# From Speaker Verification to Multispeaker Speech Synthesis, Deep Transfer with Feedback Constraint

*Zexin Cai[*], Chuxiong Zhang[†], Ming Li[*,†]*

[*]Electrial & Computer Engineering, Duke University, Durham, NC, United States
[†] Data Science Research Center, Duke Kunshan University, Kunshan, China
`ming.li369@duke.edu`

## Abstract

High-fidelity speech can be synthesized by end-to-end text-to-speech models in recent years. However, accessing and controlling speech attributes such as speaker identity, prosody, and emotion in a text-to-speech system remains a challenge. This paper presents a system involving feedback constraints for multispeaker speech synthesis. We manage to enhance the knowledge transfer from the speaker verification to the speech synthesis by engaging the speaker verification network. The constraint is taken by an added loss related to the speaker identity, which is centralized to improve the speaker similarity between the synthesized speech and its natural reference audio. The model is trained and evaluated on publicly available datasets. Experimental results, including visualization on speaker embedding space, show significant improvement in terms of speaker identity cloning in the spectrogram level. In addition, synthesized samples are available online for listening. [1]

**Index Terms**: Text-to-speech, multi-speaker speech synthesis, speaker embedding, end-to-end

## 1. Introduction

Speech synthesis, also known as text-to-speech (TTS), specifies the technique that achieves the transformation from text to audio waveform. It has been widely used in our daily life, e.g., navigation systems, audiobooks, and virtual assistants. The performance of the TTS system has been further improved recently by adopting the end-to-end neural network framework [1, 2, 3, 4]. The end-to-end principle is applied in the TTS model by a cohesive and autoregressive chain of neural network structures that are connected by well-defined input-output features. For instance, the state-of-the-art system Tacotron2 [2] consists of an encoder-decoder architecture and a neural vocoder Wavenet[5].

Extensions on these models have been developed for allowing the TTS system to control the speech characteristics. These extensions are able to enrich the expressiveness of the synthesized voice and further enhance the robustness of TTS systems. For example, Yuxuan Wang et al. proposed the style tokens to uncover the latent space regarding speech attributes that are hard to define and label [6, 7]. The models are jointly trained with the Tacotron-based TTS architecture in an unsupervised manner. On the other hand, controlling speech attributes that have easily found labels (e.g., language, emotion, and speaker identity) have also been investigated [8, 9, 10]. Typically, the speech attribute is controlled with a TTS model by conditioning the synthesizer with the vector representation called embedding.

The multispeaker TTS system is one of the extensions, which is developed to clone and manage distinct voices either seen or not seen during training. Most systems use the speaker embedding to characterize the expected voice and speaking style in the multispeaker TTS system [10, 11, 12], while speaker adaptation can also be used for speaker transfer TTS modeling [13]. Voice cloning by speaker adaptation acquires more data and computational resource for the target voice and usually is less robust compared with cloning by speaker embedding [14]. The speaker verification system plays an essential role in the multispeaker TTS system for cloning unseen voices. Eliya Nachmani et al. has proposed an approach where the speaker verification system and the synthesizer are jointly trained [15]. However, the knowledge for discriminative speaker representations is limited by the training dataset in this case. Then Ye Jia et al. further investigated the knowledge transfer in terms of speaker characteristics by decoupling these two tasks, where the speaker verification network is trained with a dataset that contains a larger amount of speakers but is not suitable for TTS training [10]. The discriminative speaker embedding extracted from the speaker verification network is used for conditioning the TTS synthesizer and leads to better performance on open-set voice cloning.

