


Expressive TTS Training With Frame and Style Reconstruction Loss

Rui Liu , *Member, IEEE*, Berrak Sisman , *Member, IEEE*, Guanglai Gao, and Haizhou Li , *Fellow, IEEE*

Abstract—We propose a novel training strategy for Tacotron-based text-to-speech (TTS) system that improves the speech styling at utterance level. One of the key challenges in prosody modeling is the lack of reference that makes explicit modeling difficult. The proposed technique doesn't require prosody annotations from training data. It doesn't attempt to model prosody explicitly either, but rather encodes the association between input text and its prosody styles using a Tacotron-based TTS framework. This study marks a departure from the style token paradigm where prosody is explicitly modeled by a bank of prosody embeddings. It adopts a combination of two objective functions: 1) frame level reconstruction loss, that is calculated between the synthesized and target spectral features; 2) utterance level style reconstruction loss, that is calculated between the deep style features of synthesized and target speech. The style reconstruction loss is formulated as a perceptual loss to ensure that utterance level speech style is taken into consideration during training. Experiments show that the proposed training strategy achieves remarkable performance and outperforms the state-of-the-art baseline in both naturalness and expressiveness. To our best knowledge, this is the first study to incorporate utterance level perceptual quality as a loss function into Tacotron training for improved expressiveness.

Index Terms—Expressive speech synthesis, tacotron, frame and style reconstruction loss, emotion recognition.

I. INTRODUCTION

WITH the advent of deep learning, neural TTS has shown many advantages over the conventional TTS

techniques [1]–[3]. For example, encoder-decoder architecture with attention mechanism, such as Tacotron [4]–[7], has consistently achieved high voice quality. The key idea is to integrate the conventional TTS pipeline [8], [9] into an unified framework that learns sequence-to-sequence mapping from text to a sequence of acoustic features [7], [10]–[15]. Furthermore, together with a neural vocoder [5], [16]–[21], neural TTS generates natural-sounding and human-like speech which achieves state-of-the-art performance. Despite the progress, the expressiveness of the synthesized speech remains to be improved.

Speech conveys information not only through phonetic content, but also through its prosody. Speech prosody can affect syntactic and semantic interpretation of an utterance [22], [23], that is called linguistic prosody. Speech prosody is also used to display one's emotional state, that is referred to as affective prosody. Both linguistic prosody and affective prosody are manifested over a segment of speech beyond short-time speech frame. Linguistically, speech prosody in general refers to stress, intonation, and rhythm in spoken words, phrases, and sentences. As speech prosody is the result of the interplay of multiple speech properties, it is not easy to define speech prosody by a simple labeling scheme [24]–[28]. Even if a labeling scheme is possible [29], [30], a set of discrete labels may not be sufficient to describe the entire continuum of speech prosody.

Besides naturalness, one of the factors that differentiate human speech from today's synthesized speech is their expressiveness. Prosody is one of the defining features of expressiveness that makes speech lively. Several recent studies successfully improve the expressiveness of Tacotron TTS framework [31]–[35]. The idea is to learn latent prosody embedding, i.e. style token, from training data [31], [36], [37]. At run-time, the style token can be used to predict the speech style from text [32], or to transfer the speech style from a reference utterance to target [33]. It is observed that such speech styling is effective and consistently improves speech quality. Sun *et al.* [34], [35] further study a hierarchical, fine-grained and interpretable latent variable model for prosody rendering. The studies show that precise control of the prosody style leads to improvement of prosody expressiveness in the Tacotron TTS framework. However, several issues have hindered the effectiveness of above prosody modeling techniques.

First, the latent embedding space of prosody is learnt in an unsupervised manner, where the style is defined as anything but speaker identity and phonetic content in speech. We note that many different styles co-exist in speech. Some are speaker dependent, such as accent and idiolect, others are speaker

Manuscript received July 19, 2020; revised February 7, 2021 and April 12, 2021; accepted April 24, 2021. Date of publication April 30, 2021; date of current version June 1, 2021. This work was supported in part by SUTD Start-up Grant Artificial Intelligence for Human Voice Conversion (SRG ISTD 2020 158) and SUTD AI Grant, titled 'The Understanding and Synthesis of Expressive Speech by AI'. The work of Haizhou Li is supported by the National Research Foundation, Singapore under its AI Singapore Programme Award AISG-GC-2019-002 and Award AISG-100E-2018-006, and its National Robotics Programme under Grant 1922500054, and in part by RIE2020 Advanced Manufacturing and Engineering Programmatic Grants A1687b0033 and A18A2b0046. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Lei Xie.

Rui Liu is with the Singapore University of Technology and Design (SUTD) and National University of Singapore, Singapore 117583, Singapore (e-mail: liurui_jmu@163.com).

Berrak Sisman is with the Singapore University of Technology and Design (SUTD), Singapore 117583, Singapore (e-mail: berraksisman@u.nus.edu).

Guanglai Gao is with the Department of Computer Science, Inner Mongolia University, Hohhot 010021, China (e-mail: csggl@imu.edu.cn).

Haizhou Li is with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore 119077, Singapore and also with Faculty 3 Computer Science/Mathematics, Enrique-Schmidt-Str. 5 Cartesium, University of Bremen, 28359 Bremen, Germany (e-mail: haizhou.li@nus.edu.sg).

Digital Object Identifier 10.1109/TASLP.2021.3076369

independent such as prosodic phrasing, lexical stress and prosodic stress. There is no guarantee that such latent embedding space of style represents only the intended prosody. Second, while the techniques don't require the prosody annotations on training data, they require a reference speech or a manual selection of style token [31] in order to explicitly control the style of output speech during run-time inference. While it is possible to automate the style token selection [32], a correct prediction of style token is subject to both the design of the style token dictionary, and the run-time style token prediction algorithm. Third, the style token dictionary in Tacotron is trained from a collection of speech utterances to represent a large range of acoustic expressiveness for a speaker or an audiobook [31]. It is not intended to provide differential prosodic details at phrase or utterance level. It is desirable for Tacotron system to learn to automate the prosody styling in response to input text at run-time, that will be the focus of this paper.

To address the above issues, we believe that Tacotron training should minimize frame level reconstruction loss [4], [5] and utterance level perceptual loss at the same time. Perceptual loss is first proposed for image stylization and synthesis [37]–[40], where feature activation patterns, or deep features, derived from pre-trained auxiliary networks are used to optimize the perceptual quality of output image. Several computational models have been proposed to approximate human perception of audio quality, such as Perceptual Evaluation of Audio Quality (PEAQ) [41], Perceptual Evaluation of Speech Quality (PESQ) [42], and Perceptual Evaluation of Audio methods for Source Separation (PEASS) [43]. However, such models are not differentiable, hence cannot be directly employed during TTS training. We believe that utterance level perceptual loss based on deep features that reflects global speech style would be useful to improve overall speech quality.

We are motivated to study a novel training strategy for TTS systems, that learns to associate prosody styles with input text implicitly. We would like to avoid the use of prosody annotations. We don't attempt to model prosody explicitly either, but rather learn the association between prosody styles and input text using existing neural TTS system, such as Tacotron. As the training strategy is only involved during training, it doesn't change the run-time inference process for neural TTS system. At run-time, we don't require any reference signal nor manual selection of prosody style.

The main contributions of this paper include: 1) we propose a novel training strategy for Tacotron TTS that improves utterance level expressiveness of speech; 2) we propose to supervise the training of Tacotron with a fully differentiable perceptual loss, which is derived from a pre-trained auxiliary network, in addition to frame reconstruction loss; and 3) we successfully implement a system that doesn't require any reference speech nor manual selection of prosody style at run-time. To our best knowledge, this is the first study to incorporate perceptual loss into Tacotron training for improved expressiveness.

This paper is organized as follows: In Section II, we present the research background and related work to motivate our study. In Section III, we propose a novel training strategy for TTS system with frame and style reconstruction loss. In Section IV,

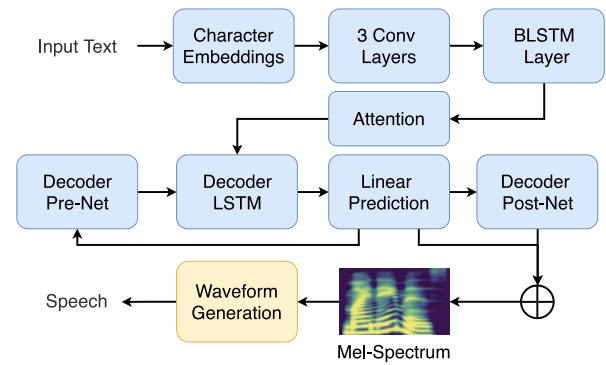


Fig. 1. Block diagram of Tacotron2-based TTS reference baseline [5].

we report the subjective and objective evaluations. Section V concludes the discussion.

II. BACKGROUND AND RELATED WORK

This work is built on several previous studies on neural TTS, prosody modeling, perceptual loss, and speech emotion recognition. Here we briefly summarize the related previous work to set the stage for our study, and to place our novel contributions in a proper context.

A. Tacotron2-Based TTS

In this paper, we adopt the Tacotron2-based [5] TTS model as a reference baseline, which is also referred to as *Tacotron* baseline for brevity.

The overall architecture of the reference baseline includes encoder, attention-based decoder and waveform generation module [44]–[46] as illustrated in Fig. 1. The encoder consists of two components, a convolutional neural network (CNN) module [47], [48] that has 3 convolutional layers, and a bidirectional LSTM (BLSTM) [49] layer. The decoder consists of four components: a 2-layer pre-net, 2 LSTM layers, a linear projection layer and a 5-convolution-layer post-net. The decoder is a standard autoregressive recurrent neural network that generates mel-spectrum features and stop tokens frame by frame. There are two common techniques to generate the audio waveform from mel-spectrum features. One is the Griffin Lim [44] algorithm, another is via a neural vocoder [5], [45], [46], [50].

Just like other TTS systems, Tacotron [4], [5] TTS system predicts mel-spectrum features from input sequence of characters by minimizing a frame level reconstruction loss. Such frame level objective function focuses on the distance between spectral features. It does not seek to optimize the perceptual quality at utterance level. To improve the suprasegmental expressiveness, there have been studies [32], [35], [51] on latent prosody representations, that make possible prosody styling in Tacotron TTS framework. However, most of the studies rely on the style tokens mechanism to explicitly model the prosody. Simply speaking, they build a Tacotron TTS system that synthesizes speech, and learns the global style tokens (GST) at the same time. At run-time inference, they apply the style tokens to control the expressive effect [31], [33], that is referred to as the GST-Tacotron paradigm.

