} = E_z(x_t^{(j)})$, $r_t^{(j)} = E_r(x_t^{(j)})$. In this paper, we use Tacotron 2 [1] as the synthesizer T , which converts the combined encoder states (including style representations, speaker representations, and text encoder states) to generate a target Mel spectrogram with the target style and the target speaker's timbre. The Tacotron 2 is composed of a text encoder and a decoder with attention, whose details will be described in Section III-D. We then employ the LPCNet neural vocoder [44], [45] to reconstruct the final speech waveforms from the generated Mel spectrograms, and the details of each individual component in LPCNet vocoder will be described in Section IV-B. Our network is trained with six objective loss functions, including: reconstruction (\mathcal{L}_{rec}), adversarial (\mathcal{L}_{adv}), style distortion (\mathcal{L}_{dis}), cycle consistency (\mathcal{L}_{cyc}), style classification ($\mathcal{L}_{stylecls}$) and speaker classification (\mathcal{L}_{spkcls}) losses. Next, we introduce the learning of latent style space, speaker encoder network, synthesizer network and six objective losses.

B. Learning Latent Style Space

The style encoder $E_z(x)$ constructs a latent style space, and then outputs a sampled style representation z to the synthesizer T to guide style generation. As discussed above, the generated style space via VAE network [25], usually a Gaussian distribution, may not be expressive enough to transfer the style effectively. In this paper, we resort to the IAF [27], [46], a potent technique to construct sophisticated

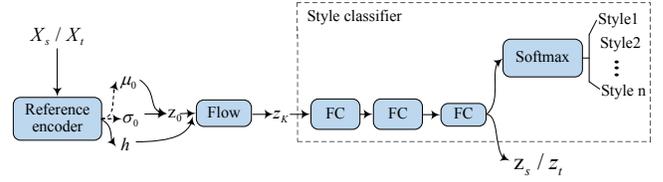


Fig. 2. Architecture of style encoder for learning style representations.

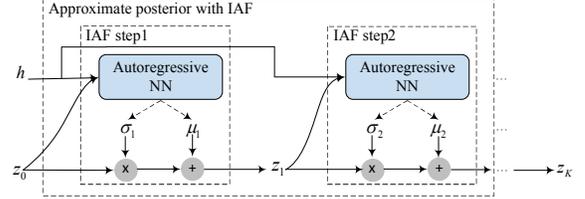


Fig. 3. The process of flow transformations.

distributions, to learn discriminative and expressive style representations. In other words, IAF can map a simple initial variable to a more complex one by applying a chain of cascaded invertible transformations. As shown in Fig. 2, we let a reference encoder network output μ_0 and σ_0 , in addition to each subsequent step in the flow. We draw a random sample $\epsilon \sim \mathcal{N}(0, I)$, and initialize the chain with:

$$z_0 = \mu_0 + \sigma_0 \odot \epsilon \quad (1)$$

Afterward, the initial variable z_0 along with hidden output h is provided to k steps of inverse autoregressive transformations to obtain flexible posterior probability distribution with latent variable z_k :

$$z_k = \mu_k + \sigma_k \odot z_{k-1} \quad (2)$$

In each step of the flow transformations, we adopt an autoregressive neural network with inputs z_{k-1} and h , and outputs μ_k and σ_k . And amortization is performed by using h as input to autoregressive networks of flow transformations [47]. These autoregressive transformations are invertible if $\sigma_i > 0$ condition is satisfied for i^{th} value of D dimension. Autoregressive structure of flow allows simple computation of the Jacobian determinant of each transformation as a change in global posterior probability density of reference encoder network denoted as $\log q(z_K|x)$, where z_K is output of the last flow step. Eq.(3) provides a tractable change of the probability density, and its detailed derivation can be found in [27]. The flexibility of the distribution of the final iteration z_K , and its ability to closely fit to the true posterior, increases with the expressivity of the autoregressive models and the depth of the chain. See Fig. 3 for an illustration of flow transformations.

$$\log q(z_K|x) = - \sum_{i=1}^D \left(\frac{1}{2} \epsilon_i^2 + \frac{1}{2} \log(2\pi) + \sum_{k=0}^K \log(\sigma_{k,i}) \right) \quad (3)$$

Different from [24], we then plug a style classifier to the style encoder shown in Fig. 2, which helps to learn more discriminative style embeddings for better discriminations of different styles. In detail, the classifier includes three fully

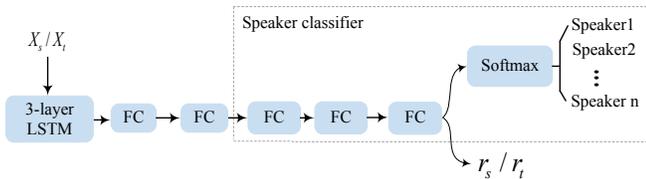


Fig. 4. Architecture of speaker encoder for generating speaker embeddings.

connected (FC) layers, all with ReLU activation, and a softmax layer to output the probability of different style types, i.e., reading, broadcasting, talking, story-telling, customer-service, poetry, and game styles. We use the output of the third FC layer as the style representation, z_s or z_t . In this way, we can obtain more discriminative and expressive style embeddings to guide style generation.

C. Developing Speaker Encoder Network

In previous works, speaker embeddings are usually obtained by using a speaker recognition model [31], [48] trained on the voxceleb corpus available in the Kaldi toolkit [49], [50] or a text-independent speaker verification task [32], [51], where the speaker extraction models need to be pre-trained. In our proposed network, we add a well-designed speaker encoder, $E_r(x)$, to learn the speaker representations. Different from [33], we propose to jointly train the speaker encoder network $E_r(x)$ with the rest of the neural TTS modules. We expect that the use of the speaker encoder can learn discriminative speaker embeddings to capture the characteristics of individual speakers, even the speaker is not seen in the training data.

Fig. 4 shows the architecture of speaker encoder in our approach, which maps a sequence of acoustic features of a speech utterance, into a fixed-dimensional embedding vector. First, we feed the features extracted from an utterance x into a 3-layer LSTM network. Two FC layers both with ReLU activation are connected to the last LSTM layer as an additional transformation of the last frame response of the network. Different from [24], a speaker classifier is added into speaker encoder to learn more discriminative speaker representations. As shown in Fig. 4, the classifier has the same structure as the style classifier, consisting of three FC layers, all with ReLU activation, and a softmax layer to output the probability of different speaker. Similarly, the output of the third FC layer is used as the speaker representation r_s or r_t , which is then used as a condition for the synthesizer T to guide speaker identity generation.

D. Synthesizer

In this paper, we leverage Tacotron 2 as the synthesizer T . Tacotron 2 is an attention-based sequence-to-sequence network [6], [7], [52], composed of a text encoder and a decoder with attention, which generates a Mel spectrogram as a function of an input text sequence and conditions the signal generated by the auxiliary encoder networks (e.g., style encoder and speaker encoder). It closely follows the network architecture of Tacotron 2 [1]. Input phonemes are represented using a learned 512-dim phoneme embedding, which is passed

through a stack of three convolutional layers each containing 512 filters with shape 5×1 , followed by a bidirectional LSTM of 256 units for each direction. In this work, to condition the output on additional attribute representations (e.g., style representations and speaker representations), the resulting text encodings are concatenated with the generated style embeddings z from the style encoder and the extracted speaker embeddings r from the speaker encoder, and then are accessed by the decoder through a location sensitive attention mechanism [53], which takes attention history into account when computing a normalized weight vector for aggregation.