Although the model proposed in [10] increases the robustness of the synthesizer for open-set multispeaker synthesis, the speaker's similarity is not close between the synthesized speech and the speaker's natural speech. Concerning the same speaker, the embeddings extracted from synthesized speech and those extracted from natural speech may have two distinct distributions. To further transfer the knowledge from a speaker verification model to the speech synthesizer, we propose a multispeaker TTS model with the feedback constraint toward the speaker embedding space. Specifically, an added score associated with the speaker similarity is performed by the verification network for forcing the synthesizer to derive the knowledge for speaker identity cloning. The proposed method is evaluated on publicly available datasets. As demonstrated in the visualization of the embedding space, speaker embeddings extracted from our synthesized speech lies in the same cluster as those from natural speech. Therefore, the model may be useful for data augmentation and the white-box spoofing attack toward speaker verification in the future.

This paper is organized as follows: section 2 describes the related works in terms of speaker verification and speech synthesis. Our proposed system is presented in section 3. Experimental setup and results are shown in section 4. Finally, we conclude the paper in section 5.

## 2. Related works

### 2.1. Speaker verification

Open-set multispeaker TTS highly depends on speaker representations for conditioning the synthesizer to copy the desired
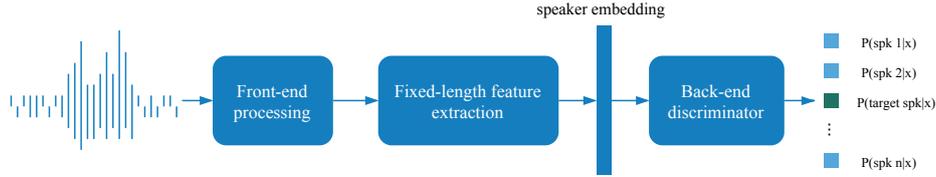
---

[1]https://caizexin.github.io/mlspk-syn-samples/index.html

Figure 1: *The overall training framework of deep speaker verification systems*

voice. To that end, speaker verification systems, especially text-independent systems, are often used for feature extraction regarding their discriminative speaker representations.

In the speaker verification field, deep speaker feature learning systems proposed in recent years have achieved comparable performance or even surpassed the classical i-vector systems [16, 17, 18]. The overall training architecture of the deep speaker verification system is shown in figure 1. Specifically, the speaker verification model takes variable-length audio signal $x = [x_1, x_2, x_3...x_n]$ as input and convert the signal into frequency-domain acoustic features, e.g., filter-bank energy or Mel frequency cepstral coefficients (MFCCs). Acoustic features are then fed into a neural network-based extractor to obtain the fixed dimensional speaker representation $z \in \mathcal{R}^d$, where $d$ is the dimension of the speaker embedding. Note that the back-end discriminator here in the training phase is different from the one in the evaluation phase. The discriminator in the training phase is to classify embeddings according to their target speaker labels in order to train a discriminative speaker embedding space, while the one used in the evaluation phase is to verify if two embeddings come from the same speaker.

Among deep speaker embedding systems that are developed with various DNN architectures [17, 19, 20], we follow the ResNet-based verification system [20] in our work for embedding extraction to extract time-invariant speaker embedding.

### 2.2. Multispeaker speech synthesis

Cloning and controlling speech attributes have been studied for decades in the text-to-speech (TTS) field. For voice synthesis, Yamagishi et al. proposed feature-space adaptive training for speaker-adaptive TTS [21]. The system aims to reduce the size of data and the cost for building different voices when developing statistical parametric speech systems based on Hidden Markov models (HMMs). After Tacotron2 demonstrated its ability to synthesize high-quality speech that can be as natural as human speech, extensions of Tacotron2 were proposed for speech attribute cloning by conditioning the linguistic encoder output with attribute embeddings. For instance, Yuxuan Wang et al. introduced global style tokens (GSTs) as the attribute embeddings to achieve style-control TTS synthesis [6]. The proposed model in [6], where GSTs are trained in an unsupervised manner, also helps improve the speech intelligibility when it is used for multispeaker TTS training.