In this paper, we advocate a new way of addressing the expressiveness issue by integrating a perceptual quality motivated objective function into the training process, in addition to the frame level reconstruction loss function. We no longer require any dedicated prosody control mechanism during run-time inference, such as style tokens in Tacotron system.

B. Prosody Modeling in TTS

Prosody conveys linguistic, para-linguistic and various types of non-linguistic information, such as speaker identity, intention, attitude and mood [52], [53]. It is inherently supra-segmental [1], [54] due to the fact that prosody patterns cannot be derived solely from short-time segments [55]. Prosody is hierarchical in nature [55]–[58] and affected by long-term dependencies at different levels such as word, phrase and utterance level [59]. Studies on hierarchical modeling of F0 in speech synthesis [1], [60], [61] suggest that utterance-level prosody modeling is more effective. Similar studies, such as continuous wavelet transform, can be found in many speech synthesis related applications [59], [62]–[65]. In this paper, we will study a novel technique to observe utterance-level prosody quality during Tacotron training to achieve expressive synthesis.

The early studies of modeling speaking styles are carried out on Hidden Markov Models (HMM) [9], [66], where we can synthesize speech with an intermediate speaking style between two speakers through model interpolation [67]. To improve the HMM-based TTS model, there have been studies to incorporate unsupervised expression cluster information during training [68]. Deep learning opens up many possibilities for expressive speech synthesis, where speaker, gender, and age codes can be used as control vectors to change TTS output in different ways [69]. The style tokens, or prosody embeddings, represent one type of such control vectors, that is derived from a representation learning network. The success of prosody embedding motivates us to further develop the idea.

Tacotron TTS framework has achieved remarkable performance in terms of spectral feature generation. With a large training corpus, it may be able to generate natural prosody and expression by remembering the training data using a large number of network parameters. However, its training process doesn't aim to optimize the system for expressive prosody rendering. As a result, Tacotron TTS system tends to generate speech outputs that represent model average, rather than the intended prosody.

The idea of global style tokens [31], [32] represents a success in controlling prosody style of Tacotron output. Style tokens learn to represent high level styles, such as speaker style, pitch range, and speaking rate across a collection of utterances or a speech database. We argue that they neither necessarily represent the useful styles to describe the continuum of prosodic expressions [70], nor provide the dynamic and differential prosodic details with the right level of granularity at utterance level. Sun *et al.* [34], [35] study a way to include a hierarchical, fine-grained prosody representation, that represents the recent attempts to address the problems in GST-Tacotron paradigm.

We would like to address three issues in the existing prosody modeling in Tacotron framework, 1) lack of prosodic supervision

during training; 2) limitation of explicit prosody modeling, such as style tokens, in describing the continuum of prosodic expressions; 3) lack of dynamic and differential prosody at utterance level.

C. Perceptual Loss for Style Reconstruction

It is noted that frame-level reconstruction loss, denoted as *frame reconstruction loss* in short, is not always consistent with human perception because it doesn't take into account human sensitivities to temporal and spectral information, such as prosody and temporal structure of the utterance. For example, if one repeatedly asks the same question two times, despite the perceptual similarity of two utterances, they would be very different as measured by frame-level losses.

Perceptual loss refers to the training loss derived from a pre-trained auxiliary network [38]. The auxiliary network is usually trained on a different task that provides perceptual quality evaluation of an input at a higher level than a speech frame. The intermediate feature representations, generated by the auxiliary network in form of hidden layer activations, are usually referred to as deep features. They are used as the high level abstraction to measure the training loss between reconstructed signals and reference signals. Such training loss is also called deep feature loss [71], [72].

In speech enhancement, perceptual loss has been used successfully in end-to-end speech denoising pipeline, with an auxiliary network pre-trained on audio classification task [73]. Kataria *et al.* [71] propose to use perceptual loss which optimizes the enhancement network with an auxiliary network pre-trained on speaker recognition task. In voice conversion, Lo *et al.* [74] propose deep learning-based assessment models to predict human ratings of converted speech. Lee [75] propose a perceptually meaningful criterion where human auditory system was taken into consideration in measuring the distances between the converted speech and the reference.

In speech synthesis, Oord *et al.* propose to train a WaveNet-like classifier with perceptual loss for phone recognition [76]. As the classifier extracts high-level features that are relevant for phone recognition, this loss term supervises the training of WaveNet to look after temporal dynamics, and penalize bad pronunciations. Cai *et al.* [77] study to use a pre-trained speaker embedding network to provide feedback constraint, that serves as the perceptual loss for the training of a multi-speaker TTS system.

In the context of prosody modeling, the perceptual loss in the above studies can be generally described as *style reconstruction loss* [38]. Following the same principle, we would like to propose a novel auxiliary network, that is pre-trained on a speech emotion recognition (SER) task, to extract high level prosody representations. By comparing prosody representations in a continuous space, we measure perceptual loss between two utterances. While perceptual loss is not new in speech reconstruction, the idea of using a pre-trained emotion recognition network for perceptual loss is a novel attempt in speech synthesis.

D. Deep Features for Perceptual Loss

Now the question is which deep features could be suitable for measuring perceptual loss. We benefit from the prior work in prosody modeling. Prosody embedding in Tacotron is a type of feature learning, that learns the representation for prediction or classification tasks. With deep learning algorithms, automatic feature learning can be achieved in either supervised, such as multilayer perceptron [78], or unsupervised manner, such as variational autoencoder [79]. Deep features are usually more generalizable, and easier to manage than hand-crafted or manually designed features [80]. There have been studies on representation learning for prosody patterns, such as speech emotion [81], and speech styles [31].

Affective prosody refers to the expression of emotion in speech [82], [83]. It is prominently exhibited in emotion speech database. Therefore, the studies in speech emotion recognition provide valuable insights into prosodic modeling. Emotion are usually characterized by discrete categories, such as happy, angry, and sad, and continuous attributes, such as activation, valence and dominance [84], [85]. Recent studies show that latent representations of deep neural networks also characterize well emotion in a continuous space [78].

There have been studies to leverage emotion speech modeling for expressive TTS [33], [68], [86]–[88]. Eyben *et al.* [68] incorporate unsupervised expression cluster information into an HMM-based TTS system. Skerry-Ryan *et al.* [33] study learning prosody representation from animated and emotive storytelling speech corpus. Wu *et al.* [86] propose a semi-supervised training of Tacotron TTS framework for emotional speech synthesis, where style tokens are defined to represent emotion categories. Gao *et al.* [87] propose to use an emotion recognizer to extract the style embedding for speech style transfer. Um *et al.* [88] study a technique to apply style embedding to Tacotron system to generate emotional speech, and to control the intensity of emotional expressiveness.

All the studies point to the fact that emotion-related deep features serve as the excellent descriptors of speech prosody and speech styles. In this paper, instead of using the style tokens to control the TTS outputs, we would like to study how to use deep style features to measure perceptual loss for the training of neural TTS system in general.

III. TACOTRON WITH FRAME AND STYLE RECONSTRUCTION LOSS

We propose a novel training strategy for Tacotron with both frame and style reconstruction loss. As the style reconstruction loss is formulated as a perceptual loss (PL) [38], the proposed frame and style training strategy is called *Tacotron-PL* in short. It seeks to optimize both frame-level spectral loss, that is *frame reconstruction loss*, as well as utterance-level style loss, that is *style reconstruction loss*, at the same time.

The overall framework is illustrated in Fig. 2, that has three stages: 1) training of style descriptor, 2) the proposed frame and style training for *Tacotron-PL* model, and 3) run-time inference. In Stage I, we train an auxiliary network to serve as the style descriptor for input speech utterances. In Stage II, the proposed

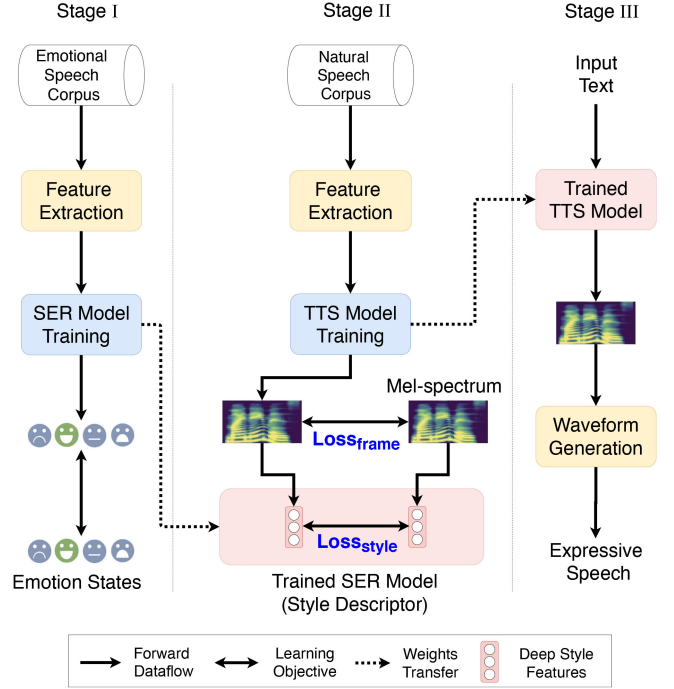


Fig. 2. Overall framework of a *Tacotron-PL* system in three stages: Stage I for training of style descriptor; Stage II for training of *Tacotron-PL*; Stage III for run-time inference.

frame and style training strategy is implemented to associate input text with acoustic features, as well as prosody style of natural speech, that is assisted by the style descriptor obtained from Stage I. In Stage III, the *Tacotron-PL* system takes input text and generates expressive speech in the same way as a standard Tacotron does. Unlike other Tacotron variants [31], *Tacotron-PL* doesn't require any add-on module or process for run-time inference.