The base Tacotron 2 autoregressive decoder network takes the attention-aggregated text encoding, and the bottlenecked previous frame (processed by a pre-net comprised of two FC layers of 256 units) at each step as input. The decoder input is then passed through a stack of two uni-directional LSTM layers with 1024 units. The output from the stacked LSTM is concatenated with the new decoder input (as a residual connection [54]), and linearly projected to predict the Mel spectrogram of the current frame, as well as the end-of-sentence token. Finally, the predicted spectrogram frames are passed to a post-net, which predicts a residual that is added to the initial decoded sequence of spectrogram frames, to predict the spectrograms by minimizing the overall mean squared errors.

E. Six Specifically-designed Objectives

For the synthesizer T , we form a reconstruction loss \mathcal{L}_{rec} to encourage the utterance from T , given style representation z and speaker representation r of an utterance x , to reconstruct x itself:

$$\begin{aligned} \mathcal{L}_{rec}(\theta_{E_z}, \theta_{E_r}, \theta_T) &= \mathbb{E}_{x_s \sim X_s} [-\log p_T(x_s | z_s, r_s)] \\ &+ \mathbb{E}_{x_t \sim X_t} [-\log p_T(x_t | z_t, r_t)] \end{aligned} \quad (4)$$

where θ denotes the parameter of the corresponding module. In this way, we can maintain reconstruction fidelity of an utterance with the synthesizer T .

Besides, for a sample x_s , we enforce the decoded sequence, given its speaker representation r_s and target style representation z_t , should be in the target domain X_t . Following GAN [37], [55], we introduce an adversarial loss \mathcal{L}_{adv} to be minimized in decoding and adopt a discriminator D , as shown in Fig. 1, to distinguish from $T(r_s, z_t)$ and $T(r_t, z_t)$. The task of the synthesizer is to fool the discriminator. Specifically, the adversarial loss \mathcal{L}_{adv} is defined as

$$\begin{aligned} \mathcal{L}_{adv}(\theta_{E_r}, \theta_{E_z}, \theta_T, \theta_D) &= \mathbb{E}_{x_s \sim X_s} [-\log(1 - D(T(r_s, z_t)))] \\ &+ \mathbb{E}_{x_t \sim X_t} [-\log D(T(r_t, z_t))] \end{aligned} \quad (5)$$

But, for a sample $x_s \in X_s$, its z_s can be an arbitrary value that minimizes the above reconstruction loss and adversarial loss, which may not necessarily capture the utterance style. This will affect the speaker representation, which is critical for representing the speaker, and it should be invariant against the transferred style. To address the issue, we introduce a

style distortion loss \mathcal{L}_{dis} to constrain style representation of a source utterance to be closer to the target style representation. As shown in Fig. 1, a discriminator, D_s , is first trained to predict whether a given utterance x has the target style with an output probability, $p_{D_s}(x \in X_t)$. When learning the style representation z_s , we then force the distortion between this style representation z_s and target style representation z_t to be consistent with the output probability of D_s . Here, we use the L_2 norm to measure the style distortion, $d(z_s, z_t) = \|z_s - z_t\|_2$, and make the style distortion positively correlated with $1 - p_{D_s}(x_s \in X_t)$. To incorporate this idea into our model, we model it with a standard normal distribution to evaluate the style distortion loss. Intuitively, when an utterance x_s have a large output probability $p_{D_s}(x_s \in X_t)$, our model can result in a small style distortion. That is, z_s will be closer to z_t , and the style distortion loss \mathcal{L}_{dis} is:

$$\mathcal{L}_{dis}(\theta_{E_z}) = \mathbb{E}_{x_s \sim X_s} [p_{D_s}(x_s \in X_t) d(z_s, z_t)^2] \quad (6)$$

where D_s is a pre-trained model trained with a portion of the training data. Here, if we integrate D_s into our training, we may start with a D_s with a low accuracy, and then our model is inclined to optimize a wrong style distortion loss for many epochs and gets stuck into a poor local optimum.

Unfortunately, the style distortion loss \mathcal{L}_{dis} can only constrain the generated utterance to be aligned with the target style, but cannot guarantee to keep the speaker identity of a source utterance intact. To address this issue, we introduce a cycle consistency loss [56], [57], \mathcal{L}_{cyc} , to our model shown in Fig. 1, which requires a transferred utterance to preserve the speaker identity of its source utterance, and enables the recovery of the source utterance in a cyclic manner. The cycle consistency loss \mathcal{L}_{cyc} is defined as

$$\begin{aligned} \mathcal{L}_{cyc}(\theta_{E_r}, \theta_{E_z}, \theta_T) \\ = \mathbb{E}_{x_s \sim X_s} [-\log p_T(x_s | E_r(\tilde{x}_s), z_s)] \\ + \mathbb{E}_{x_t \sim X_t} [-\log p_T(x_t | E_r(\tilde{x}_t), z_t)] \end{aligned} \quad (7)$$

where \tilde{x}_s is the transferred utterance from a source sample x_s and has the target style z_t . We encode the \tilde{x}_s with the speaker encoder $E_r(\tilde{x}_s)$ to obtain its speaker representation \tilde{r}_s , which is then combined with its source style z_s for decoding. Here, we expect that the source utterance can be generated with a high probability. For a target sample x_t , although we do not aim at changing its style in our model, similar to x_s , we still calculate its cycle consistency loss for additional regularization.

In addition, to learn more discriminative style and speaker representations that can distinguish different styles and capture the characteristics of different speakers, we introduce a style classification loss \mathcal{L}_{stycls} and a speaker classification loss \mathcal{L}_{spkcls} , as shown in Fig. 1. In this paper, we let the style classifier and the speaker classifier have the same architecture, and train the two classifiers with the cross-entropy loss, or the softmax loss. The details of the two classifiers will be described in Section IV-B. Specially, the softmax loss

$\mathcal{L}_{softmax}$ (i.e., the style classification loss \mathcal{L}_{stycls} and the speaker classification loss \mathcal{L}_{spkcls}) is formulated as

$$\mathcal{L}_{softmax} = - \sum_{i,j} y_{i,j} \log(\hat{y}_{i,j}) \quad (8)$$

where i refers to the style/speaker index of the input style/speaker embedding and j refers to the style/speaker index upon which the classification occurs. The $y_{i,j}$ is the ground-truth style/speaker class for the i -th embedding in the j -th style/speaker dimension, and $\hat{y}_{i,j}$ is the predicted style/speaker class. For $i = j$, the two classifiers encourage a style/speaker embedding to contain the correct information of the i -th style/speaker. For $i \neq j$, the two classifiers discourage the use of information about the other style/speaker.

By joint-training to minimize the weighted six specific objectives, or the overall objective function \mathcal{L}

$$\begin{aligned} \mathcal{L} = \alpha \mathcal{L}_{rec} + \beta \mathcal{L}_{adv} + \gamma \mathcal{L}_{dis} \\ + \lambda \mathcal{L}_{cyc} + \kappa \mathcal{L}_{stycls} + \omega \mathcal{L}_{spkcls} \end{aligned} \quad (9)$$

where $\alpha/\beta/\gamma/\lambda/\kappa/\omega$ are coefficients of the reconstruction term \mathcal{L}_{rec} in Eq.(4), the adversarial loss \mathcal{L}_{adv} in Eq.(5), the style distortion loss \mathcal{L}_{dis} in Eq.(6), the cycle consistency loss \mathcal{L}_{cyc} in Eq.(7), the style classification loss \mathcal{L}_{stycls} and the speaker classification loss \mathcal{L}_{spkcls} in Eq.(8), respectively.

