

AdaVITS: Tiny VITS for Low Computing Resource Speaker Adaptation

Kun Song^{1,3}, Heyang Xue², Xinsheng Wang¹, Jian Cong¹, Yongmao Zhang¹, Lei Xie^{1,2*}, Bing Yang³, Xiong Zhang³, Dan Su³

Audio, Speech and Language Processing Group, ¹School of Computer Science, ²School of Software, Northwestern Polytechnical University, Xi'an, China

³Cloud and Smart Industries Group, Tencent Technology Co., Ltd., China

kunsong.npu.se@gmail.com, lxie@nwpu.edu.cn

Abstract

Speaker adaptation in text-to-speech synthesis (TTS) is to fine-tune a pre-trained TTS model to adapt to new target speakers with limited data. While much effort has been conducted towards this task, seldom work has been performed for low computational resource scenarios due to the challenges raised by the requirement of the lightweight model and less computational complexity. In this paper, a tiny VITS-based [1] TTS model, named AdaVITS, for low computing resource speaker adaptation is proposed. To effectively reduce the parameters and computational complexity of VITS, an inverse short-time Fourier transform (iSTFT)-based wave construction decoder is proposed to replace the upsampling-based decoder which is resource-consuming in the original VITS. Besides, NanoFlow is introduced to share the density estimate across flow blocks to reduce the parameters of the prior encoder. Furthermore, to reduce the computational complexity of the textual encoder, scaled-dot attention is replaced with linear attention. To deal with the instability caused by the simplified model, we use phonetic posteriorgram (PPG) as a frame-level linguistic feature for supervising the model process from phoneme to spectrum. Experiments show that AdaVITS can generate stable and natural speech in speaker adaptation with 8.97M model parameters and 0.72 GFlops computational complexity.¹

Index Terms: speaker adaptation, low computing resource, adversarial learning, normalizing flows

1. Introduction

To create human-like natural speech, modern neural network-based text-to-speech (TTS) models are usually large and contain mass of parameters [2, 3, 4, 5]. To train such a model, sufficient data and computational resources are necessary. However, to realize customized TTS, a corpus with sufficient samples recorded by a new target speaker is not always available in practice, which makes the few/one/zero-shot methods gain much interest. As compared with one/zero-shot methods [6, 7] which usually rely on an extra speaker embedding module, *speaker adaptation*, i.e. fine-tuning a well-trained base model with limited data to adapt to the new speaker, is still a practical approach with better speaker similarity. Considering the limited computational resources in many real-world scenarios, such as personalized voice services on edge devices, a lightweight TTS model with a small model size and low computation consumption is essential for the speaker adaptation task.

Due to the significant role of speaker adaptation in many scenarios, i.e., virtual avatars and personal assistants,

much effort has been conducted in this field. For example, AdaSpeech [8] reduces the number of adaptive parameters to alleviate the memory usage and serving cost by using conditional layer normalization. Besides, some approaches aim to reduce the model training time in the adaptation process to get a better user experience. For instance, Meta-Voice [9] uses meta-learning to obtain better model initialization for faster adaptation. Furthermore, some studies address the noise-robust speaker adaptation problem, aiming to obtain a noise-invariant TTS model for the target speaker with only noisy samples at hand [10].

While these efforts have successfully improved the performance of speaker adaptation models with limited data, reducing the TTS model size and computation complexity is still desired because the current popular TTS models are still too large. For instance, a typical Fastspeech 2 [11] model has 28M parameters, while the use of a neural vocoder adds extra. It is non-trivial to obtain a lightweight solution for speaker adaptation with decent performance. First, effectively reducing the parameters and computational complexity is an intuitive challenge. Second, the simplified model structure could arise instability in the generation of speech, resulting in low speaker similarity and naturalness with obvious artifacts and even pronunciation errors. Model compression via distillation and quantization [12] can be directly adopted but it may induct apparent performance loss. Recent effort on neural architecture search (NAS) [13] is another promising solution while the search process itself consumes much computation power and time effort.