As discussed in Section II-A, traditional Tacotron architecture contains a text encoder and an attention-based decoder. We first encode input character embedding into hidden state, from which the decoder generates mel-spectrum features. During training, we adopt a frame-level mel-spectrum loss as in [5], which is a L_2 loss between the synthesized mel-spectrum $\hat{\mathbf{Y}} = \{\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_t, \dots, \hat{\mathbf{y}}_T\}$ and target mel-spectrum $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_t, \dots, \mathbf{y}_T\}$. We have $Loss_{frame}$ as follows,

$$Loss_{frame}(\mathbf{Y}, \hat{\mathbf{Y}}) = \sum_{t=1}^T L_2(\mathbf{y}_t, \hat{\mathbf{y}}_t) \quad (1)$$

which is designed to minimize frame level distortion. As it doesn't guarantee utterance level similarity concerning speech expressions, such as speech prosody and speech styles. We will study a new loss function $Loss_{style}$ next, that measures the utterance-level style reconstruction loss.

A. Stage I: Training of Style Descriptor

One of the great difficulties of prosody modeling is the lack of reference samples. In linguistics, we usually describe prosody styles qualitatively. However, precise annotation of

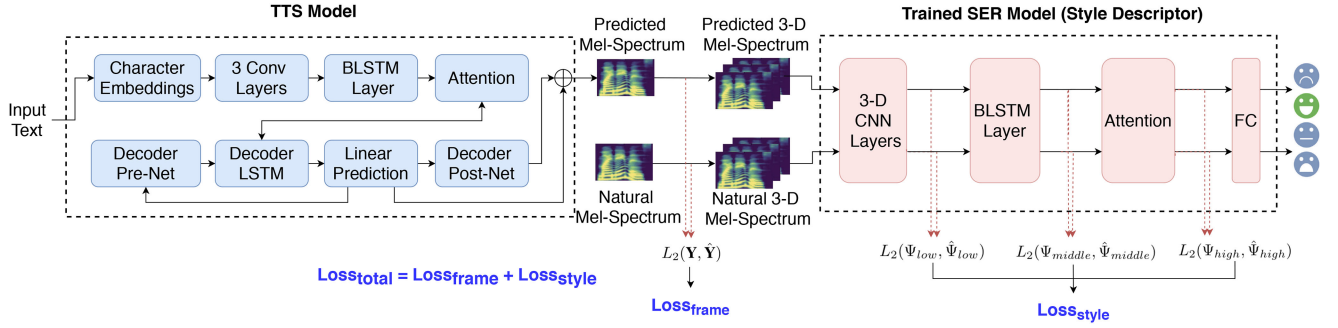


Fig. 3. Block diagram of the proposed training strategy, *Tacotron-PL*. A speech emotion recognition (SER) model is trained separately to serve as an auxiliary model to extract deep style features. A *style reconstruction loss*, $Loss_{style}$, is computed between the deep style features of the generated and reference speech at utterance-level.

speech prosody is not straightforward. One of the ways to describe a prosody style is to show by example. The idea of style token [31] shows a way to compare two prosody styles quantitatively using deep features.

Manual prosodic annotations of recorded speech [29] provide quantifiable prosodic labels that allow us to associate speech styles with actual acoustic features. Prosody labeling schemes often attempt to describe prosodic phenomena, such as the supra-segmental features of intonation, stress, rhythm and speech rate, in discrete categories. Categorical labels of speech emotion [89] also seek to achieve a similar goal. The prosody labeling schemes serve as a type of style descriptor. With deep neural network, one is able to learn the feature representation of the data at different level of abstraction in a continuous space [90]. As speech styles naturally spread over a continuum rather than forced-fitting into a finite set of categorical labels, we believe that deep neural network learned from animated and emotive speech serves as a more suitable style descriptor.

We propose to use a speech emotion recognizer (SER) [82], [83] as a style descriptor $F(\cdot)$, which extracts deep style features Ψ from an utterance \mathbf{Y} , or $\Psi = F(\mathbf{Y})$. We use neuronal activations of hidden units in a deep neural network as the deep style features to represent high level prosodic abstraction at utterance level. In practice, we first train an SER network with highly animated and emotive speech with supervised learning. We then derive deep style features from a small intermediate layer. As the intermediate layer is small relative to the size of the other layers, it creates a constriction in the network that forces the information pertinent to emotion classification into a low dimensional prosody representation [91]. Such low dimensional prosody representation is expected to describe the prosody style of speech signals as the SER network relies on the prosody representation for accurate emotion classification.

We follow the SER implementation in [36], [92] as illustrated in Fig. 3, that forms part of Fig. 2. The SER network includes 1) a three-dimensional (3-D) CNN layer; 2) a BLSTM layer [93]; 3) an attention layer; and 4) a fully connected (FC) layer. The 3-D CNN [92] first extracts a latent representation from mel-spectrum, its delta and delta-delta values from input utterance, converting the input utterance of variable length into a fixed size latent representation, denoted as deep features sequence Ψ_{low} , that reflects the semantics of emotion. The BLSTM summarizes

TABLE I
THE MCD, RMSE AND FD RESULTS OF DIFFERENT SYSTEMS

| System | MCD [dB] | RMSE [Hz] | FD [frame] |
|------------------|-------------|-------------|--------------|
| Tacotron | 7.01 | 1.53 | 15.59 |
| Tacotron-PL(L) | 6.37 | 0.94 | 13.96 |
| Tacotron-PL(M) | 6.70 | 1.21 | 14.20 |
| Tacotron-PL(H) | 6.88 | 1.42 | 15.41 |
| Tacotron-PL(LMH) | 6.62 | 1.16 | 14.15 |

TABLE II
THE AB PREFERENCE TEST FOR EXPRESSIVENESS AND NATURALNESS
EVALUATION BY 15 LISTENERS, WITH 95% CONFIDENCE INTERVALS
COMPUTED FROM THE T-TEST

| Contrastive pair | Preference(%) | | | p-value |
|-------------------------------------|---------------|---------|--------|---------|
| | Former | Neutral | Latter | |
| Expressiveness | | | | |
| Tacotron vs. Tacotron-PL(L) | 32.44 | 13.33 | 54.23 | 0.00119 |
| Tacotron-PL(LMH) vs. Tacotron-PL(L) | 37.78 | 11.56 | 50.66 | 0.00124 |
| Naturalness | | | | |
| Tacotron vs. Tacotron-PL(L) | 36.44 | 18.22 | 45.34 | 0.00101 |
| Tacotron-PL(LMH) vs. Tacotron-PL(L) | 39.11 | 15.56 | 45.33 | 0.00096 |

TABLE III
BEST WORST SCALING (BWS) LISTENING EXPERIMENTS THAT COMPARE
FOUR DEEP STYLE FEATURES IN FOUR *TACOTRON-PL* MODELS

| System | Best (%) | Worst (%) |
|------------------|-----------|-----------|
| Tacotron-PL(L) | 80 | 5 |
| Tacotron-PL(M) | 8 | 26 |
| Tacotron-PL(H) | 2 | 48 |
| Tacotron-PL(LMH) | 10 | 21 |

the temporal information of Ψ_{low} into another latent representation Ψ_{middle} . Finally, the attention layer assigns weights to Ψ_{middle} and generates Ψ_{high} for emotion prediction.

The question is which of the latent representations, Ψ_{low} , Ψ_{middle} , and Ψ_{high} , is suitable to be the deep style features. To validate the descriptiveness of deep style features, we perform an analysis on LJ-Speech corpus [94]. Specifically, we randomly select five utterances from each of the six style groups from the database, each group having a distinctive speech style, namely, 1) Short question; 2) Long question; 3) Short answer; 4) Short statement; 5) Long statement and 6) Digit string. The complete list of utterances can be found at Table V in Appendix A.

TABLE IV
THE AB PREFERENCE TEST FOR EXPRESSIVENESS AND NATURALNESS
EVALUATION BY 15 LISTENERS, WITH 95% CONFIDENCE INTERVALS
COMPUTED FROM THE T-TEST

| Contrastive pair | Preference(%) | | | <i>p</i> -value |
|-----------------------------|---------------|---------|--------|-----------------|
| | Former | Neutral | Latter | |
| Expressiveness | | | | |
| Tacotron vs. Tacotron-ST | 30.22 | 20.44 | 49.34 | 0.00135 |
| Tacotron-ST vs. Tacotron-PL | 28.89 | 15.11 | 56.00 | 0.00129 |
| Naturalness | | | | |
| Tacotron vs. Tacotron-ST | 31.11 | 22.22 | 46.67 | 0.00133 |
| Tacotron-ST vs. Tacotron-PL | 30.67 | 20.89 | 48.44 | 0.00108 |

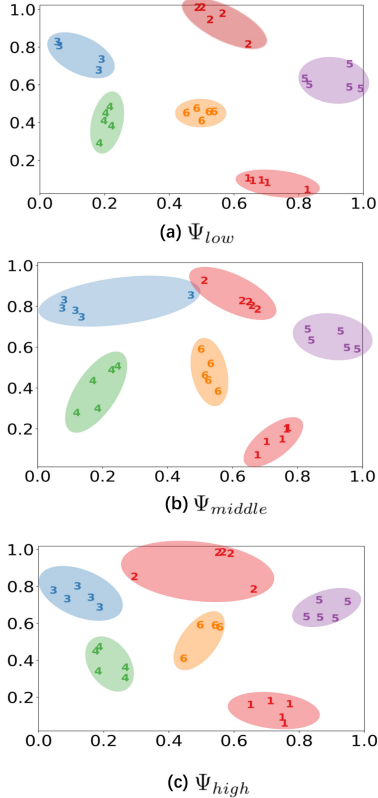


Fig. 4. t-SNE plot of the distributions of deep style features Ψ_{low} , Ψ_{middle} and Ψ_{high} for six groups of utterances in LJ-Speech corpus. The list of utterances can be found at Table V in Appendix A.

We visualize the Ψ_{low} , Ψ_{middle} and Ψ_{high} of utterances using the t-SNE algorithm in a two dimensional plane [95], as shown in Fig. 4. Please note that the distributions of digits 1 to 6 represent those of groups 1 to 6 in the two dimensional space. As illustrated in Table V, the utterances within the same group form a cluster, while the utterances between groups distance from one another. To visualize, we color the clusters to highlight their distributions. It is observed that Ψ_{low} , Ψ_{middle} and Ψ_{high} of utterances form clear style groups in terms of feature distributions, that correspond to the six different utterance styles summarized in Table V. Furthermore, it is clear that Fig. 4(a) shows a better clustering than Fig. 4(b) and Fig. 4(c). We will further compare the performance of different deep style features through TTS experiments in Section IV.