On the other hand, in order to achieve zero-shot voice cloning, the speaker representation is commonly extracted by a separated model and used as the conditioned feature in multispeaker TTS models [10, 15]. In this case, the multispeaker TTS model developed for zero-shot voice cloning consists of two models, one for speaker embedding extraction and the other for TTS conversion. When the two models are trained jointly, the TTS system yields moderate performance in synthesizing voices that are unseen in the training data [15]. The reason might be because the datasets collected for speech synthesis have limited speakers, and the datasets collected for speaker analysis have no transcriptions for TTS training. Jia Ye et al. chose to train the two models individually, where the TTS model learns the speaker representation knowledge by the embedding extracted from the speaker verification model [10]. This method improves the robustness with the ability to clone unseen voices. However, two distinct clusters, representing synthesized speech and natural speech from the same speaker, are observed in the embedding space as shown in [10]. To further investigate this problem and enhance the knowledge transfer, we propose a model with a feedback constraint that engages the speaker embedding extractor. We show that by showing that embeddings from different speakers result in distinct distributions in the vector space, while embeddings from the same speaker, whether synthesized or natural, lie in the same cluster.

In our work, we use a speaker embedding extractor that is different from the model described in [10]. By the time we finalized our work, Erica et al. published a study investigating how different speaker embedding networks affects the multispeaker synthesis system [22]. In that study, the author claim that LDE-based embedding could improve speaker similarity and naturalness. Our model has a similar speech encoder as the learnable dictionary encoding-based (LDE-based) systems described in [22].

## 3. Methods

Our proposed multispeaker TTS framework is shown in figure 2. We follow the baseline end-to-end speaker verification system presented in [20] as our embedding extraction network. The Mel-spectrogram is used as the acoustic feature for both the speaker embedding extraction system and the multispeaker TTS system. As for the speaker embedding network, ResNet34 architecture is used as the encoder network, followed by a pooling layer that calculates the mean and the standard deviation of encoder outputs. Then the speaker embedding is obtained by concatenating the mean and the standard deviation. While in the training phase, a back-end classification network consisting of two fully connected layers maps embeddings to target speakers.

We use the tacotron-based model as the Mel-spectrogram prediction model. The input character sequence is converted into a vector sequence by a trainable lookup table. Then the encoder, which consists of 5 convolutional layers and a bi-directional long short-term memory (BLSTM) layer, consumes the embedding sequence and delivers the memory that represents the context and linguistic characteristics of the input text. Speaker embedding, extracted from the target audio signal, is then concatenated with the encoder output memory globally as the final encoding states.

The decoder takes steps to predict Mel-spectrograms with three modules, which are the attention mechanism, the RNN decoder and the PostNet. The attention mechanism provides the
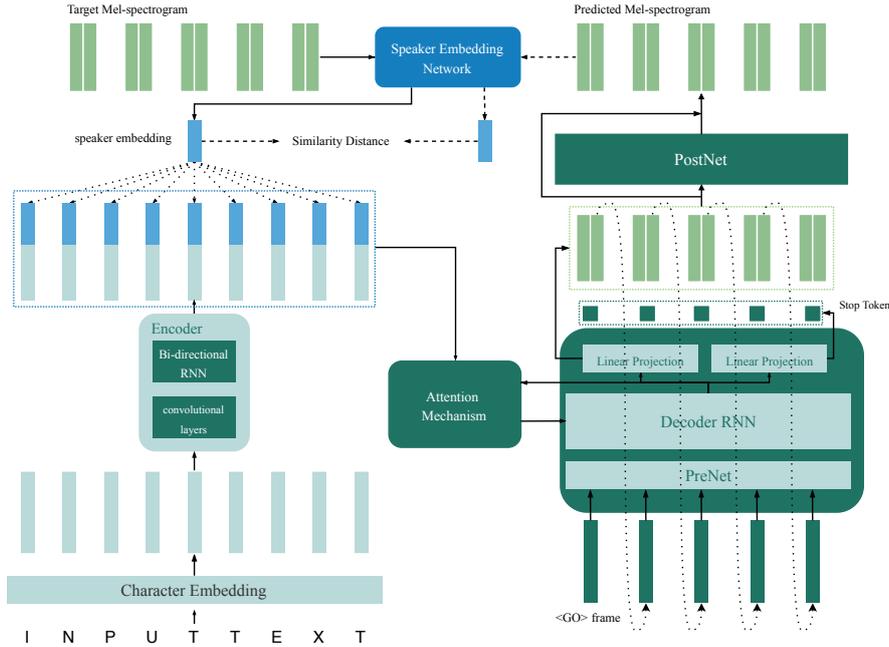
Figure 2: *Proposed multi-speaker speech synthesis model*