To face the above challenges, in this paper we propose a tiny TTS model for low computing resource speaker adaptation. Inspired by the superiority of VITS [1], which is a fully end-to-end TTS model, on eliminating the mismatch between acoustic feature generation and wave construction in typical two-stage based methods, our lightweight solution, named *AdaVITS*, is built upon VITS with substantial modification to fit the speaker adaption scenario with fewer parameters, lower computational complexity, and stable performance. First, considering the resource-consuming characteristic of the upsampling-based decoder, an inverse short-time Fourier transform (iSTFT)-based wave construction decoder is proposed. Besides, to reduce the parameters of the prior encoder, flow indication embedding (FLE) is utilized to share the density estimate across flow blocks. Moreover, for the FFT blocks, scaled-dot attention is replaced with linear attention to reduce computational complexity. To deal with the instability caused by the simplified model, phonetic posteriorgram (PPG) is used as a frame-level linguistics feature to constrain the phoneme to spectrum modeling process. Extensive experiments demonstrate the good performance of AdaVITS on the speaker adaptation task with only 8.97M model parameters and 0.72 GFlops computation.

¹Audio samples are available at <https://AdaVITS.github.io/AdaVITS/>

*Lei Xie is the Corresponding author.

2. Method

VITS [1] is an end-to-end model of state-of-the-art, which uses variational autoencoder to learn latent variable as an intermediate representation between acoustic model and vocoder in end-to-end learning. In order to make the prior distribution close to the latent variable z , VITS added normalizing flow to the prior network to improve the representation ability of the prior distribution. Inspired by VITS, in this paper, we use the VITS approach to model the process from text features to speech waveforms, with the difference of using frame-level text features PPG as the intermediate constraint between phoneme and z . The advantage of using PPG is that it can explicitly decouple the timbre and content information of speech, which will make modeling more flexible [14]. In addition, in speaker adaptation, we can only learn the target speaker’s timbre rather than specific speech characteristics, for many people do not have reasonable control of speed and prosody when recording. As illustrated in Figure 1, similar to VITS, the proposed AdaVITS is composed of a posterior encoder, a prior encoder, and a decoder. The posterior encoder is used to extract latent variable z from the waveform when training and is not used in inference. The prior encoder is used to extract the prior distribution $p(z|c)$ of z from the phoneme, and the decoder is performed to generate waveform by the z and speaker embedding.

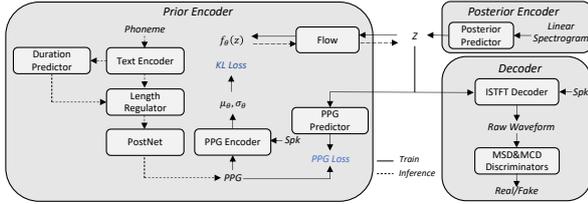


Figure 1: Architecture of AdaVITS.

2.1. Posterior Encoder

In AdaVITS, the posterior encoder is similar to that in the original VITS, which takes the linear spectrum as input to extract the mean and variance of the posterior distribution $p(z|y)$, and obtains the latent variable z . Due to the already existing speaker information in the linear spectrum, no extra speaker embedding is attached.

2.2. Prior Encoder

Conditioned on the conditional information c , including phoneme and speaker embedding, the prior encoder is to get the prior distribution $p(z|c)$ of the conditional variational autoencoder (CVAE). Compared with the original VITS, we use PPG as an intermediate constraint from phoneme to z . PPG is extracted from the acoustic model of a speaker-independent automatic speech recognition system, which is a frame-level linguistic feature not contain speaker information. Compared with spectrums, which contain not only linguistic information but also rich acoustic information, less information conveyed by PPG allows us to greatly simplify the complexity of the model and decouple the content and speaker information.

For text processing, fewer FFT layers are used for the text encoder, and then a length regulator is used to expand features from phoneme level to frame level, the construction of PPG is directly performed by the post-net rather than the structure with decoder and post-net. Since the computational complexity of scaled-dot attention in FFT is not linear with the sequence

length n , it has incredibly high computational complexity in long sentences. Here, referring to [15], linear attention is used to replace scaled-dot attention in FFT blocks, which will ensure the attention effect while reducing computational complexity. The modeling process from phoneme to PPG is pre-trained and not finetuned in speaker adaptation.

After obtaining PPG, the PPG encoder is used to get a prior normal distribution with mean μ_θ and variance σ_θ from PPG and speaker embedding. To be specific, the PPG encoder is composed of FFT blocks, in which linear attention is utilized here to reduce the computational complexity.