B. Stage II: Tacotron-PL Training

During the training of *Tacotron-PL*, the SER-based style descriptor $F(\cdot)$ is used to extract the deep style features Ψ . We define a style reconstruction loss that compares the prosody style between the reference speech \mathbf{Y} and the generated speech $\hat{\mathbf{Y}}$.

$$Loss_{style}(\mathbf{Y}, \hat{\mathbf{Y}}) = L_2(\Psi, \hat{\Psi}) \quad (2)$$

where $\Psi = F(\mathbf{Y})$ and $\hat{\Psi} = F(\hat{\mathbf{Y}})$. As illustrated in Fig. 3, the proposed training strategy involves two loss functions: 1) $Loss_{frame}$ that minimizes the loss between synthesized and original mel-spectrum at frame level; and 2) $Loss_{style}$ that minimizes the style differences between the synthesized and reference speeches at utterance level.

$$Loss_{total}(\mathbf{Y}, \hat{\mathbf{Y}}) = Loss_{frame}(\mathbf{Y}, \hat{\mathbf{Y}}) + Loss_{style}(\mathbf{Y}, \hat{\mathbf{Y}}) \quad (3)$$

where $Loss_{frame}$ is also the loss function of a traditional Tacotron [5] system.

Style reconstruction loss can be seen as perceptual quality feedback at utterance level to supervise the training of prosody style. All parameters in the TTS model are updated with the gradients of the total loss through back-propagation. We expect that mel-spectrum generation will learn from local and global viewpoint through the frame and style reconstruction loss.

C. Stage III: Run-Time Inference

The inference stage follows exactly the same Tacotron workflow, that only involves the TTS Model in Fig. 3. The difference between *Tacotron-PL* and other global style tokens variation of Tacotron is that *Tacotron-PL* encodes prosody styling inside the standard Tacotron architecture. It doesn't require any add-on module.

At run-time, the Tacotron architecture takes text as input and generate expressive mel-spectrum features as output, that is followed by Griffin-Lim algorithm [44] and WaveRNN vocoder [45] in this paper to generates audio signals.

IV. EXPERIMENTS

We train a SER as the style descriptor on IEMOCAP dataset [89], which consists of five sessions. The dataset contains a total of 10039 utterances, with an average duration of 4.5 seconds at a sampling rate of 16 kHz. We only use a subset of the improvised data with four emotional categories, namely, happy, angry, sad, and neutral, which are recorded in the hypothetical scenarios designed to elicit specific types of emotions.

With the style descriptor, we further train a Tacotron system on LJ-Speech database [94], which consists of 13100 short clips with a total of nearly 24 hours of speech from one single speaker reading 7 non-fiction books. The speech samples are available from the demo link.¹

¹**Speech Samples:** <https://ttslr.github.io/Expressive-TTS-Training-with-Frame-and-Style-Reconstruction-Loss/>

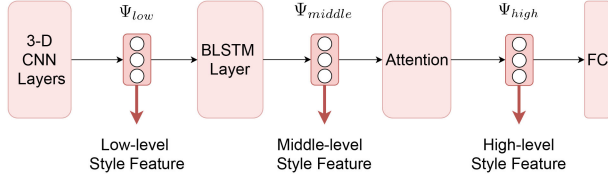


Fig. 5. Three level (low, middle and high) of deep style features extracted from SER-based style descriptors for computing style construction loss.

A. Comparative Study

We develop five Tacotron-based TTS systems for a comparative study, that includes the Tacotron baseline, and four variants of Tacotron with the proposed training strategy, *Tacotron-PL*.

To study the effect of different style descriptors, we compare the use of four deep style features, which includes three single features and a combination of them, in $Loss_{style}$, as illustrated in Fig. 5, and summarized as follows:

- *Tacotron*: Tacotron [5] trained with $Loss_{frame}$ as in Eq. (1), that doesn't explicitly model speech style.
- *Tacotron-PL(L)*: *Tacotron-PL* which uses Ψ_{low} in $Loss_{style}$.
- *Tacotron-PL(M)*: *Tacotron-PL* which uses Ψ_{middle} in $Loss_{style}$.
- *Tacotron-PL(H)*: *Tacotron-PL* which uses Ψ_{high} in $Loss_{style}$.
- *Tacotron-PL(LMH)*: *Tacotron-PL* which uses $\{\Psi_{low}, \Psi_{middle}, \Psi_{high}\}$ in $Loss_{style}$.

B. Experimental Setup

For SER training, we first split the speech signals into segments of 3 seconds as in [92]. We then extract 40-channel mel-spectrum features with a frame size of 50 ms and 12.5 ms frame shift. The first convolution layer has 128 feature maps, while the remaining convolution layers have 256 feature maps. The filter size for all convolution layers is 5×3 , with 5 along the time axis, and 3 along the frequency axis, and the pooling size for the max pooling layer is 2×2 . We add a linear layer with 200 output units after 3-D CNN for dimension reduction.

In this way, the 3-D CNN extracts a fixed size of latent representation with 150×200 dimension from the input utterance, that we use as the deep style features $\Psi_{low} = F_{low}(\cdot)$ to represent a temporal sequence of 150 segment, each having an embedding of 200 elements. As each direction of BLSTM layer contains 128 cells, in two directions, we obtain 256 output activations for each input segment, that are further mapped to 200 output units via a linear layer. BLSTM summarizes the temporal information of Ψ_{low} into another fixed size latent representation $\Psi_{middle} = F_{middle}(\cdot)$ of 150×200 dimension. The attention layer assigns the weights to Ψ_{middle} and generate a new latent representation $\Psi_{high} = F_{high}(\cdot)$. All latent representation $\Psi_{low}, \Psi_{middle}, \Psi_{high}$ have the same dimension.

The fully connected layer contains 64 output units. Batch normalization [96] is applied to the fully connected layer to accelerate training and improve the generalization performance. The parameters of the SER model were optimized by minimizing

the cross-entropy objective function, with a minibatch of 40 samples, using the Adam optimizer with Nestorov momentum. The initial learning rate is set to 10^{-4} and the momentum is set to 0.9. In this way, we obtain a SER style descriptor that is reported with an average classification accuracy of 73.2% for all emotions on the test set.

The SER-based style descriptor is used to extract deep style features for the computing of $Loss_{style}$. For TTS training, the encoder takes a 256-dimensions character sequence as input and the decoder generates the 40-channel mel-spectrum. The training utterances from LJ-Speech database are of variable length. Mel-spectrum features are also extracted with a frame size of 50 ms and 12.5 ms frame shift. They are normalized to zero-mean and unit-variance to serve as the reference target. The decoder predicts only one non-overlapping output frame at each decoding step. We use the Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$ and a learning rate of 10^{-3} exponentially decaying to 10^{-5} starting at 50 k iterations. We also apply L_2 regularization with weight 10^{-6} . All models are trained with a batch size of 32 and 150 k steps.

C. Objective Evaluation

We conduct objective evaluation experiments to compare the systems in a comparative study. The results are summarized in Table I.

1) *Performance Evaluation Metrics*: Mel-cepstral distortion (MCD) [97] is used to measure the spectral distance between the synthesized and reference mel-spectrum features that is known to correlate well with human perception [97]. MCD is calculated as:

$$MCD = \frac{10\sqrt{2}}{\ln 10} \frac{1}{N} \sqrt{\sum_{k=1}^N (y_{t,k} - \hat{y}_{t,k})^2} \quad (4)$$

where N represents the dimension of the mel-spectrum, $y_{t,k}$ denotes the k^{th} mel-spectrum component in t^{th} frame for the reference target mel-spectrum, and $\hat{y}_{t,k}$ for the synthesized mel-spectrum. Lower MCD value indicates smaller distortion.

We use Root Mean Squared Error (RMSE) as the evaluation metrics for F0 modeling, that is calculated as:

$$RMSE = \sqrt{\frac{1}{T} \sum_{t=1}^T (F0_t - \widehat{F0}_t)^2} \quad (5)$$

where $F0_t$ and $\widehat{F0}_t$ denote the reference and synthesized F0 at t^{th} frame. We note that lower RMSE value suggests that the two F0 contours are more similar.

Moreover, we propose to use frame disturbance, denoted as FD, to calculate the deviation in the dynamic time warping (DTW) alignment path [98]–[100]. FD is calculated as:

$$FD = \sqrt{\frac{1}{T} \sum_{t=1}^T (a_{t,x} - a_{t,y})^2} \quad (6)$$

where $a_{t,x}$ and $a_{t,y}$ denote the x-coordinate and the y-coordinate of the t^{th} frame in the DTW alignment path. As FD represents

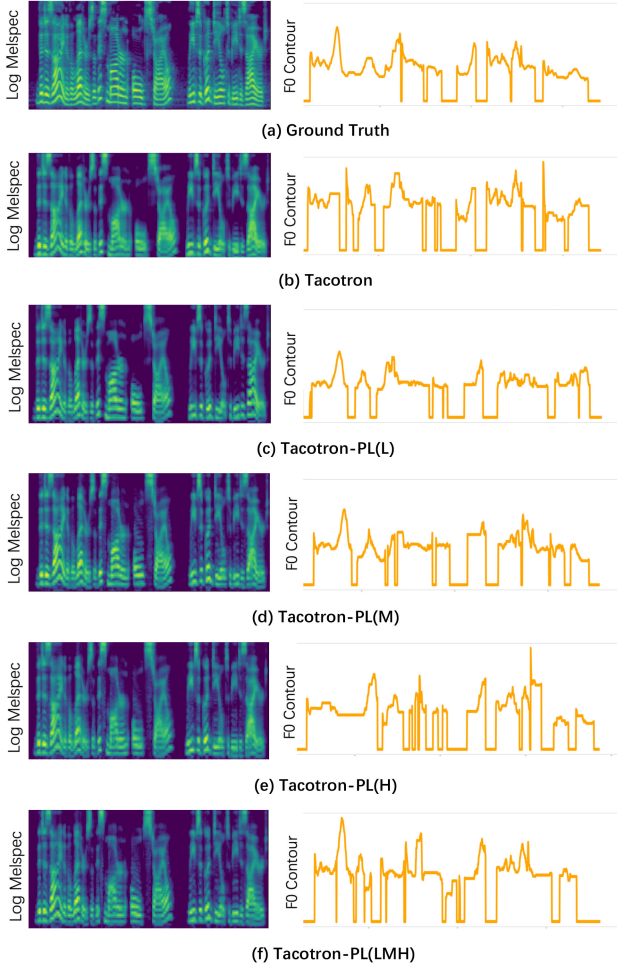


Fig. 6. Spectrogram (left) and F0 contour (right) of an utterance “The design of the letters of this modern ‘old style’ leaves a good deal to be desired.” from LJ-Speech database between the reference natural speech, labelled as Ground Truth, and five Tacotron systems. It is observed that *Tacotron-PL* models produce finer spectral details, prosodic phrasing and F0 contour that are closer to those of the reference than *Tacotron* baseline.

the duration deviation of the synthesized speech from the target, it is a proxy to show the duration distortion. A larger value indicates poor duration modeling performance and a smaller value indicates otherwise.