IV. EXPERIMENTAL SETUP

A. Datasets and Comparison Models

We have carried out experiments to evaluate the performance of the proposed approach. We focus on disjoint, multi-style datasets, where datasets of different styles are recorded, and each style is recorded by one speaker with multiple utterances. Here, an internal Chinese corpus is used in experiments, which is divided into source data and target data. Source data contains examples of four styles: *reading* (standard reading speech), *broadcasting* (news broadcasting speech), *talking* (conversational speech) and *story-telling* (audiobook speech) styles, from four different speakers whereas target data contains samples of three styles: *customer-service* (spontaneous speech with fast speed), *poetry* (classical Chinese poetry speech with rich prosody variations), and *game* (exaggerated speech for role dubbing in the game) styles, from three other different speakers. This represents a minimalistic scenario of the disjoint, multi-style datasets: a single model must be able to properly transfer an arbitrary or unknown style to target style while keeping a minimal perceived change in the speaker's timbre. The corpus contains 22,212 samples (~ 25 hours) and each style contains 4,000 samples except for poetry and game styles. There are 698 samples in poetry style while 1514 samples in game style. For samples of each style, we respectively use 90% as the training set, 5% as the validation set and the rest 5% for the test. We remove long silence (> 0.1 sec) at the beginning and ending of each utterance. 80-dim Mel spectrogram is extracted as target speech representations with a Hanning window of 50 ms and 12.5 ms frame shift. Phoneme sequences are used as the input, which are obtained by text normalization and G2P pre-processing modules.

For all different systems in our experiments, we train $\sim 220k$ steps with a single Nvidia Tesla P40 GPU, and a batch size is 32. The models trained and tested in our experiments include:

- GST: we introduce “global style tokens” (GST) [20] into Tacotron 2 [1] to uncover expressive factors of variation in speaking style to perform the style transfer task, and make a fair comparison;
- VAE: we incorporate VAE [21] into Tacotron 2 [1] to learn the latent representation of speaking style to guide the style in synthesizing speech;
- MRF-IT: we augment a multi-reference encoder into GST-Tacotron 2 [22] to model multiple styles simultaneously, and adopt intercross training to extract and separate different classes of speech styles, thus achieving the transfer for desired speech styles;
- MRF-ACC: we adopt an adversarial cycle consistency training scheme for multi-reference neural TTS stylization [23] to ensure the use of information from all style classes, thus performing multi-reference style transfer on disjoint datasets;
- Proposed model: we introduce an IAF technique [27] to improve variational inference and learn expressive style representations, and develop a joint-training speaker encoder network to obtain discriminative speaker representations. Six loss functions with different purposes in network training are used together for enhancing the performance of seen and unseen style transfer on disjoint, multi-style datasets.

B. Model Details

The style encoder contains a reference encoder, an IAF flow, and a style classifier, which forms a more discriminative and expressive latent style representation. Similar to [18], [21], the reference encoder consists of a stack of six 2-D convolutional layers cascaded with one unidirectional 128-unit GRU layer. In each IAF step, we use the structure proposed in [58] as the autoregressive neural network. As shown in Fig. 2, the style classifier contains three FC layers, all with ReLU activation, and a softmax output layer. We use the output of the third FC layer in the style classifier as the style embeddings. For the speaker encoder, we use a 3-layer LSTM with the projection operations followed by a speaker classifier shown in Fig. 4. In this paper, we let the style classifier and the speaker classifiers have the same architecture, to learn discriminative style representations and speaker embeddings that can make both style and speaker to be more distinguishable.

In our model, we adopt Tacotron 2 [1] as the synthesizer, which takes the concatenation of the speaker and style representations as the initial hidden state. As for the discriminator D , we follow the architecture of the discriminator in [59]. The pre-trained discriminator D_s used in the style distortion loss has the same structure as the style encoder followed by a sigmoid output layer. As for the objective function in Eq.(9), we empirically set the coefficients α , β , λ , κ and ω to 1.0, and γ to 5.0, respectively. The six

coefficients ($\alpha/\beta/\gamma/\lambda/\kappa/\omega$) are preset weights for balancing the different loss terms ($\mathcal{L}_{rec}/\mathcal{L}_{adv}/\mathcal{L}_{dis}/\mathcal{L}_{cyc}/\mathcal{L}_{stycls}/\mathcal{L}_{spkcls}$).

In this paper, we use LPCNet, a variant of WaveRNN [60], as a neural vocoder to convert the rendered Mel spectrograms by the synthesizer network into time speech waveforms. The architecture is the same as that described in [45], composed of a sample rate network and a frame rate network. Different from [45], we adopt the Mel spectrogram as the input, rather than the Bark-scale cepstral coefficients [61] and pitch parameters, thus performing a direct connection with the synthesizer network. The training pipeline of our LPCNet is the same as that of [45]. Here, we use the same LPCNet for fair comparisons across all different systems.

C. Evaluation

The performance of seen and unseen speech style transfer is evaluated in speech naturalness, style similarity, and speaker similarity. We conduct Mean Opinion Score (MOS) and preference listening tests (ABX) of speech naturalness to evaluate the reconstruction performance of different experimental systems. An ABX test of style similarity is also conducted to assess the style conversion performance, where subjects are asked to choose which speech sample (A or B) sounds closer to the target style (X) in terms of style. We further conduct a Comparative Mean Opinion Score (CMOS) test of speaker similarity to evaluate how well the transferred speech matches that of the source speaker. For each system, we randomly select the reading style as the seen style and make three experiments: from Reading style to Customer-service style (R2C), Reading style to Poetry style (R2P), and Reading style to Game style (R2G). We randomly choose an unique Taiwanese-reading style from a new female speaker as the unseen style, and conduct three tests from Taiwanese-reading style to Customer-service style (TR2C), Taiwanese-reading style to Poetry style (TR2P), and Taiwanese-reading style to Game style (TR2G) to assess the performance of unseen style transfer.

In addition, style classification accuracy and visualization of style embedding space are adopted to further evaluate the style similarity. Speaker classification accuracy and cosine similarity are calculated to measure the speaker similarity objectively.

V. EXPERIMENTAL RESULTS

A. Speech Naturalness

We use the seen and unseen style evaluation sets, which respectively contain 20 sentences in each style, randomly selected from the test set, to compare the performance of all models in speech naturalness with the MOS and ABX listening tests¹. Table I summarizes the results of MOS from different models, where 15 subjects are requested to carefully listen and evaluate with rating scores from 1 to 5 in 0.5 point increments. It can be seen that our proposed approach outperforms the reference models on both seen and unseen style transfer tasks.

¹Samples can be found at <https://xiaochunan.github.io/disentangling/index.html>

TABLE I
MOS RESULTS WITH 95 % CONFIDENCE INTERVAL OF SEEN AND UNSEEN STYLE TRANSFER FROM DIFFERENT MODELS.

System	Seen style transfer			Unseen style transfer		
	R2C	R2P	R2G	TR2C	TR2P	TR2G
GST	3.35±0.04	3.27±0.02	3.24±0.01	2.92±0.02	2.83±0.03	2.81±0.05
VAE	3.43±0.02	3.39±0.04	3.35±0.02	2.97±0.03	2.92±0.07	2.89±0.06
MRF-IT	3.56±0.05	3.49±0.06	3.45±0.07	3.09±0.05	3.01±0.12	2.97±0.08
MRF-ACC	3.78±0.03	3.70±0.04	3.63±0.06	3.58±0.04	3.49±0.04	3.47±0.05
Proposed	3.99±0.01	3.95±0.02	3.91±0.01	3.86±0.02	3.83±0.03	3.81±0.02

TABLE II
ABX PREFERENCE RESULTS BETWEEN THE PROPOSED MODEL AND EACH REFERENCE SYSTEM FOR SPEECH NATURALNESS ON SEEN AND UNSEEN STYLE TRANSFER.