Due to the lack of explicit constraint of the pronunciation to z , the model tends to raise pronunciation errors such as mispronunciation and abnormal tone. To face this issue, a PPG predictor is introduced to provide the pronunciation constraint. The architecture of the PPG predictor is consistent with the phoneme predictor of VISinger [16]. With the input of z , the PPG produced by the PPG predictor is used to obtain the following constraint loss:

$$L_{\text{ppg}} = \left\| \hat{\text{PPG}} - \text{PPG} \right\|_1 \quad (1)$$

where PPG is the input. The PPG predictor is only trained on the pre-trained model and will be frozen during the adaptation.

The distribution z is then transformed into a more complex distribution using the normalized flow f_θ . This normalized flow includes multiple layers of affine coupling, and each layer consists of a stack of WaveNet [17] residual blocks following VITS. Due to the use of multiple layers, the flow has a large number of parameters. Referring to NanoFlow [18], we share the parameter in each affine coupling layer of flows, and each layer is distinguished by FLE. The amount of parameters in flow is controlled to be a single layer through this method. Similar to VITS, latent variable z is transformed to $f(z)$ by flow during training. During inference, the output of the PPG encoder is transformed into a latent variable \hat{z} by the inverse flow. As the speaker embedding is added to the input PPG, no extra speaker embedding will be added to the flow.

$$p(z|c) = N(f_\theta(z); \mu_\theta(c), \sigma_\theta(c)) \left| \det \frac{\partial f_\theta(z)}{\partial z} \right|. \quad (2)$$

2.3. Decoder

In VITS, the reconstruction of waveform is performed by a typical vocoder-like decoder, which consists of a series of upsampling layers. While this upsampling layers-based decoder generally has strong modeling capabilities, the gradual increasing process to transfer the input to the time domain is computation consuming. Because our model is equivalent to the joint training of the acoustic model and the vocoder, it is feasible to use iSTFT to generate the waveform directly. In practice, the real and imaginary parts of the waveform are predicted based on the features in the frequency domain, which can effectively reduce the computational cost.

As illustrated in Figure 2(a), we propose decoder-v1. We use multiple convolutions to gradually increase the input dimension to $(f/2 + 1) * 2$ to make the output fit the total dimension of real and imaginary parts, where f indicates the fast Fourier transform size. A stack of residual blocks follow each one-dimensional convolution for more information on the corresponding scale. Due to the frequency domain dimension modeling, we do not use dilated convolution but use a smaller kernel size with the aim to ensure that the receptive field will not be too large. The group convolution is used in one-dimensional convolution to save computation. Then, the output is split into real and imaginary parts, based on which the final waveform can be

produced via iSTFT. Note that, following VITS, the input condition includes speaker embedding and latent variables z , as we found that the speaker similarity will be degraded significantly if the speaker embedding is not added to the decoder.

As illustrated in Figure 2(b), to accommodate the computing resource requirements of different scenarios, we also provide an alternative v2 version for a trade-off between computation complexity and sound quality. In decoder-v2, we only use iSTFT-based decoder-v1 to model the high-frequency part while use the upsampling layer with a residual network in the GAN-based vocoder to model the low-frequency part. Because the upsampling method can synthesize high-quality harmonics and the high-frequency part requires less modeling capability, iSTFT is sufficient to meet its modeling requirements. Then, we take the signals generated by the upsampling network as low-frequency bands, and the signals generated by the decoder-v1 part as high-frequency bands, and adopt the pseudo quadrature Mirror filter-bank (PQMF) for subband modeling.

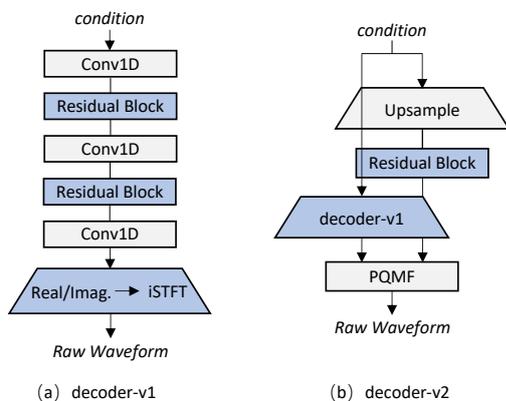


Figure 2: Architecture of decoder.