2) *Spectral Modeling*: We observe that all implementations of *Tacotron-PL* model consistently provide lower MCD values than *Tacotron* baseline, with *Tacotron-PL(L)* representing the lowest MCD, as can be seen in Table I. We also visualize the spectrograms of same speech content synthesized by five different models, together with that of the reference natural speech in Fig. 6. A visual inspection of the spectrograms suggests that *Tacotron-PL* models consistently provide finer spectral details than *Tacotron* baseline.

3) *F0 Modeling*: Fundamental frequency, or F0, is an essential prosodic feature of speech [32], [35]. As there is no guarantee that synthesized speech and reference speech have the same length, we apply DTW [101] to align speech pairs and calculate RMSE between the F0 contour of them. The results

are reported in Table I. It is observed that *Tacotron-PL* models consistently generate F0 contours which are closer to reference speech than *Tacotron* baseline.

We note that both F0 and prosody style contributes to RMSE measurement. To show the effect of various deep style features on the F0 contours, we also plot the F0 contours of the utterances in Fig. 6. A visual inspection suggests that the *Tacotron-PL* models benefit from the perceptual loss training, and produce F0 contour with a better fit to that of the reference speech, with *Tacotron-PL(L)* producing the best fit (see Fig. 6(c)).

4) *Duration Modeling*: Frame disturbance is a proxy to the duration difference [100] between synthesized speech and reference natural speech. We report frame disturbance of five systems in Table I. As shown in Table I, *Tacotron-PL* models obtain significantly lower FD value than *Tacotron* baseline, with *Tacotron-PL(L)* giving the lowest FD. From Fig. 6, we can also observe that *Tacotron-PL(L)* example clearly provides a better duration prediction than other models. We can conclude that perceptual loss training with style reconstruction loss helps Tacotron to achieve a more accurate rendering of prosodic patterns.

5) *Deep Style Features*: We compare four different deep style features by evaluating the performance of their use in *Tacotron-PL* models, namely *Tacotron-PL(L)*, *Tacotron-PL(M)*, *Tacotron-PL(H)* and *Tacotron-PL(LMH)*.

In supervised feature learning, the features that are near the input layer are related to the low level features, while those that are near the output are related to the supervision target, that are the categorical labels of the emotion. While we expect the style descriptors to capture utterance level prosody style, we don’t want the style reconstruction loss function to directly relate to emotion categories. Hence, the lower level deep features, Ψ_{low} , as illustrated in Fig. 4, would be more appropriate than the higher level deep features, such as Ψ_{middle} and Ψ_{high} .

We observe that Ψ_{low} is more descriptive than other deep style features for perceptual loss evaluation, as reported in spectral modeling (MCD), F0 modeling (RMSE), duration modeling (FD) for *Tacotron-PL* experiment in Table I. The observations confirm our intuition and the analysis in Fig. 4.

D. Subjective Evaluation

We conduct listening experiments to evaluate several aspects of the synthesized speech, and the choice of deep style features for $Loss_{style}$. Griffin-Lim algorithm [44] and neural vocoder are employed to generate the speech waveform. We choose WaveRNN vocoder which follows the same parameter settings as [45] since it’s the first sequential neural model for real-time audio synthesis [45].

1) *Voice Quality*: Each audio is listened by 15 subjects, each of which listens to 150 synthesized speech samples. We first evaluate the voice quality in terms of mean opinion score (MOS) among *Tacotron*, *Tacotron-PL(L)*, *Tacotron-PL(M)*, *Tacotron-PL(H)*, and *Tacotron-PL(LMH)*. As shown in Fig. 7, *Tacotron-PL* models consistently outperforms *Tacotron* baseline with either Griffin-Lim algorithm or WaveRNN vocoder, while *Tacotron-PL(L)* achieves the best result. Note that WaveRNN

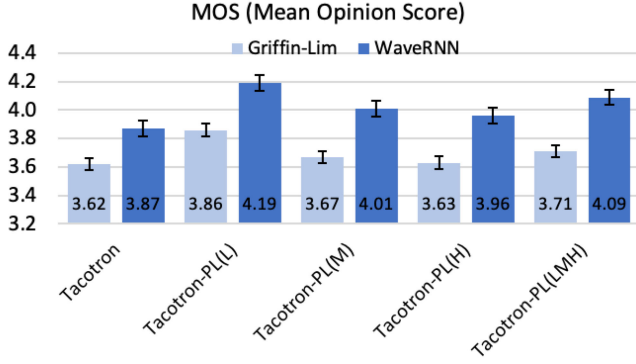


Fig. 7. The mean opinion scores (MOS) of five systems evaluated by 15 listeners, with 95% confidence intervals computed from the t-test.

vocoder achieves better speech quality than Griffin-Lim algorithm, we conduct the subsequent listening experiments only with the speech samples generated by WaveRNN vocoder.

2) *Expressiveness*: In the objective evaluations and MOS listening tests, *Tacotron-PL(L)* and *Tacotron-PL(LHM)* consistently offer better results. We next focus on comparing *Tacotron-PL(L)* and *Tacotron-PL(LHM)* with *Tacotron* baseline. We first conduct the AB preference test to assess speech expressiveness of the systems. Each audio is listened by 15 subjects, each of which listens to 120 synthesized speech samples. Table II reports the speech expressiveness evaluation results. We note that *Tacotron-PL(L)* outperforms both *Tacotron* baseline and *Tacotron-PL(LMH)* in the preference test. The results suggest that Ψ_{low} is more effective than other deep style features to inform the speech style.

3) *Naturalness*: We further conduct the AB preference test to assess the naturalness of the systems. Each audio is listened by 15 subjects, each of which listens to 120 synthesized speech samples. Table II reports the naturalness evaluation results. Just like in the expressiveness evaluation, we note that *Tacotron-PL(L)* outperforms both *Tacotron* baseline and *Tacotron-PL(LMH)* in the preference test. The results confirm that Ψ_{low} is more effective to inform the speech style.

4) *Deep Style Features*: We finally conduct Best Worst Scaling (BWS) listening experiments to compare the four different *Tacotron-PL* systems with different deep style features. The subjects are invited to evaluate multiple samples derived from the different models, and choose the best and the worst sample. We perform this experiment for 18 different utterances, and each subject listens to 72 speech samples in total. Each audio is listened by 15 subjects.

Table III summarizes the results. We can see that *Tacotron-PL(L)* is selected for 80% of time as the best model and only 5% of time as the worst model, that shows Ψ_{low} is the most effective deep style features.

E. Comparison With GST-Tacotron Paradigm

We further compare *Tacotron-PL* with the state-of-the-art expressive TTS framework, i.e., GST-Tacotron [31]. The original

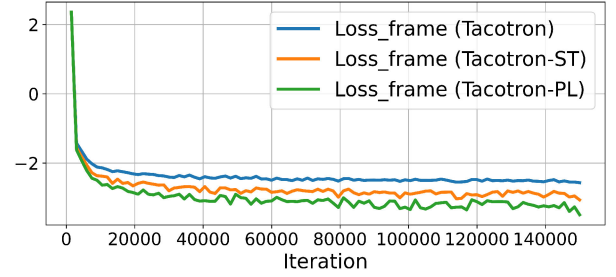


Fig. 8. The convergence trajectories of three loss values on LJ-Speech training data over the iteration steps, namely $Loss_{frame}$ for *Tacotron* baseline, *Tacotron-ST*, and $Loss_{frame}$ component as part of the $Loss_{total}$ for *Tacotron-PL*.

GST-Tacotron model [31] is focused on style control and transfer, which differs from *Tacotron-PL*. For a fair comparison, we modify the GST-Tacotron framework and build a comparative system, denoted as *Tacotron-ST*. Specifically, the reference encoder of the GST-Tacotron model is replaced with a pre-trained SER-based style descriptor as described in Sec. III-A. The style features Ψ extracted by the reference encoder informs *Tacotron-ST* the style information as GST-Tacotron does [31]. We then jointly train the whole *Tacotron-ST* framework including the pre-trained SER-based reference encoder with $Loss_{frame}$.

Tacotron-ST and *Tacotron-PL* share a similar architecture with *Tacotron* baseline [31] except that *Tacotron-ST* is augmented by a reference encoder derived from a pre-trained SER model, while *Tacotron-PL* is augmented by the proposed style reconstruction loss. In other words, both *Tacotron-ST* and *Tacotron-PL* incorporate style representations into the TTS training. We take *Tacotron-ST* under the parallel style transfer scenario [31] as the contrastive model for *Tacotron-PL*. We also use the *Tacotron* model [5] as another baseline.

We use the low-level style feature Ψ_{low} as the style embedding for *Tacotron-ST* and the deep style feature for *Tacotron-PL* in this section. We then conduct a set of experiments, following the previous experiment setup in Sec. IV-B.

1) *Convergence Trajectories of $Loss_{frame}$* : To examine the effect of the proposed training strategy, and the influence of and reference encoder and perceptual loss $Loss_{style}$, we would like to observe how $Loss_{frame}$ converges with different training schemes on the same training data. We only compare the convergence trajectories of $Loss_{frame}$ between *Tacotron* baseline, *Tacotron-ST* and the $Loss_{frame}$ component of $Loss_{total}$ for the training of *Tacotron-PL* in Fig. 8.