context vector for the decoder RNN to generate spectrograms that associate with specific encoder states for each decoding time step. In addition, it provides soft alignment between the input encoder states and the target Mel-spectrogram. For each decoding step, the decoder RNN predicts the spectrogram with respect to the context vector and the predicted result from the previous time step, where the previous predicted frame is taken by the PreNet module. Two linear projection layers are followed by the decoder RNN for predicting Mel-spectrograms and stop tokens, respectively. Stop tokens are the binary sequence that specifies the valid decoding frames, where 0 denotes a valid frame, and 1 indicates the end of the decoding process. The PostNet takes the predicted Mel-spectrogram as input to obtain the residual parameters that are related to future context since the decoder RNN is unable to foresee future frames. The predicted Mel-spectrograms are finetuned with the PostNet, which leads to high-quality outputs.

The speaker embedding network is engaged after the Post-Net during the training phase. It is added as a feedback constraint to force the TTS model to learn the speaker variety knowledge sufficiently so that the speaker characteristics extracted from synthesized Mel-spectrogram lays in the same distribution as those extracted from the natural speech from the same person. In this case, the parameters of the speaker embedding network are not updated during the training phase.

We use the cosine distance between the ground truth speaker embedding and the one extracted from the predicted Mel-spectrogram by speaker embedding network as one of the loss functions for optimizing the TTS network. Other than that, mean square error (MSE) between predicted Mel-spectrograms and the ground truth spectrogram, classification loss of the stop tokens, and the regularization loss for encoder-decoder parameters are applied to ensure correct predictions.

The Mel-spectrogram is converted back to the audio signal by the neural vocoder WaveRNN [23], which is able to generate high-quality speech at fast speed.

# 4. Experiments

We used four publicly available datasets for training and evaluation. All data from Voxceleb1 [24] and Voxceleb2 [25], with more than 7,000 speakers, are used for training the speaker verification system. The VCTK English dataset [26], which contains 109 speakers with various accents, is used for TTS model training, while data from 8 speakers are randomly excluded as the VCTK test set. For each speaker in the training set, 8 utterances are randomly picked out as the VCTK validation set. 7 speakers from the Librispeech dataset [27] are randomly selected as the cross-domain test set. All audios are downsampled to 16 kHz in our experiments.
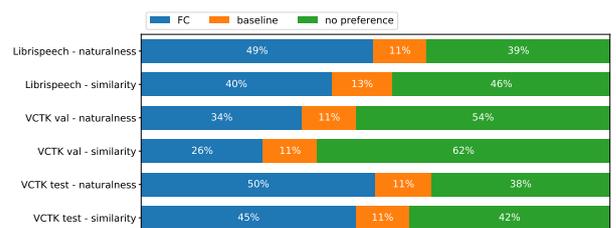