Preference (%) for seen style transfer									
System A vs System B	R2C			R2P			R2G		
	System A	Neutral	System B	System A	Neutral	System B	System A	Neutral	System B
GST vs Proposed	28.5	31.2	40.3	28.9	31.0	40.1	29.0	31.2	39.8
VAE vs Proposed	29.4	31.0	39.6	30.1	30.1	39.8	30.7	29.9	39.4
MRF-IT vs Proposed	31.2	29.1	39.7	31.5	29.3	39.2	31.9	29.1	39.0
MRF-ACC vs Proposed	33.5	27.5	39.0	33.3	28.1	38.6	33.5	28.1	38.4
Preference (%) for unseen style transfer									
System A vs System B	TR2C			TR2P			TR2G		
	System A	Neutral	System B	System A	Neutral	System B	System A	Neutral	System B
GST vs Proposed	24.9	31.9	43.2	25.2	32.5	42.3	25.8	32.3	41.9
VAE vs Proposed	25.6	33.1	41.3	26.1	33.0	40.9	27.0	31.9	41.1
MRF-IT vs Proposed	29.6	29.7	40.7	29.4	30.2	40.4	29.2	30.8	40.0
MRF-ACC vs Proposed	31.8	29.0	39.2	31.5	29.0	39.5	31.7	28.7	39.6

The performance of the proposed model on unseen style transfer is much better than other models. Specially, most of subjects find that there are more unintelligible parts in an utterance synthesized by using the GST, VAE, MRF-IT and MRF-ACC models. This phenomenon is more obvious on unseen style transfer task. The results show a better generalization of the proposed model on the unseen style transfer. Seen style transfer performs with better speech naturalness than the unseen style transfer. Partially due to the small size of the speech data set in poetry and game styles, hence their MOS scores are slightly lower than that of the customer-service style.

Table II shows the results of ABX tests between the proposed model and all reference systems for speech naturalness on seen and unseen style transfer, where the same 15 subjects are asked to choose the speech samples with higher speech naturalness. We can observe that our proposed framework consistently outperforms four other systems of the prior art, which is also consistent with the MOS results. These observations verify the effectiveness of the proposed approach in terms of speech naturalness.

We also calculate the character error rate (CER) of all models on the same evaluation sets via a Conformer based ASR system [62] to evaluate speech intelligibility objectively. As shown in Table III, we show that our proposed model is superior to the four reference models, i.e., GST, VAE, MRF-IT and MRF-ACC, on both seen and unseen style transfer. This verifies objectively the effectiveness of our proposed approach in terms of speech intelligibility.

B. Style Similarity

We first investigate the style classification accuracy via a speech style classifier, which is independently trained using

TABLE III
CER RESULTS (%) OF DIFFERENT MODELS ON SEEN AND UNSEEN STYLE TRANSFER.

System	Seen style transfer			Unseen style transfer		
	R2C	R2P	R2G	TR2C	TR2P	TR2G
GST	20.5	20.7	20.8	22.8	23.4	23.7
VAE	19.9	20.3	20.5	22.4	22.9	23.1
MRF-IT	19.2	19.6	19.7	21.9	22.1	22.3
MRF-ACC	17.8	18.2	18.6	19.1	19.4	19.5
Proposed	15.4	15.6	15.8	16.2	16.4	16.7

TABLE IV
RESULTS OF STYLE CLASSIFICATION (%) FOR SEEN AND UNSEEN STYLE TRANSFER USING DIFFERENT MODELS.

System	Seen style transfer			Unseen style transfer		
	R2C	R2P	R2G	TR2C	TR2P	TR2G
GST	62.8	62.2	61.5	58.7	57.6	56.2
VAE	64.2	63.8	63.0	59.9	59.2	58.4
MRF-IT	71.3	70.4	69.2	67.7	67.1	66.5
MRF-ACC	76.8	76.3	75.4	73.2	72.5	71.6
Proposed	90.5	90.2	89.9	87.6	86.9	85.8

the seven style samples (i.e., reading, broadcasting, talking, story-telling, customer-service, poetry and game styles) from the training data, to objectively evaluate the style conversion performance. The classifier has the same architecture as the style encoder in our model. Its final validation accuracy is 95.1%. We then synthesize the transferred speech adopting the same evaluation sets, and predict their style labels using the trained classifier. Table IV shows the results of style classification for seen and unseen style transfer using different models. Our proposed model achieves greater accuracy on both seen and unseen style transfer tasks, showing its ability to transfer style in synthesized samples. However, the reference models perform poorly on style classification. Specially, for

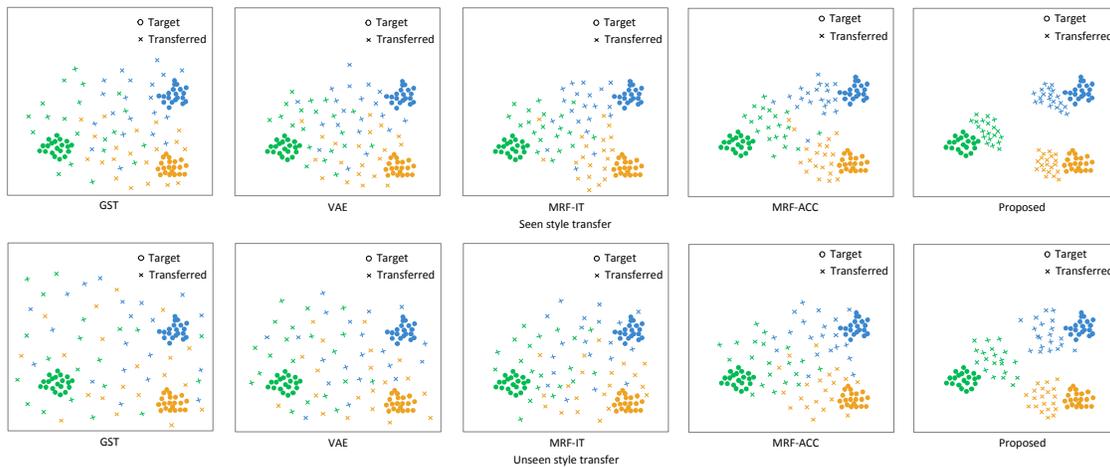


Fig. 5. t-SNE plots of style embeddings for seen and unseen style transfer using different models. Each color corresponds to a different style. Real and transferred utterances appear nearby when they are from the same style, however real and transferred utterances consistently form distinct clusters.

TABLE V
ABX PREFERENCE RESULTS BETWEEN THE PROPOSED MODEL AND EACH REFERENCE SYSTEM FOR STYLE SIMILARITY ON SEEN AND UNSEEN STYLE TRANSFER.

Preference (%) for seen style transfer									
System A vs System B	R2C			R2P			R2G		
	System A	Neutral	System B	System A	Neutral	System B	System A	Neutral	System B
GST vs Proposed	26.8	32.9	40.3	27.2	32.0	40.8	27.6	32.1	40.3
VAE vs Proposed	27.7	31.7	40.6	28.1	31.7	40.2	28.5	31.5	40.0
MRF-IT vs Proposed	28.9	31.2	39.9	29.3	30.6	40.1	29.4	30.7	39.9
MRF-ACC vs Proposed	29.6	31.2	39.2	30.2	30.0	39.8	31.4	29.1	39.5
Preference (%) for unseen style transfer									
System A vs System B	TR2C			TR2P			TR2G		
	System A	Neutral	System B	System A	Neutral	System B	System A	Neutral	System B
GST vs Proposed	24.8	33.2	42.0	25.3	32.9	41.8	25.7	32.4	41.9
VAE vs Proposed	25.7	32.0	42.3	26.1	31.9	42.0	26.3	32.1	41.6
MRF-IT vs Proposed	26.4	32.3	41.3	27.0	31.6	41.4	27.6	31.7	40.7
MRF-ACC vs Proposed	27.5	31.7	40.8	28.2	29.9	41.9	28.7	30.8	40.5

the unseen style transfer (i.e., TR2C, TR2P, and TR2G), many transferred samples from the GST, VAE, and MRF-IT models are difficult to distinguish as their samples after the transfer always have poor style expression, demonstrating a poor unseen style transfer performance. The best reference system, MRF-ACC, still obtains a much lower rate of seen and unseen style transfer than the proposed model. These results validate the effectiveness of the proposed approach for seen and unseen style transfer.