A lower frame-level reconstruction loss, $Loss_{frame}$, indicates a better convergence, thus a better frame level spectral prediction. We observe that the $Loss_{frame}$ component in $Loss_{total}$ achieves a lower convergence value than $Loss_{frame}$ in traditional *Tacotron* and *Tacotron-ST* training. This suggests that utterance-level style objective function of *Tacotron-PL* and reference signal supervision of *Tacotron-ST* not only optimizes style reconstruction loss, but also reduces frame-level reconstruction loss over the *Tacotron* baseline.

Finally, *Tacotron-PL* obtains the best convergence trajectories during training, that further validates the proposed

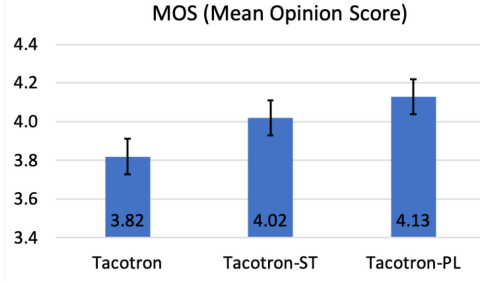


Fig. 9. The mean opinion scores (MOS) of three systems evaluated by 15 listeners, with 95% confidence intervals computed from the t-test.

frame and style training strategy. We note that the trajectories of *Tacotron-PL(M)* vs. *Tacotron-ST(M)*, *Tacotron-PL(H)* vs. *Tacotron-ST(H)*, *Tacotron-PL(LMH)* vs. *Tacotron-ST(LMH)* follow a similar pattern as *Tacotron-PL(L)* vs. *Tacotron-ST(L)*.

2) *Objective and Subjective Evaluation:* We also conduct objective and subjective evaluation experiments to compare the systems. In objective evaluation of *Tacotron-ST*, we obtain 6.58, 1.14 and 14.18 of MCD, RMSE and FD respectively. The *Tacotron-ST* results are consistently lower than those of *Tacotron*, but higher than those of *Tacotron-PL(L)* in Table I, which further confirms the effectiveness of the frame and style training strategy.

In subjective evaluation, we conduct the MOS and AB preference tests to assess the overall performance of the systems. The MOS scores are reported in Fig. 9. Each audio is listened by 15 subjects, each of which listens to 75 synthesized speech samples. It is observed that *Tacotron-PL* outperforms the *Tacotron* and *Tacotron-ST* baselines, that shows the clear advantage of frame and style training strategy. Table IV reports the AB preference test results. Each audio is listened by 15 subjects, each of which listens to 120 synthesized speech samples. All results show that *Tacotron-PL* outperforms both *Tacotron* baseline and *Tacotron-ST* significantly in terms of expressiveness and naturalness.

All the above experiments confirm that the proposed frame and style training strategy is more effective in informing the speech style than GST-Tacotron paradigm, which is encouraging.

V. CONCLUSION

We have studied a novel training strategy for Tacotron-based TTS system that includes frame and style reconstruction loss. We implement an SER model as the style descriptor to extract deep style features to evaluate the style reconstruction loss. We have conducted a series of experiments and demonstrated that the proposed Tacotron-PL training strategy outperforms the start-of-the-art Tacotron and GST-Tacotron-based baselines without the need of any add-on mechanism at run-time. While we conduct the experiments only on Tacotron, the proposed idea is applicable to other end-to-end neural TTS systems, that will be the future work in our plan.

APPENDIX

TABLE V
THE SCRIPTS OF UTTERANCES IN SIX DISTINCTIVE STYLE GROUPS FROM LJ-SPEECH DATABASE, THE DEEP STYLE FEATURES OF WHICH ARE VISUALIZED IN FIG. 4

| | |
|-------------------------------------|---|
| Group 1 (Short Question) | (1) What did he say to that? (2) Where would be the use? (3) Where is it? (4) The soldiers then? (5) What is my proposal? |
| Group 2 (Long Question) | (1) Could you advise me as to the general view we have on the American Civil Liberties Union? (2) Why not relieve Newgate by drawing more largely upon the superior accommodation which Millbank offered? (3) Who ever heard of a criminal being sentenced to catch the rheumatism or the typhus fever? (4) Why not move the city prison bodily into this more rural spot, with its purer air and greater breathing space? (5) Great Britain in many ways has advanced further along lines of social security than the United States? |
| Group 3 (Short Answer) | (1) Answer: Yes. (2) Answer: No. (3) Answer: Thank you. (4) Answer: No, sir. (5) Answer: By not talking to him. |
| Group 4 (Short Statement) | (1) In September he began to review Spanish. (2) They agree that Hosty told Revill. (3) Hardly any one. (4) They are photographs of the same scene. (5) and other details in the picture. |
| Group 5 (Long Statement) | (1) I only know that his basic desire was to get to Cuba by any means, and that all the rest of it was window dressing for that purpose. End quote. (2) He tried to start a conversation with me several times, but I would not answer. And he said that he didn't want me to be angry at him because this upsets him. (3) Several of the publications furnished the Commission with the prints they had used, or described by correspondence the retouching they had done. (4) From an examination of one of the photographs, the Commission determined the dates of the issues of the Militant and the Worker which Oswald was holding in his hand. (5) He later wrote to another official of the Worker, seeking employment, and mentioning the praise he had received for submitting his photographic work. |
| Group 6 (Digit String) | (1) Nineteen sixty-three. (2) Fourteen sixty-nine, fourteen seventy. (3) March nine, nineteen thirty-seven. Part one. (4) Section ten. March nine, nineteen thirty-seven. Part two. (5) On November eight, nineteen sixty-three. |

REFERENCES

- [1] K. Tokuda, Y. Nankaku, T. Toda, H. Zen, J. Yamagishi, and K. Oura, "Speech synthesis based on hidden Markov models," *IEEE Proc. IRE*, vol. 101, no. 5, pp. 1234–1252, May 2013.
- [2] H. Zen, A. Senior, and M. Schuster, "Statistical parametric speech synthesis using deep neural networks," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2013, pp. 7962–7966.
- [3] R. Liu, F. Bao, G. Gao, and Y. Wang, "Mongolian text-to-speech system based on deep neural network," in *Proc. Nat. Conf. Man-Mach. Speech Commun.*, 2017, pp. 99–108.
- [4] Y. Wang *et al.*, "Tacotron: A fully end-to-end text-to-speech synthesis model," in *Proc. INTERSPEECH*, 2017, pp. 4006–4010.
- [5] J. Shen *et al.*, "Natural TTS synthesis by conditioning wavenet on MEL spectrogram predictions," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2018, pp. 4779–4783.
- [6] R. Liu, B. Sisman, F. Bao, G. Gao, and H. Li, "Wavetts: Tacotron-based tts with joint time-frequency domain loss," in *Proc. Odyssey 2020 The Speaker and Language Recognition Workshop*, 2020, pp. 245–251.
- [7] Y. Lee and T. Kim, "Robust and fine-grained prosody control of end-to-end speech synthesis," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2019, pp. 5911–5915.
- [8] A. J. Hunt and A. W. Black, "Unit selection in a concatenative speech synthesis system using a large speech database," in *ICASSP IEEE Int. Conf. Acoust., Speech, Signal Process.*, 1996, pp. 373–376.
- [9] K. Tokuda, H. Zen, and A. W. Black, "An HMM-based speech synthesis system applied to english," in *Proc. IEEE Speech Synth. Workshop*, 2002, pp. 227–230.