Figure 3: *Subjective preference result*

We evaluate the performance by comparing the proposed system, which has added feedback constraint (FC), with the multispeaker TTS baseline system without FC. The two systems are identified by **'baseline'** and **'FC'** in the following subjective and objective results. We first trained the baseline model until it can synthesize intelligible speech. Then the FC model is trained from the pre-trained baseline model while engaging the speaker embedding network. Both models are then trained to the same total training steps with the same batch size.
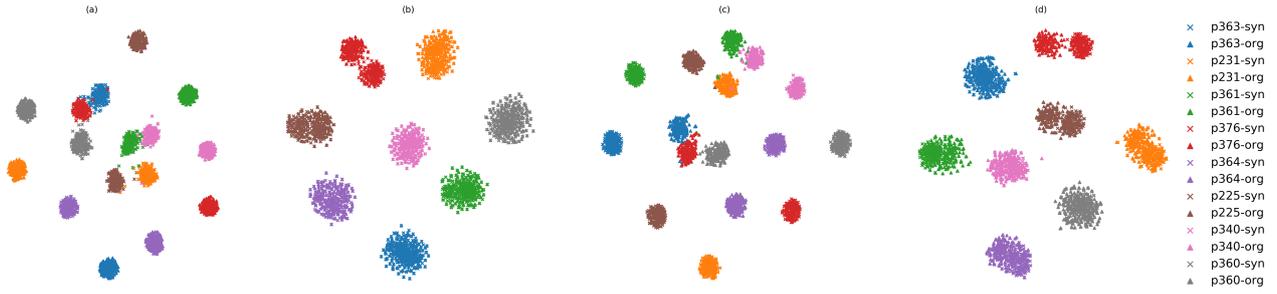
Figure 4: *Speaker embedding visualization by t-SNE for the VCTK test set. (a) baseline & text-dependent speaker embedding; (b) FC & text-dependent speaker embedding; (c) baseline & text-independent speaker embedding; (d) FC & text-independent speaker embedding;*

Table 1: *Objective evaluation results*

| | Systems | SV-EER (%) Dep / Indep | Average cosine similarity Dep / Indep |
|---|---|---|---|
| VCTK test | natural | 1.76 | - |
| | baseline | 14.72 / 13.18 | 0.403 / 0.333 |
| | FC | 8.22 / 7.68 | 0.764 / 0.577 |
| VCTK val | natural | 1.61 | - |
| | baseline | 9.22 / 8.23 | 0.472 / 0.394 |
| | FC | 5.02 / 3.42 | 0.842 / 0.67 |
| Librispeech | natural | 5.26 | - |
| | baseline | 26.84 / 26.46 | 0.222 / 0.139 |
| | FC | 16.54 / 16.11 | 0.626 / 0.389 |

### 4.1. Subjective evaluation

We asked 12 people to choose their preferable speech for pairs that contain speech synthesized by both systems. Audios in each pair are synthesized with the same text content and the conditioned embedding from the same reference audio. Each person chose their preference concerning speaker similarity and naturalness from 38 pairs, which are randomly selected from the VCTK test set, the VCTK validation set, and the Librispeech set. Preference results are shown in figure 3. For subjective evaluation, the FC system outperforms the baseline system on all three evaluation sets. For seen speakers in the training data, the speech synthesized by both systems is close. Hence people do not have a preferred choice for more than $50\%$ pairs in the VCTK validation set. Given these points, both systems can copy seen voices well, while the FC system obtains better performance on unseen voice cloning.

### 4.2. Objective evaluation

For each utterance from all evaluation datasets, we synthesized speech according to the given transcript and the embedding extracted from the original speaker's voice. Two different synthesized results were collected for each utterance. Although both are synthesized with the same reference voice, one is synthesized based on the speaker embedding extracted from the utterance that has the exact same content, while the other result is synthesized with the speaker embedding extracted from a randomly selected utterance with different content. These are identified as text-dependent **(Dep)** result and text-independent

**(Indep)** result in table 1. Speaker verification equal error rate **(SV-EER)** is used to evaluate the speaker discrimination performance for a set of embeddings. We randomly generate enrollment-verification pairs for each experiment, where half of the trials are cross-speaker pairs. We also compute the average cosine similarity between embeddings extracted from synthesized speech and the ground truth embeddings to measure the speaker similarity performance objectively. As shown in table 1, the FC system obtains significantly lower EERs than the baseline system on all evaluation sets, whether text-dependent or text-independent. The FC system also has higher average cosine similarities than the baseline system. In either case, we can conclude that the voice synthesized by the FC system is more close to the reference voice than the baseline system. The similarity is improved with the feedback constraint network.