We further visualize the style embedding space using t-SNE [63] to evaluate the style similarity, where different colors correspond to different styles. The style embeddings are respectively extracted from 20 real target speech and 20 transferred utterances for each style. As shown in Fig. 5, our proposed model produces much closer and more separable clusters than the other four reference systems. Particularly, in our model, different styles are well separated from each other in the style embedding space, and transferred utterances tend to lie close to real target speech from the same style in the embedding space. However, the transferred utterances are still easily distinguishable from the real human speech as demonstrated by the t-SNE visualization, where utterances from each transferred style form a nice cluster which is adjacent to a cluster of real target utterances of

the corresponding style. For the GST, VAE, MRF-IT, and MRF-ACC models, real sentences form distinct clusters, but transferred utterances cannot be well separated by style, and appear distant from the corresponding real target style utterances. These observations suggest that our approach can learn effective and discriminative representations in the style space, thus enhancing the performance of both seen and unseen style transfer.

We also conduct an ABX test of style similarity between the proposed model and each reference system to subjectively assess the style conversion performance, where the same 15 listeners are asked to choose the speech samples which sound closer to the target style in terms of style expression. Here, the listeners follow the instructions: “You should not judge the content, grammar, audio quality, or speaker identity of the sentences; instead, just focus on the similarity of the style to one another.” Higher preference means more style similarity is perceived. The results of style similarity on seen and unseen style transfer are shown in Table V, where the listeners give higher preference to the proposed system, showing the proposed approach improves the performance of seen and unseen style transfer. For the unseen style transfer, we find that the GST, VAE and MRF-IT models, in most cases, fail to transfer unseen style of the Taiwanese-reading style to the

target style of customer-service style or poetry style or game style. The MRF-ACC system, the best of all references, is still significantly inferior to the proposed model in the style similarity test. These results demonstrate the effectiveness of our proposed approach for both seen and unseen style transfer.

C. Speaker Similarity

To evaluate how well the transferred speech matches that of the source speaker’s timbre, we respectively conduct CMOS tests between the proposed model and each reference system (i. e., proposed vs GST, proposed vs VAE, proposed vs MRF-IT, and proposed vs MRF-ACC) by using the same evaluation sets. Here, CMOS is used to make a comparison in speaker similarity between two voices, ranging from -3 to 3. Generally, a positive score indicates the reference voice is better than the proposed voice in speaker similarity, worse for a negative score. The same 15 listeners are asked to just focus on the speaker identity of utterances while not judging the content, grammar, audio quality, or style information of the utterances. Table VI reports CMOS results for speaker similarity, where score of the proposed model is fixed to 0 on both seen and unseen style transfer tasks. We can observe that each reference system obtains negative CMOS scores, demonstrating that they perform worse in speaker similarity than the proposed model. This is also evident that our proposed approach delivers better speaker similarity performance than all reference systems, on both seen and unseen style transfer tasks.

We also adopt the cosine similarity to calculate the similarity between the speaker embedding of a transferred sample and the speaker embedding of a randomly selected ground truth utterance from the same speaker to objectively measure the speaker similarity. The results are shown in Table VII, where we can find that the proposed model delivers a higher speaker similarity than the other reference systems, reflecting that our approach has learned a more discriminative speaker representation. Specially, our proposed model obtains a greater cosine similarity, which demonstrates that the transferred utterances via our method are nearly always even closer to the speaker identity of source speaker. On the unseen style transfer of TR2C, TR2P and TR2G, the GST, VAE and MRF-IT models are not capable of keeping the speaker’s timbre, resulting in lower speaker similarity. The best reference system, MRF-ACC, still performs significantly worse than the proposed approach.

We further train a speaker verification model to calculate speaker classification accuracy for the same evaluation sets to objectively assess the speaker similarity. In this experiments, we follow the architecture of speaker verification model in [64], which consists of a small size Thin ResNet-34 [65] with SE-block [66]. Table VIII summarizes the results of speaker classification accuracy on both seen and unseen style transfer tasks. Our proposed approach obtains a higher speaker classification accuracy than the other reference systems. Specially, on the unseen style transfer task (i. e., TR2C, TR2P, and TR2G), the proposed model achieves a much higher accuracy than the GST, VAE, and MRF-IT models, showing its ability to retain speaker identity in the

TABLE VI
CMOS RESULTS FOR SPEAKER SIMILARITY BETWEEN THE PROPOSED MODEL AND EACH REFERENCE SYSTEM.

System	Seen style transfer			Unseen style transfer		
	R2C	R2P	R2G	TR2C	TR2P	TR2G
Proposed	0	0	0	0	0	0
GST	-0.83	-0.85	-0.86	-0.94	-0.95	-0.97
VAE	-0.70	-0.68	-0.67	-0.76	-0.78	-0.80
MRF-IT	-0.45	-0.48	-0.49	-0.52	-0.53	-0.56
MRF-ACC	-0.22	-0.24	-0.25	-0.30	-0.33	-0.34

TABLE VII
COSINE SIMILARITY RESULTS FOR SPEAKER SIMILARITY USING DIFFERENT MODELS.

System	Seen style transfer			Unseen style transfer		
	R2C	R2P	R2G	TR2C	TR2P	TR2G
GST	0.30	0.27	0.26	0.20	0.18	0.17
VAE	0.36	0.34	0.33	0.24	0.23	0.21
MRF-IT	0.44	0.43	0.41	0.35	0.34	0.33
MRF-ACC	0.58	0.56	0.53	0.44	0.42	0.41
Proposed	0.78	0.76	0.75	0.71	0.70	0.68

TABLE VIII
RESULTS OF SPEAKER CLASSIFICATION (%) FOR SEEN AND UNSEEN STYLE TRANSFER USING DIFFERENT MODELS.

System	Seen style transfer			Unseen style transfer		
	R2C	R2P	R2G	TR2C	TR2P	TR2G
GST	62.9	62.4	61.5	59.8	59.1	58.0
VAE	64.7	63.6	63.0	61.9	60.5	60.3
MRF-IT	71.4	70.8	70.1	68.2	67.4	67.1
MRF-ACC	78.2	78.0	77.3	75.6	74.8	73.5
Proposed	92.1	91.6	91.2	90.3	90.0	89.4



Fig. 6. ABX preference results between P-woIAF and proposed model for style similarity on seen and unseen style transfer.

transferred samples. The MRF-ACC system, the best reference model, is still significantly inferior to the proposed method in the speaker classification test. From these results we can conclude that our proposed model can preserve speaker identity and generate speech that resembles the source speaker, but still not as good as the real source speaker.

D. Investigation on Style Transfer without IAF

It is meaningful to investigate the effect of not using IAF in learning style representations for speech style transfer. In this experiments, we remove the IAF structure in the style encoder and test the performance of its seen and unseen style transfer, which is denoted as P-woIAF. Adopting the same evaluation sets, we conduct an ABX test of style similarity between the P-woIAF and the proposed approach to subjectively compare their style conversion performance.