2.4. Discriminator

A multi-resolution spectrum discriminator (MSD) [19] and a multi-resolution complex-valued spectrum discriminator (MCD) are adopted for the adversarial training. The models frequency domain information at different levels, and the discriminator is particularly effective in the iSTFT decoder and has noticeable gains in high-frequency harmonic reconstruction. In addition to MSD, the proposed MCD is to model the relationship between the real and imaginary parts of the signal, which is useful for improving the phase accuracy. To be specific, MCD divides the signal into real and imaginary parts by short-time Fourier transform (STFT) at multiple scales and then applies 2-d complex convolutions to the input, which has been shown to work well in the complex-valued domain [20]. Its architecture is consistent with the multi-resolution spectrum discriminator.

2.5. Loss

The training of AdaVITS includes CVAE and GAN training. The CVAE loss is expressed as

$$L_{cvae} = L_{kl} + \lambda_{recon} * L_{recon} + \lambda_{ppg} * L_{ppg} \quad (3)$$

where L_{kl} is Kullback-Leibler divergence following VITS, and L_{recon} calculates the L1 distance between the mel-spectrum of the waveform generated by the decoder and the ground truth. λ_{recon} and λ_{ppg} are 45 and 10, respectively. With GAN training, the final objectives are expressed as

$$L_G = L_{adv}(G) + \lambda_{fm} * L_{fm}(G) + L_{cvae} \quad (4)$$

$$L_D = L_{adv}(D) \quad (5)$$

where $L_{adv}(G)$ and $L_{adv}(D)$ are the GAN loss of G and D, the feature matching loss L_{fm} is used to improve the training stability, and the λ_{fm} is 2. Since the discriminator consists of multiple sub-discriminators in the MSD and MCD, the above GAN loss and feature matching loss are the sum of the losses of multiple sub-discriminators.

3. Experiments

3.1. Datasets

We use *train-clean360* and *train-clean100* subsets of LibriTTS [21] for pre-training model, which contains around 242 hours of utterances from 1151 speakers. To evaluate the performance of AdaVITS in the speaker adaption task, VCTK [22], which is another commonly used multi-speaker TTS corpus with different acoustic conditions from LibriTTS, is adopted to fine-tune the pre-trained model. In practice, five males and five females are randomly selected from VCTK to work as the target speakers for speaker adaptation. For each speaker, 20 utterances are randomly selected. Another 10 extra sentences from each speaker are randomly selected, resulting in a testing set with a total of 100 sentences from 10 speakers.

All audio samples are downsampled to 16kHz, and are then represented as frame level with 12.5ms hop length and 50ms window length. We use the bottleneck feature extracted from the pre-trained WeNet [23] model as PPG with dimension 256.

3.2. Model Configuration

To evaluate the performance of AdaVITS, some representative models, including *Fastspeech 2 with HiFiGAN* and *VITS* are compared in the experiments. In practice, two *Fastspeech 2 with HiFiGAN* systems are compared, which are referred to as *Fs2-o+HiFiGAN v1* and *Fs2-l+HiFiGAN v2*, respectively. *Fs2-o* is a standard *Fastspeech 2* model, which follows the basic architecture in *Fastspeech 2* [11]. The difference is that the duration predictor and pitch predictor use 5 Conv1D layers with kernel size 5 for more prediction accuracy. Compared with the original *Fastspeech 2*, the energy predictor is not used in *Fs2-o*. *Fs2-l* is a lightweight version with reduced filter size and layers. Compared with *Fs2-o*, we use two FFT layers in both encoder and decoder of *Fs2-l* and set the filter size and hidden dimension to 768 and 128. As for the vocoder, compared with *HiFiGANv1*, *HiFiGANv2* is a lightweight version, and the details can be found from [5]. This *Fs2-l+HiFiGANv2* has similar model parameters and computational complexity with the proposed AdaVITS, designed for comparison purpose.

In the AdaVITS, *AdaVITS-v1* means to use *decoder-v1*, *AdaVITS-v2* means to use *decoder-v2* described in Section 2.2.2. The duration predictor of AdaVITS and all FFT blocks follow the configuration of *Fs2-l*. All encoders in the proposed approach consists of 2 FFT blocks and the post-net follows [3]. In the decoder, for the *decoder-v1*, conv1D channels are [256, 384, 1026], and all kernels are set to 3; and for the *decoder-v2*, conv1D channels are [256, 384, 774] while the upsample rates are [5, 5, 2] and upsample hidden channels are [256, 192, 64]. Through this method, *decoder-v2* upsampling layers model low frequency from 0 to 4Khz, and iSTFT modeling high frequency from 4 to 16Khz. We follow [24] for all upsampling layers and residual network structures. The MCD and MSD follow the architecture of MSD in *Glow-WaveGAN*[25]. Other settings are the same as the original VITS[1].