Likewise, we can visualize the results from the embedding space by the t-Distributed Stochastic Neighbor Embedding (t-SNE), as shown in figure 4. The utterances synthesized by the baseline system with reference embeddings from the same speaker is in the same cluster, but do not have the same distribution with reference embeddings, even is closer to another speaker. For example, as shown in figure 4 $(a)$, embeddings from 'p225-syn' have a distribution that is close to 'p376-org' other than its reference speaker 'p225-org'. For the FC system, the embeddings extracted from the same voice, either synthesized or natural, lie in the same distribution in the embedding space. Therefore, the synthesized voice is more close to the original speaker for utterances synthesized by the FC system.

## 5. Conclusion

In this paper, a multispeaker TTS approach that explores the use of a speaker verification system is presented. A trained speaker verification system is incorporated into the TTS framework acting as the feedback constraint to facilitate voice cloning. Experimental evaluations, including both subjective and objective evaluations, demonstrate that our proposed system enhances the knowledge transfer from speaker verification to speech synthesis. Accordingly, our proposed method achieves significant improvement regarding voice cloning, which can be used for data augmentation or white-box spoofing attack in the future.

# 6. References

[1] S. Mehri, K. Kumar, I. Gulrajani, R. Kumar, S. Jain, J. Sotelo, A. Courville, and Y. Bengio, "Samplernn: An unconditional end-to-end neural audio generation model," in *International Conference on Learning Representations 2017*, 2017.

[2] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Skerrv-Ryan *et al.*, "Natural TTS synthesis by conditioning WaveNet on mel spectrogram predictions," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2018, pp. 4779–4783.

[3] W. Ping, K. Peng, and J. Chen, "Clarinet: Parallel wave generation in end-to-end text-to-speech," in *International Conference on Learning Representations 2019*, 2019.

[4] J. Sotelo, S. Mehri, K. Kumar, J. F. Santos, K. Kastner, A. C. Courville, and Y. Bengio, "Char2wav: End-to-end speech synthesis," in *International Conference on Learning Representations 2017*, 2017.

[5] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "WaveNet: A generative model for raw audio," *arXiv preprint arXiv:1609.03499*, 2016.

[6] Y. Wang, D. Stanton, Y. Zhang, R.-S. Ryan, E. Battenberg, J. Shor, Y. Xiao, Y. Jia, F. Ren, and R. A. Saurous, "Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis," in *Proceedings of the 35th International Conference on Machine Learning*, 2018, pp. 5180–5189.

[7] Y. Wang, R. Skerry-Ryan, Y. Xiao, D. Stanton, J. Shor, E. Battenberg, R. Clark, and R. A. Saurous, "Uncovering latent style factors for expressive speech synthesis," *arXiv preprint arXiv:1711.00520*, 2017.

[8] Y. Zhang, R. J. Weiss, H. Zen, Y. Wu, Z. Chen, R. Skerry-Ryan, Y. Jia, A. Rosenberg, and B. Ramabhadran, "Learning to Speak Fluently in a Foreign Language: Multilingual Speech Synthesis and Cross-Language Voice Cloning," in *Proc. Interspeech 2019*, 2019, pp. 2080–2084.

[9] C. Yu, H. Lu, N. Hu, M. Yu, C. Weng, K. Xu, P. Liu, D. Tuo, S. Kang, G. Lei *et al.*, "Durian: Duration informed attention network for multimodal synthesis," *arXiv preprint arXiv:1909.01700*, 2019.

[10] Y. Jia, Y. Zhang, R. Weiss, Q. Wang, J. Shen, F. Ren, P. Nguyen, R. Pang, I. L. Moreno, Y. Wu *et al.*, "Transfer learning from speaker verification to multispeaker text-to-speech synthesis," in *Advances in neural information processing systems*, 2018, pp. 4480–4490.