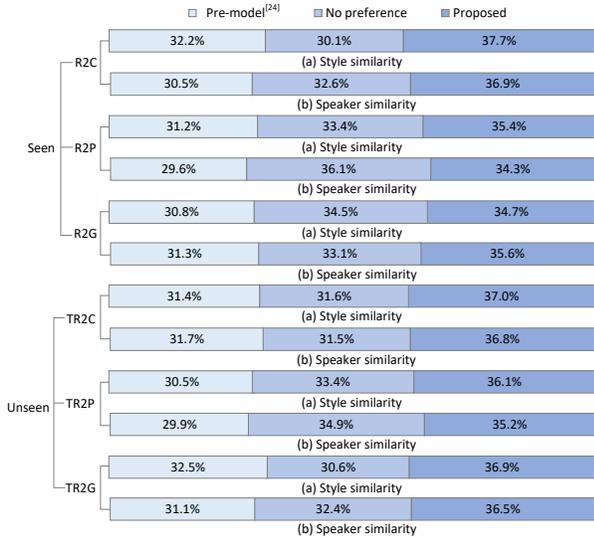


Fig. 7. ABX preference results between Pre-model^[24] and proposed model for style similarity and speaker similarity on seen and unseen style transfer.

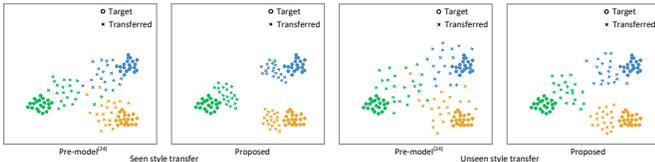


Fig. 8. t-SNE plots of style embeddings for Pre-model^[24] and proposed model on seen and unseen style transfer. Each color corresponds to a different style. Real and transferred utterances appear nearby when they are from the same style, however real and transferred utterances consistently form distinct clusters.

Here the same 15 listeners are also asked to choose the speech samples that sound closer to the target style in terms of style expression. Fig. 6 presents the ABX results for style similarity on both seen and unseen style transfer tasks. We can find that the P-woIAF model obtains lower preference, showing lower style similarity is perceived. However, when we plug the IAF to the style encoder, i.e., our proposed model, the listeners give higher preference to the proposed approach. The performance gain is essentially contributed by the IAF scheme in learning expressive style representations. In summary, we should construct sophisticated distributions with IAF to learn discriminative style representations, for possible improved performance in seen and unseen style transfer.

E. Comparison of Style Transfer without Style and Speaker Classifiers

As mentioned in Section I, our proposed model can be considered as a variant of our previous model [24], which also uses style encoder and speaker encoder but in a different way. Here, we denote the previous model as the Pre-model^[24]. Specially, in the Pre-model^[24], no style classifier and speaker classifier are used in the corresponding style and speaker encoders, hence no constraints of style and speaker classification losses are imposed in the model training process. The Pre-model^[24] is therefore unable to distinguish

TABLE IX
ACCURACY OF SUBJECTIVE STYLE CATEGORY CLASSIFICATION BASED ON THE ABLATION STUDY OF OUR PROPOSED MODEL.

Objective	Seen style transfer			Unseen style transfer		
	R2C	R2P	R2G	TR2C	TR2P	TR2G
\mathcal{L}_{rec}	0.64	0.63	0.61	0.56	0.54	0.53
$+\mathcal{L}_{adv}$	0.70	0.68	0.67	0.63	0.62	0.60
$+\mathcal{L}_{dis}$	0.82	0.80	0.78	0.74	0.73	0.71
$+\mathcal{L}_{cyc}$	0.85	0.83	0.82	0.79	0.75	0.74
$+\mathcal{L}_{stycls}+\mathcal{L}_{spkcls}$	0.94	0.93	0.92	0.88	0.87	0.84

TABLE X
ACCURACY OF SUBJECTIVE SPEAKER CATEGORY CLASSIFICATION BASED ON THE ABLATION STUDY OF OUR PROPOSED MODEL.

Objective	Seen style transfer			Unseen style transfer		
	R2C	R2P	R2G	TR2C	TR2P	TR2G
\mathcal{L}_{rec}	0.62	0.60	0.59	0.57	0.56	0.54
$+\mathcal{L}_{adv}$	0.69	0.68	0.66	0.63	0.61	0.60
$+\mathcal{L}_{dis}$	0.76	0.75	0.73	0.70	0.69	0.67
$+\mathcal{L}_{cyc}$	0.86	0.84	0.83	0.81	0.80	0.78
$+\mathcal{L}_{stycls}+\mathcal{L}_{spkcls}$	0.95	0.93	0.92	0.90	0.89	0.87

the different styles and speakers in a direct manner. In the experiments, we compare the new proposed approach and the Pre-model^[24] with ABX subjective tests of style similarity and speaker similarity. The ABX results are shown in Fig. 7. We show that listeners give higher preference to the new proposed model which has improved the corresponding style and speaker similarities. Similar to Section V-B, we also show the scatter plots of the style embedding space for the Pre-model^[24] and the proposed model. As shown in Fig. 8, our proposed model produces more compact and highly separable clusters than the Pre-model^[24] on both seen and unseen style transfer. All evidence shows that the style and speaker classifiers proposed in this paper contribute to better style and speaker similarities for both seen and unseen style transfer. The results confirm that the style and speaker classifiers are advantageous to make the transfer process more discriminative in styles and speakers.

F. Effect of Different Objectives on Style Transfer

In this paper, we conduct ablation studies on seen and unseen speech style transfer to validate the effectiveness of different objectives. Specifically, we respectively conduct two subjective tests: a style classification test and a speaker classification test, on the same evaluation sets where the same 15 listeners are asked to focus on the style or speaker identity. The style classification accuracy and the speaker classification accuracy are shown in Table IX and Table X, respectively. For each testing sample, we let the listeners select one style from eight style categories (i.e., reading, broadcasting, talking, story-telling, customer-service, poetry, game, and Taiwanese-reading styles), or select one speaker from eight speaker categories (i.e., seven speakers seen in training and a new speaker unseen in training). We can notice that the addition of the adversarial loss (\mathcal{L}_{adv}), the style distortion loss (\mathcal{L}_{dis}), the cycle consistency loss (\mathcal{L}_{cyc}), and the style and speaker classification losses (\mathcal{L}_{stycls} and \mathcal{L}_{spkcls}) can bring substantial

gain for the style classification accuracy and the speaker classification accuracy. The combination of all losses achieves the best accuracy, which outperforms the Pre-model^[24] model (i.e., without the $\mathcal{L}_{stylcls}$ and \mathcal{L}_{spkcls}) by a large margin. In addition, the listeners find that it is much easier to distinguish the category of each style and each speaker after the addition of $\mathcal{L}_{stylcls}$ and \mathcal{L}_{spkcls} .

VI. CONCLUSION

We proposed a novel approach to improve performance of seen and unseen speech style transfer on disjoint, multi-style datasets. More precisely, the introduction of IAF technique has significantly improved the variational inference performance to learn a distinctive and expressive style representation, and the introduction of a well-designed speaker encoder has performed a joint training for learning a discriminative speaker representation. A reconstruction term is used to measure the distortions in both source and target reconstructions, and the adversarial loss for “fooling” a well-trained discriminator. A style distortion loss has made the style representation of a source utterance closer to the target style representation while a cycle consistency loss ensures the transferred utterance to preserve the speaker identity of the source utterance. The addition of a style classifier and a speaker classifier has further enhanced the style and speaker representations, and made the category of each style and each speaker more distinguishable. We have conducted extensive experiments to analyze the IAF structure, the transfer performance without the style and speaker classifiers, and style transfer in different objectives. Experimental results have demonstrated the superior performance of the proposed approach on both seen and unseen style transfer. In the future, we would incorporate our findings to disentangle style and speaker representations, and further improve the speech style transfer performance. We will also apply our seen and unseen style transfer to other TTS applications.