- [10] R. Liu, B. Sisman, Y. Lin, and H. Li, "Fasttalker: A neural text-to-speech architecture with shallow and group autoregression," *Neural Netw.*, vol. 141, pp. 306–314, 2021.
- [11] Y.-A. Chung, Y. Wang, W.-N. Hsu, Y. Zhang, and R. Skerry-Ryan, "Semi-supervised training for improving data efficiency in end-to-end speech synthesis," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2019, pp. 6940–6944.
- [12] M. He, Y. Deng, and L. He, "Robust sequence-to-sequence acoustic modeling with stepwise monotonic attention for neural TTS," in *Proc. INTERSPEECH*, 2019, pp. 1293–1297.
- [13] H.-T. Luong, X. Wang, J. Yamagishi, and N. Nishizawa, "Training multi-speaker neural text-to-speech systems using speaker-imbalanced speech corpora," in *Proc. INTERSPEECH*, 2019, pp. 1303–1307.
- [14] R. Liu, B. Sisman, J. Li, F. Bao, G. Gao, and H. Li, "Teacher-student training for robust tacotron-based tts," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2020, pp. 6274–6278.
- [15] R. Liu, B. Sisman, and H. Li, "Graphspeech: Syntax-aware graph attention network for neural speech synthesis," 2020, *arXiv:2010.12423*.
- [16] T. Hayashi, A. Tamamori, K. Kobayashi, K. Takeda, and T. Toda, "An investigation of multi-speaker training for wavenet vocoder," in *Proc. IEEE Autom. Speech Recognit. Understanding Workshop*, 2017, pp. 712–718.
- [17] K. Chen, B. Chen, J. Lai, and K. Yu, "High-quality voice conversion using spectrogram-based wavenet vocoder," in *Proc. INTERSPEECH*, 2018, pp. 1993–1997.
- [18] T. Okamoto, T. Toda, Y. Shiga, and H. Kawai, "Real-time neural text-to-speech with sequence-to-sequence acoustic model and WaveGlow or single Gaussian WaveRNN vocoders," in *Proc. INTERSPEECH*, 2019, pp. 1308–1312.
- [19] B. Sisman, M. Zhang, and H. Li, "A voice conversion framework with tandem feature sparse representation and speaker-adapted wavenet vocoder," in *Proc. INTERSPEECH*, 2018, pp. 1978–1982.
- [20] B. Sisman, M. Zhang, and H. Li, "Group sparse representation with WaveNet vocoder adaptation for spectrum and prosody conversion," *IEEE/ACM Trans. Audio, Speech Lang. Process.*, vol. 27, no. 6, pp. 1085–1097, Jun. 2019.
- [21] B. Sisman, M. Zhang, S. Sakti, H. Li, and S. Nakamura, "Adaptive wavenet vocoder for residual compensation in GAN-based voice conversion," in *Proc. IEEE Spoken Lang. Technol. Workshop*, 2018, pp. 282–289.
- [22] J. Hirschberg, "Pragmatics and intonation," *The Handbook of Pragmatics*, pp. 515–537, 2004.
- [23] R. Liu, B. Sisman, F. Bao, J. Yang, G. Gao, and H. Li, "Exploiting morphological and phonological features to improve prosodic phrasing for mongolian speech synthesis," in *Proc. IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 29, pp. 274–285, 2021, doi: [10.1109/TASLP.2020.3040523](https://doi.org/10.1109/TASLP.2020.3040523).
- [24] H. Luong, S. Takaki, G. E. Henter, and J. Yamagishi, "Adapting and controlling DNN-based speech synthesis using input codes," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, 2017, pp. 4905–4909.
- [25] W.-C. Lin, Y. Tsao, F. Chen, and H.-M. Wang, "Investigation of neural network approaches for unified spectral and prosodic feature enhancement," in *Proc. Asia-Pacific Signal Inf. Process. Assoc. Annu. Summit Conf.*, 2019, pp. 1179–1184.
- [26] Z. Hodari, O. Watts, and S. King, "Using generative modelling to produce varied intonation for speech synthesis," in *Proc. 10th ISCA Speech Synth. Workshop*, 2019, pp. 239–244.
- [27] Y. Zhao, H. Li, C.-I. Lai, J. Williams, E. Cooper, and J. Yamagishi, "Improved prosody from learned F0 codebook representations for VQ-VAE speech waveform reconstruction," in *Proc. Interspeech*, 2020, pp. 4417–4421.
- [28] Z. Hodari, C. Lai, and S. King, "Perception of prosodic variation for speech synthesis using an unsupervised discrete representation of F0," in *Proc. 10th Int. Conf. Speech Prosody*, 2020, pp. 965–969.
- [29] K. Silverman *et al.*, "ToBI: A standard for labeling english prosody," in *Proc. 2nd Int. Conf. Spoken Lang. Process.*, 1992, pp. 867–870.
- [30] P. Taylor and A. W. Black, "Assigning phrase breaks from part-of-speech sequences," *Comput. Speech Lang.*, vol. 12, no. 2, pp. 99–117, 1998.
- [31] Y. Wang *et al.*, "Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 5180–5189.
- [32] D. Stanton, Y. Wang, and R. Skerry-Ryan, "Predicting expressive speaking style from text in end-to-end speech synthesis," in *Proc. IEEE Spoken Lang. Technol. Workshop*, 2018, pp. 595–602.
- [33] R. Skerry-Ryan *et al.*, "Towards end-to-end prosody transfer for expressive speech synthesis with tacotron," in *Proc. 35th Int. Conf. Mach. Learn. PMLR*, 2018, pp. 4693–4702.
- [34] G. Sun, Y. Zhang, R. J. Weiss, Y. Cao, H. Zen, and Y. Wu, "Fully-hierarchical fine-grained prosody modeling for interpretable speech synthesis," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2020, pp. 6264–6268.
- [35] G. Sun *et al.*, "Generating diverse and natural text-to-speech samples using a quantized fine-grained vae and autoregressive prosody prior," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2020, pp. 6699–6703.
- [36] K. Zhou, B. Sisman, R. Liu, and H. Li, "Seen and unseen emotional style transfer for voice conversion with a new emotional speech dataset," 2020, *864arXiv:2010.14794*.
- [37] A. Dosovitskiy and T. Brox, "Generating images with perceptual similarity metrics based on deep networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2016, pp. 658–666.
- [38] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 694–711.
- [39] Q. Chen and V. Koltun, "Photographic image synthesis with cascaded refinement networks," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 1511–1520.
- [40] A. Wright and V. Válimäki, "Perceptual loss function for neural modeling of audio systems," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2020, pp. 251–255.
- [41] T. Thiede *et al.*, "Peq-a-the itu standard for objective measurement of perceived audio quality," *J. Audio Eng. Soc.*, vol. 48, no. 1/2, pp. 3–29, 2000.
- [42] A. W. Rix, J. G. Beerends, M. P. Hollier, and A. P. Hekstra, "Perceptual evaluation of speech quality (PESQ)-A new method for speech quality assessment of telephone networks and codecs," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2001, pp. 749–752.
- [43] V. Emiya, E. Vincent, N. Harlander, and V. Hohmann, "The peass toolkit-perceptual evaluation methods for audio source separation," in *Proc. 9th Int. Conf. on Latent Variable Analysis and Signal Separation*, 2010.
- [44] D. Griffin and J. Lim, "Signal estimation from modified short-time fourier transform," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 32, no. 2, pp. 236–243, Apr. 1984.
- [45] N. Kalchbrenner *et al.*, "Efficient neural audio synthesis," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 2410–2419.
- [46] A. v. d. Oord *et al.*, "Wavenet: A generative model for raw audio," in *Proc. 9th ISCA Speech Synthesis Workshop*, 2016, p. 125.
- [47] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.
- [48] K. Emir Ak, A. Kassim, J. Hwee Lim, and J. Yew Tham, "Learning attribute representations with localization for flexible fashion search," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 7708–7717.
- [49] K. Emir Ak, J. Hwee Lim, J. Yew Tham, and A. Kassim, "Semantically consistent hierarchical text to fashion image synthesis with an enhanced-attentional generative adversarial network," in *Proc. IEEE Int. Conf. Comput. Vis. Workshops*, 2019, pp. 3121–3124.
- [50] B. Sisman, J. Yamagishi, S. King, and H. Li, "An overview of voice conversion and its challenges: From statistical modeling to deep learning," in *Proc. IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 29, 2021, pp. 132–157, doi: [10.1109/TASLP.2020.3038524](https://doi.org/10.1109/TASLP.2020.3038524).
- [51] Y. Yasuda, X. Wang, S. Takaki, and J. Yamagishi, "Investigation of enhanced tacotron text-to-speech synthesis systems with self-attention for pitch accent language," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2019, pp. 6905–6909.
- [52] M. S. Ribeiro and R. A. J. Clark, "A multi-level representation of F0 using the continuous wavelet transform and the discrete cosine transform," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2015, pp. 4909–4913.
- [53] A. Wennerstrom, *The Music of Everyday Speech Prosody and Discourse Analysis*. London, U.K.: Oxford, 2001, pp. 153–158.
- [54] D. R. Ladd, "Intonational Phonology," Cambridge, U.K.: Cambridge, 2008, pp. 153–158.
- [55] Y. XU, "Speech prosody: A methodological review," *J. Speech*, vol. 1, no. 1, pp. 85–115, 2011.
- [56] B. Şişman, H. Li, and K. C. Tan, "Transformation of prosody in voice conversion," in *Proc. Asia-Pacific Signal Inf. Process. Assoc. Annu. Summit Conf.*, 2017, pp. 1537–1546.