[11] J. Park, K. Zhao, K. Peng, and W. Ping, "Multi-speaker end-to-end speech synthesis," *arXiv preprint arXiv:1907.04462*, 2019.

[12] D. lvarez, S. Pascual, and A. Bonafonte, "Problem-Agnostic Speech Embeddings for Multi-Speaker Text-to-Speech with SampleRNN," in *Proc. 10th ISCA Speech Synthesis Workshop*, 2019, pp. 35–39.

[13] Y. Fan, Y. Qian, F. K. Soong, and L. He, "Multi-speaker modeling and speaker adaptation for DNN-based TTS synthesis," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2015, pp. 4475–4479.

[14] S. Arik, J. Chen, K. Peng, W. Ping, and Y. Zhou, "Neural voice cloning with a few samples," in *Advances in Neural Information Processing Systems*, 2018, pp. 10 019–10 029.

[15] E. Nachmani, A. Polyak, Y. Taigman, and L. Wolf, "Fitting new speakers based on a short untranscribed sample," in *Proceedings of the 35th International Conference on Machine Learning*, 2018, pp. 3683–3691.

[16] W. Cai, J. Chen, and M. Li, "Analysis of length normalization in end-to-end speaker verification system," in *Proc. Interspeech 2018*, 2018, pp. 3618–3622.

[17] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, "X-vectors: Robust dnn embeddings for speaker recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2018, pp. 5329–5333.

[18] W. Cai, J. Chen, and M. Li, "Exploring the encoding layer and loss function in end-to-end speaker and language recognition system," in *Proc. Odyssey 2018 The Speaker and Language Recognition Workshop*, 2018, pp. 74–81.

[19] D. Snyder, P. Ghahremani, D. Povey, D. Garcia-Romero, Y. Carmiel, and S. Khudanpur, "Deep neural network-based speaker embeddings for end-to-end speaker verification," in *2016 IEEE Spoken Language Technology Workshop*, 2016, pp. 165–170.

[20] W. Cai, J. Chen, J. Zhang, and M. Li, "On-the-fly data loader and utterance-level aggregation for speaker and language recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 1038–1051, 2020.

[21] J. Yamagishi, T. Nose, H. Zen, Z.-H. Ling, T. Toda, K. Tokuda, S. King, and S. Renals, "Robust speaker-adaptive HMM-based text-to-speech synthesis," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 6, pp. 1208–1230, 2009.

[22] E. Cooper, C. Lai, Y. Yasuda, F. Fang, X. Wang, N. Chen, and J. Yamagishi, "Zero-shot multi-speaker text-to-speech with state-of-the-art neural speaker embeddings," in *2020 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2020, pp. 6184–6188.

[23] N. Kalchbrenner, E. Elsen, K. Simonyan, S. Noury, N. Casagrande, E. Lockhart, F. Stimberg, A. van den Oord, S. Dieleman, and K. Kavukcuoglu, "Efficient neural audio synthesis," in *Proceedings of the 35th International Conference on Machine Learning*, 2018, pp. 2410–2419.

[24] A. Nagrani, J. S. Chung, and A. Zisserman, "Voxceleb: A large-scale speaker identification dataset," in *Proc. Interspeech 2017*, 2017, pp. 2616–2620.

[25] J. S. Chung, A. Nagrani, and A. Zisserman, "Voxceleb2: Deep speaker recognition," in *Proc. Interspeech 2018*, 2018, pp. 1086–1090.

[26] C. Veaux, J. Yamagishi, K. MacDonald *et al.*, "Superseded-CSTR VCTK corpus: English multi-speaker corpus for CSTR voice cloning toolkit," *University of Edinburgh, The Centre for Speech Technology Research*, 2016.

[27] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2015, pp. 5206–5210.