ACKNOWLEDGMENT

This work was supported by the National Key Research and Development Program of China (No. 2020AAA0108600).

REFERENCES

- [1] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Skerry-Ryan, R. A. Saurous, Y. Agiomyrgiannakis, and Y. Wu, “Natural TTS synthesis by conditioning WaveNet on mel spectrogram predictions,” in *Proc. ICASSP*, 2018, pp. 4779–4783.
- [2] Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, Q. Le, Y. Agiomyrgiannakis, R. Clark, and R. A. Saurous, “Tacotron: Towards end-to-end speech synthesis,” in *Proc. INTERSPEECH*, 2017, pp. 4006–4010.
- [3] W. Ping, K. Peng, A. Gibiansky, S. O. Arik, A. Kannan, S. Narang, J. Raiman, and J. Miller, “Deep voice 3: 2000-speaker neural text-to-speech,” in *Proc. ICLR*, 2018, pp. 1–16.
- [4] F. Yang, S. Yang, P. Zhu, P. Yan, and L. Xie, “Improving mandarin end-to-end speech synthesis by self-attention and learnable gaussian bias,” in *Proc. ASRU*, 2019, pp. 208–213.
- [5] G. Sun, Y. Zhang, R. J. Weiss, Y. Cao, H. Zen, A. Rosenberg, B. Ramabhadran, and Y. Wu, “Generating diverse and natural text-to-speech samples using a quantized fine-grained VAE and autoregressive prosody prior,” in *Proc. ICASSP*, 2020, pp. 6699–6703.
- [6] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” in *Proc. NIPS*, 2014, pp. 3104–3112.
- [7] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv preprint arXiv:1409.0473*, 2014.
- [8] P. Wu, Z. Ling, L. Liu, Y. Jiang, H. Wu, and L. Dai, “End-to-end emotional speech synthesis using style tokens and semi-supervised training,” in *Proc. APSIPA ASC*, 2019, pp. 623–627.
- [9] D. Stanton, Y. Wang, and R. Skerry-Ryan, “Predicting expressive speaking style from text in end-to-end speech synthesis,” in *Proc. SLT*, 2018, pp. 595–602.
- [10] R. Liu, B. Sisman, G. Gao, and H. Li, “Expressive TTS training with frame and style reconstruction loss,” *IEEE Transactions on Audio, Speech, and Language Processing*, pp. 1806–1818, 2021.
- [11] Y. Lei, S. Yang, and L. Xie, “Fine-grained emotion strength transfer, control and prediction for emotional speech synthesis,” in *Proc. SLT*, 2021, pp. 423–430.
- [12] K. Akuzawa, Y. Iwasawa, and Y. Matsuo, “Expressive speech synthesis via modeling expressions with variational autoencoder,” in *Proc. INTERSPEECH*, 2018, pp. 3067–3071.
- [13] X. An, Y. Wang, S. Yang, Z. Ma, and L. Xie, “Learning hierarchical representations for expressive speaking style in end-to-end speech synthesis,” in *Proc. ASRU*, 2019, pp. 184–191.
- [14] W.-N. Hsu, Y. Zhang, R. J. Weiss, H. Zen, Y. Wu, Y. Wang, Y. Cao, Y. Jia, Z. Chen, J. Shen, P. Nguyen, and R. Pang, “Hierarchical generative modeling for controllable speech synthesis,” *arXiv preprint arXiv:1810.07217*, 2018.
- [15] R. Habib, S. Mariooryad, M. Shannon, E. Battenberg, R. Skerry-Ryan, D. Stanton, D. T. H. Kao, and T. Bagby, “Semi-supervised generative modeling for controllable speech synthesis,” *arXiv preprint arXiv:1910.01709*, 2019.
- [16] G. Sun, Y. Zhang, R. J. Weiss, Y. Cao, H. Zen, and Y. Wu, “Fully-hierarchical fine-grained prosody modelling for interpretable speech synthesis,” in *Proc. ICASSP*, 2020, pp. 6264–6268.
- [17] S. Pan and L. He, “Cross-speaker style transfer with prosody bottleneck in neural speech synthesis,” in *Proc. INTERSPEECH*, 2021, pp. 4678–4682.
- [18] R. Skerry-Ryan, E. Battenberg, Y. Xiao, Y. Wang, D. Stanton, J. Shor, R. J. Weiss, R. Clark, and R. A. Saurous, “Towards end-to-end prosody transfer for expressive speech synthesis with Tacotron,” in *Proc. ICML*, 2018, pp. 4693–4702.
- [19] T. Li, S. Yang, L. Xue, and L. Xie, “Controllable emotion transfer for end-to-end speech synthesis,” in *Proc. ICSLP*, 2021, pp. 1–5.
- [20] Y. Wang, D. Stanton, Y. Zhang, R. Skerry-Ryan, E. Battenberg, J. Shor, Y. Xiao, F. Ren, Y. Jia, and R. A. Saurous, “Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis,” in *Proc. ICML*, 2018, pp. 5180–5189.
- [21] Y. Zhang, S. Pan, L. He, and Z. Ling, “Learning latent representations for style control and transfer in end-to-end speech synthesis,” in *Proc. ICASSP*, 2019, pp. 6945–6949.
- [22] Y. Bian, C. Chen, Y. Kang, and Z. Pan, “Multi-reference tacotron by intercross training for style disentangling, transfer and control in speech synthesis,” *arXiv preprint arXiv:1904.02373*, 2019.
- [23] M. Whitehill, S. Ma, D. McDuff, and Y. Song, “Multi-reference neural TTS stylization with adversarial cycle consistency,” in *Proc. INTERSPEECH*, 2020, pp. 4442–4446.
- [24] X. An, F. K. Soong, and L. Xie, “Improving performance of seen and unseen speech style transfer in end-to-end neural TTS,” in *Proc. INTERSPEECH*, 2021, pp. 4688–4692.
- [25] D. P. Kingma and M. Welling, “Auto-encoding variational Bayes,” *arXiv preprint arXiv:1312.6114*, 2013.
- [26] A. Atanov, A. Ashukha, K. Struminsky, D. Vetrov, and M. Welling, “The deep weight prior,” *arXiv preprint arXiv:1810.06943*, 2018.
- [27] D. P. Kingma, T. Salimans, R. Jozefowicz, X. Chen, I. Sutskever, and M. Welling, “Improving variational inference with inverse autoregressive flow,” in *Proc. NIPS*, 2016, pp. 4743–4751.
- [28] A. Oord, Y. Li, I. Babuschkin, K. Simonyan, O. Vinyals, K. Kavukcuoglu, G. van den Driessche, E. Lockhart, L. C. Cobo, F. Stimberg, N. Casagrande, D. Grewe, S. Noury, S. Dieleman, E. Elsen, N. Kalchbrenner, H. Zen, A. Graves, H. King, D. B. Tom Walters, and D. Hassabis, “Parallel WaveNet: Fast high-fidelity speech synthesis,” in *Proc. ICML*, 2018, pp. 3915–3923.
- [29] P. Esling, N. Masuda, A. Bardet, R. Despres, and A. Chemla-Romeu-Santos, “Universal audio synthesizer control with normalizing flows,” *arXiv preprint arXiv:1907.00971*, 2019.