- [57] J. Latorre, "Multilevel parametric-base F0 model for speech synthesis," in *Proc. Ann. Conf. Int. Speech Commun. Assoc., INTERSPEECH*, 2008, pp. 2274–2277.
- [58] Z. Wu, T. Kinnunen, E. S. Chng, and H. Li, "Text-independent F0 transformation with non-parallel data for voice conversion," in *Proc. INTERSPEECH*, 2010, pp. 1732–1735.
- [59] G. Sanchez, H. Silen, J. Nurminen, and M. Gabbouj, "Hierarchical modeling of F0 contours for voice conversion," in *Proc. INTERSPEECH*, 2014, pp. 2318–2321.
- [60] M. Vainio *et al.*, "Continuous wavelet transform for analysis of speech prosody," TRASP 2013-Tools and resources for the analysis of speech prosody, an interspeech 2013 satellite event, Aug. 30, 2013, *Laboratoire Parole et Lang.*, Aix-en-Provence, France, Proceedings, pp. 78–81, 2013.
- [61] A. Suni, D. Aalto, T. Raitio, P. Alku, and M. Vainio, "Wavelets for intonation modeling in HMM speech synthesis," in *Proc. 8th ISCA Speech Synth. Workshop*, 2014, pp. 285–290.
- [62] H. Ming, D. Huang, L. Xie, S. Zhang, M. Dong, and H. Li, "Exemplar-based sparse representation of timbre and prosody for voice conversion," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech, Signal Process. Conf. Proc.*, 2016, pp. 5175–5179.
- [63] Z. Luo, J. Chen, T. Takiguchi, and Y. Ariki, "Emotional voice conversion with adaptive scales F0 based on wavelet transform using limited amount of emotional data," in *Proc. INTERSPEECH*, 2017, pp. 3399–3403.
- [64] Z. Luo, J. Chen, T. Takiguchi, and Y. Ariki, "Emotional voice conversion using neural networks with arbitrary scales F0 based on wavelet transform," *EURASIP J. Audio, Speech, Music Process.*, vol. 2017, no. 1, pp. 1–13, 2017.
- [65] H. Ming, D. Huang, L. Xie, J. Wu, M. Dong, and H. Li, "Deep bidirectional LSTM modeling of timbre and prosody for emotional voice conversion," in *Proc. INTERSPEECH*, 2016, pp. 2453–2457.
- [66] J. Yamagishi, K. Onishi, T. Masuko, and T. Kobayashi, "Modeling of various speaking styles and emotions for HMM-based speech synthesis," in *Proc. 8th Eur. Conf. Speech Commun. Technol.*, 2003, pp. 2461–2464.
- [67] M. Tachibana, J. Yamagishi, K. Onishi, T. Masuko, and T. Kobayashi, "HMM-based speech synthesis with various speaking styles using model interpolation," in *Proc. Speech Prosody, Int. Conf.*, 2004, pp. 1–4.
- [68] F. Eyben *et al.*, "Unsupervised clustering of emotion and voice styles for expressive TTS," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2012, pp. 4009–4012.
- [69] H.-T. Luong, S. Takaki, G. E. Henter, and J. Yamagishi, "Adapting and controlling DNN-based speech synthesis using input codes," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2017, pp. 4905–4909.
- [70] T. Kenter, V. Wan, C.-A. Chan, R. Clark, and J. Vit, "Chive: Varying prosody in speech synthesis with a linguistically driven dynamic hierarchical conditional variational network," in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 3331–3340.
- [71] S. Kataria, P. S. Nidadavolu, J. Villalba, N. Chen, P. Garc  Perera, and N. Dehak, "Feature enhancement with deep feature losses for speaker verification," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2020, pp. 7584–7588.
- [72] M. Kawanaka, Y. Koizumi, R. Miyazaki, and K. Yatabe, "Stable training of dnn for speech enhancement based on perceptually-motivated black-box cost function," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2020, pp. 7524–7528.
- [73] F. G. Germain, Q. Chen, and V. Koltun, "Speech denoising with deep feature losses," in *Proc. INTERSPEECH*, 2019, pp. 2723–2727.
- [74] C.-C. Lo *et al.*, "Mosnet: Deep learning-based objective assessment for voice conversion," in *Proc. INTERSPEECH*, 2019, pp. 1541–1545.
- [75] K.-S. Lee, "Voice conversion using a perceptual criterion," *Appl. Sci.*, vol. 10, no. 8, 2020, Art. no. 2884.
- [76] A. van den Oord *et al.*, "Parallel WaveNet: Fast high-fidelity speech synthesis," in *Proc. 35th Int. Conf. Mach. Learn., Ser. Proc. Mach. Learn. Res.*, J. Dy and A. Krause, Eds., vol. 80. Stockholmms  ssan, Stockholm Sweden: PMLR, 10–15 Jul. 2018, pp. 3918–3926.
- [77] Z. Cai, C. Zhang, and M. Li, "From speaker verification to multispeaker speech synthesis, deep transfer with feedback constraint," in *Proc. Interspeech*, 2020, pp. 3974–3978.
- [78] E. Kim and J. W. Shin, "DNN-based emotion recognition based on bottleneck acoustic features and lexical features," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2019, pp. 6720–6724.
- [79] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *Statist.*, vol. 1050, p. 1, 2014.
- [80] G. Zhong, L. Wang, and J. Dong, "An overview on data representation learning: From traditional feature learning to recent deep learning," *J. Finance Data Sci.*, vol. 2, no. 4, pp. 265–278, 2016.
- [81] S. Latif, R. Rana, J. Qadir, and J. Epps, "Variational autoencoders for learning latent representations of speech emotion: A preliminary study," in *Proc. INTERSPEECH*, 2018, pp. 3107–3111.
- [82] S. Zhang, S. Zhang, T. Huang, and W. Gao, "Speech emotion recognition using deep convolutional neural network and discriminant temporal pyramid matching," *IEEE Trans. Multimedia*, vol. 20, no. 6, pp. 1576–1590, Jun. 2018.
- [83] R. Lotfian and C. Busso, "Curriculum learning for speech emotion recognition from crowdsourced labels," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 27, no. 4, pp. 815–826, Apr. 2019.
- [84] I. R. Murray and J. L. Arnott, "Toward the simulation of emotion in synthetic speech: A review of the literature on human vocal emotion," *J. Acoust. Soc. Amer.*, vol. 93, no. 2, pp. 1097–1108, 1993.
- [85] O. Pierre-Yves, "The production and recognition of emotions in speech: Features and algorithms," *Int. J. Human-Comput. Stud.*, vol. 59, no. 1–2, pp. 157–183, 2003.
- [86] 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. IEEE Asia-Pacific Signal Inf. Process. Assoc. Annu. Summit Conf.*, 2019, pp. 623–627.
- [87] Y. Gao, W. Zheng, Z. Yang, T. Kohler, C. Fuegen, and Q. He, "Interactive Text-to-Speech System via Joint Style Analysis," in *Proc. Interspeech*, 2020, pp. 4447–4451.
- [88] S.-Y. Um, S. Oh, K. Byun, I. Jang, C. Ahn, and H.-G. Kang, "Emotional speech synthesis with rich and granularized control," in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process.*, 2020, pp. 7254–7258.
- [89] C. Busso *et al.*, "Iemocap: Interactive emotional dyadic motion capture database," *Lang. Resour. Eval.*, vol. 42, no. 4, 2008, Art. no. 335.
- [90] I. Goodfellow, Y. Bengio, and A. Courville, D. Learning. Cambridge, MA, USA: MIT Press, 2016.
- [91] D. Yu and M. L. Seltzer, "Improved bottleneck features using pretrained deep neural networks," in *Proc. INTERSPEECH*, 2011, pp. 237–240.
- [92] M. Chen, X. He, J. Yang, and H. Zhang, "3-D convolutional recurrent neural networks with attention model for speech emotion recognition," *IEEE Signal Process. Lett.*, vol. 25, no. 10, pp. 1440–1444, Oct. 2018.
- [93] K. Greff, R. K. Srivastava, J. Koutn  k, B. R. Steunebrink, and J. Schmidhuber, "LSTM: A search space odyssey," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 28, no. 10, pp. 2222–2232, Oct. 2017.
- [94] K. Ito, "The LJ speech dataset," 2017. [Online]. Available: <https://keithito.com/LJ-Speech-Dataset/>
- [95] L. v. d. Maaten and G. Hinton, "Visualizing data using t-SNE," *J. Mach. Learn. Res.*, vol. 9, no. 11, pp. 2579–2605, 2008.
- [96] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proc. Int. Conf. Mach. Learn.*, 2015, pp. 448–456.
- [97] R. Kubichek, "Mel-Cepstral distance measure for objective speech quality assessment," in *Proc. IEEE Pacific Rim Conf. Commun. Comput. Signal Process.*, 1993, pp. 125–128.
- [98] B. Sisman, G. Lee, H. Li, and K. C. Tan, "On the analysis and evaluation of prosody conversion techniques," in *Proc. Int. Conf. Asian Lang. Process.*, 2017, pp. 44–47.
- [99] A. Z. Jusoh, R. Togneri, S. Nordholm, N. Sulaiman, and M. H. Khairolanuar, "The investigation of frame disturbance (FD) in perceptual evaluation speech quality (PESQ) as a perceptual metric," *ARPN J. Eng. Appl. Sci.*, vol. 10, no. 15, pp. 6365–6369, 2015.
- [100] C. Gupta, H. Li, and Y. Wang, "Perceptual evaluation of singing quality," in *Proc. IEEE Asia-Pacific Signal Inf. Process. Assoc. Annu. Summit Conf.*, 2017, pp. 577–586.
- [101] M. M  ller, "Dynamic time warping," *Inf. Retrieval Music Motion*, pp. 69–84, 2007.



Rui Liu (Member, IEEE) received the B.S. degree from the Department of Software, Taiyuan University of Technology, Taiyuan, China, in 2014, and the Ph.D. degree in computer science and technology from the Inner Mongolia Key Laboratory of Mongolian Information Processing Technology, Inner Mongolia University, Hohhot, China, in 2020. He is also an exchange Ph.D. Candidate with the Department of Electrical and Computer Engineering, National University of Singapore (NUS), Singapore, funded by China Scholarship Council. He is currently a joint Postdoctoral Research Fellow with NUS and Singapore University of Technology and Design, Singapore. His research interests include prosody and acoustic modeling for speech synthesis, machine learning, and natural language processing.



Berrak Sisman (Member, IEEE) received the Ph.D. degree in electrical and computer engineering from the National University of Singapore, Singapore, in 2020, fully funded by A*STAR Graduate Academy under Singapore International Graduate Award (SINGA). She is currently an Assistant Professor with the Singapore University of Technology and Design (SUTD), Singapore. She is also an Affiliated Researcher with the National University of Singapore, Singapore. Prior to joining SUTD, she was a Postdoctoral Research Fellow with the National

University of Singapore, and a Visiting Researcher with Columbia University, New York City, NY, USA. She was also an exchange Ph.D. Student with the University of Edinburgh, Edinburgh, U.K., and a Visiting Scholar with The Centre for Speech Technology Research, University of Edinburgh in 2019. She was attached to RIKEN Advanced Intelligence Project, Japan in 2018. Her research interests include machine learning, signal processing, speech synthesis, and voice conversion. She was the General Coordinator of the Student Advisory Committee of International Speech Communication Association.



Guanglai Gao received the B.S. degree from Inner Mongolia University, Hohhot, China, in 1985 and the M.S. degree from the National University of Defense Technology, Changsha, China, in 1988. He was a Visiting Researcher with the University of Montreal, Montreal, QC, Canada. He is currently a Professor with the Department of Computer Science, Inner Mongolia University, Hohhot, China. His research interests include artificial intelligence and pattern recognition.



Haizhou Li (Fellow, IEEE) received the B.Sc., M.Sc., and Ph.D. degrees in electrical and electronic engineering from the South China University of Technology, Guangzhou, China, in 1984, 1987, and 1990, respectively. He is currently a Professor with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore. Prior to joining NUS, he taught with the University of Hong Kong, Hong Kong, from 1988 to 1990 and South China University of Technology, Guangzhou, China, from 1990 to 1994. He was a Visiting Professor with

CRIN, France, from 1994 to 1995, the Research Manager with the Apple-ISS Research Centre from 1996 to 1998, the Research Director of Lernout & Hauspie Asia Pacific from 1999 to 2001, the Vice President of InfoTalk Corp. Ltd., from 2001 to 2003, and the Principal Scientist and Department Head of Human Language Technology in the Institute for Infocomm Research, Singapore, from 2003 to 2016. His research interests include automatic speech recognition, speaker and language recognition, and natural language processing. From 2015 to 2018, he was the Editor-in-Chief of the IEEE/ACM TRANSACTIONS ON AUDIO, SPEECH AND LANGUAGE PROCESSING and from 2012 to 2018, a Member of the Editorial Board of Computer Speech and Language. He was an elected Member of IEEE Speech and Language Processing Technical Committee from 2013 to 2015, the President of the International Speech Communication Association from 2015 to 2017, the President of Asia Pacific Signal and Information Processing Association from 2015 to 2016, and the President of Asian Federation of Natural Language Processing from 2017 to 2018. He was the General Chair of ACL 2012, INTERSPEECH 2014 and ASRU 2019. Dr Li is a Fellow of the IEEE and the ISCA. He was the recipient of the National Infocomm Award 2002 and the President's Technology Award 2013 in Singapore. He was named one of the two Nokia Visiting Professors in 2009 by the Nokia Foundation and Bremen Excellence Chair Professor in 2019.