- [30] V. Aggarwal, M. Cotesco, N. Prateek, J. Lorenzo-Trueba, and R. Barra-Chicote, "Using VAEs and normalizing flows for one-shot text-to-speech synthesis of expressive speech," in *Proc. ICASSP*, 2020, pp. 6179–6183.
- [31] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, "X-vectors: Robust DNN embeddings for speaker recognition," in *Proc. ICASSP*, 2018, pp. 5329–5333.
- [32] L. Wan, Q. Wang, A. Papir, and I. L. Moreno, "Generalized end-to-end loss for speaker verification," in *Proc. ICASSP*, 2018, pp. 4879–4883.
- [33] Y. Jia, Y. Zhang, R. J. Weiss, Q. Wang, J. Shen, F. Ren, Z. Chen, P. Nguyen, R. Pang, I. L. Moreno, and Y. Wu, "Transfer learning from speaker verification to multispeaker text-to-speech synthesis," in *Proc. NIPS*, 2018, pp. 4485–4495.
- [34] A. Gibiansky, S. Ö. Arik, G. F. Diamos, J. Miller, K. Peng, W. Ping, J. Raiman, and Y. Zhou, "Deep voice 2: Multi-speaker neural text-to-speech," in *Proc. NIPS*, 2017, pp. 2962–2970.
- [35] S. Gururani, K. Gupta, D. Shah, Z. Shakeri, and J. Pinto, "Prosody transfer in neural text to speech using global pitch and loudness features," *arXiv preprint arXiv:1911.09645*, 2019.
- [36] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. NIPS*, 2017, pp. 5998–6008.
- [37] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Proc. NIPS*, 2014, pp. 2672–2680.
- [38] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena, "Self-attention generative adversarial networks," *arXiv preprint arXiv:1805.08318*, 2018.
- [39] S. R. Bowman, L. Vilnis, O. Vinyals, A. M. Dai, R. Jozefowicz, and S. Bengio, "Generating sentences from a continuous space," *arXiv preprint arXiv:1511.06349*, 2015.
- [40] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, " β -VAE: Learning basic visual concepts with a constrained variational framework," in *Proc. ICLR*, 2017, pp. 1–22.
- [41] C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins, and A. Lerchner, "Understanding disentangling in β -VAE," *arXiv preprint arXiv:1804.03599*, 2018.
- [42] W.-N. Hsu, Y. Zhang, and J. Glass, "Learning latent representations for speech generation and transformation," in *Proc. INTERSPEECH*, 2019, pp. 1273–1277.
- [43] H. Daume III and D. Marcu, "Domain adaptation for statistical classifiers," *Journal of artificial Intelligence research*, vol. 26, pp. 101–126, 2006.
- [44] J. Valin and J. Skoglund, "A real-time wideband neural vocoder at 1.6 kb/s using LPCNet," *arXiv preprint arXiv:1903.12087*, 2019.
- [45] J. Valin and J. Skoglund, "LPCNET: Improving neural speech synthesis through linear prediction," in *Proc. ICASSP*, 2019, pp. 5891–5895.
- [46] D. Rezende and S. Mohamed, "Variational inference with normalizing flows," in *Proc. ICML*, 2015, pp. 1530–1538.
- [47] R. v. d. Berg, L. Hasenclever, J. M. Tomczak, and M. Welling, "Sylvester normalizing flows for variational inference," *arXiv preprint arXiv:1803.05649*, 2018.
- [48] J. Villalba, N. Chen, D. Snyder, D. Garcia-Romero, A. McCree, G. Sell, J. Borgstrom, F. Richardson, S. Shon, F. Grondin *et al.*, "State-of-the-art speaker recognition for telephone and video speech: The JHU-MIT submission for NIST SRE18," in *Proc. INTERSPEECH*, 2019, pp. 1488–1492.
- [49] J. S. Chung, A. Nagrani, and A. Zisserman, "Voxceleb2: Deep speaker recognition," *arXiv preprint arXiv:1806.05622*, 2018.
- [50] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz, S. Jan, S. Georg, and V. Karel, "The Kaldi speech recognition toolkit," in *Proc. ASRU*, 2011.
- [51] N. Li, D. Tuo, D. Su, Z. Li, and D. Yu, "Deep discriminative embeddings for duration robust speaker verification," in *Proc. INTERSPEECH*, 2018, pp. 2262–2266.
- [52] O. Vinyals, L. Kaiser, T. Koo, S. Petrov, I. Sutskever, and G. Hinton, "Grammar as a foreign language," in *Proc. NIPS*, 2015, pp. 2773–2781.
- [53] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," in *Proc. NIPS*, 2015, pp. 577–585.
- [54] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, J. Klingner, A. Shah, M. Johnson, X. Liu, L. Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa, A. Rudnick, O. Vinyals, G. Corrado, M. Hughes, and J. Dean, "Google's neural machine translation system: Bridging the gap between human and machine translation," *arXiv preprint arXiv:1609.08144*, 2016.
- [55] M. Mathieu, J. Zhao, P. Sprechmann, A. Ramesh, and Y. LeCun, "Disentangling factors of variation in deep representations using adversarial training," *arXiv preprint arXiv:1611.03383*, 2016.
- [56] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proc. ICCV*, 2017, pp. 2223–2232.
- [57] L. Xue, S. Pan, L. He, L. Xie, and F. K. Soong, "Cycle consistent network for end-to-end style transfer TTS training," *Neural Networks*, vol. 140, pp. 223–236, 2021.
- [58] M. Germain, K. Gregor, I. Murray, and H. Larochelle, "Made: Masked autoencoder for distribution estimation," in *Proc. ICML*, 2015, pp. 881–889.
- [59] H. Guo, F. K. Soong, L. He, and L. Xie, "A new GAN-based end-to-end TTS training algorithm," in *Proc. INTERSPEECH*, 2019, pp. 1288–1292.
- [60] N. Kalchbrenner, E. Elsen, K. Simonyan, S. Noury, N. Casagrande, E. Lockhart, F. Stimberg, A. Oord, S. Dieleman, and K. Kavukcuoglu, "Efficient neural audio synthesis," in *Proc. ICML*, 2018, pp. 2410–2419.
- [61] B. C. Moore, *An introduction to the psychology of hearing*. Brill, 2012.
- [62] P. Guo, F. Boyer, X. Chang, T. Hayashi, Y. Higuchi, H. Inaguma, N. Kamo, C. Li, D. Garcia-Romero, J. Shi, J. Shi, S. Watanabe, K. Wei, and W. Z. Y. Zhang, "Recent developments on ESPnet toolkit boosted by conformer," in *Proc. ICASSP*, 2021, pp. 5874–5878.
- [63] L. Van der Maaten and G. Hinton, "Visualizing data using t-SNE," *Journal of machine learning research*, vol. 9, no. 11, pp. 2579–2605, 2008.
- [64] L. Zhang, Q. Wang, K. A. Lee, L. Xie, and H. Li, "Multi-level transfer learning from near-field to far-field speaker verification," in *Proc. INTERSPEECH*, 2021, pp. 1094–1098.
- [65] H. S. Heo, B.-J. Lee, J. Huh, and J. S. Chung, "Clova baseline system for the voxceleb speaker recognition challenge 2020," *arXiv preprint arXiv:2009.14153*, 2020.
- [66] J. Zhou, T. Jiang, Z. Li, L. Li, and Q. Hong, "Deep speaker embedding extraction with channel-wise feature responses and additive supervision softmax loss function," in *Proc. INTERSPEECH*, 2019, pp. 2883–2887.