In addition to AdaVITS, a variation referred to as *AdaVITS-e* is also compared. In *AdaVITS-e*, the model is

Table 1: Experimental results in terms of MOS and WER. Model parameters and computational complexity are also shown.

Model	Params (M)	Com. (GFlops)	Naturalness	Similarity	WER (%)
Fs2-o+HiFiGAN v1	40.16	15.85	3.08 (± 0.13)	3.21 (± 0.10)	8.90
FS2-l+HiFiGAN v2	8.67	0.98	2.63 (± 0.11)	3.08 (± 0.14)	10.53
VITS [1]	29.36	15.76	3.59 (± 0.13)	3.53 (± 0.12)	15.29
AdaVITS-e	8.70	0.66	2.82 (± 0.15)	3.16 (± 0.16)	11.11
AdaVITS-v1	8.97	0.72	2.94 (± 0.14)	3.10 (± 0.14)	8.19
AdaVITS-v2	11.55	3.63	3.15 (± 0.13)	3.12 (± 0.12)	8.17
Recording	-	-	3.70 (± 0.12)	3.62 (± 0.11)	4.68

trained via an end-to-end way with text as input instead of using PPG as an intermediate constraint. The architecture of AdaVITS-e is similar to VISinger [16] but without the F0 predictor. The number of FFT block layer of the text encoder and frame prior network is set to 2, and other configurations in AdaVITS are applied to AdvaTTS-e.

In all the above models, the dimension of speaker embedding is set as 256. Fs2-o/Fs2-l/VITS/AdaVITS-e/AdaVITS pre-trained models and HiFiGAN v1/HiFiGAN v2 are trained up to 800k steps on 2080Ti GPU with batch size of 32. In the adaptation process, we finetune Fs2-o/Fs2-l/VITS/AdaVITS/AdaVITS-e on 2080Ti GPU for 2000 steps, and HiFiGAN v1/HiFiGAN v2 will not be updated further.

3.3. Experimental Results

To evaluate the performance of different models, a mean opinion score (MOS) test is conducted in terms of the naturalness and speaker similarity. A good synthesized sample should have high quality in naturalness and similarity with the target speaker. In this human rating test, each utterance is listened by 20 listeners, and the participants are asked to rate the sample with a score ranging from one to five for the naturalness and speaker similarity respectively. As for the evaluation of computational complexity, the GFlops required for generating speech per second is utilized. Since the computational complexity required by scaled-dot attention is not linear with the sentence length, the average of the test set is used as the result. In addition, the word error rate (WER) of each system is calculated to show the stability of each model, especially concerning pronunciation and intonation. WER is calculated by pre-trained WeNet model[†]. Note that this model is different from the model used to extract PPG. Results are shown in Table 1.

As can be seen from Table 1, compared with FS2-l+HiFiGAN v2 which has a similar model size with AdaVITS-v1, the proposed AdaVITS achieves better naturalness and less computational complexity. As for the WER of samples synthesized by AdaVITS is only 52.6% of that synthesized by FS2-l+HiFiGAN v2, indicating the good stability of AdvaVITS. Compared with FS2-o+HiFiGAN v1, AdaVITS-v2 has a similar naturalness but a smaller model size. Compared with the original VITS, AdaVITS still has a gap to bridge in terms of naturalness and speaker similarity. However, AdaVITS achieves much better WER compared with other methods, which is mainly attributed to the utilization of PPG-based linguistic features, which can be proved by the performance of AdaVITS-e, in which regular text is used as input. It should be noted that the higher MOS score in terms of naturalness has no necessary relationship with the WER, as the participants paid more attention

[†]<https://github.com/wenet-e2e/wenet/tree/main/examples/librispeech/>

to the prosody and quality of speech and our human beings have a higher tolerance for the pronunciation than an ASR model.

3.4. Ablation study

To evaluate the effectiveness of each component in the proposed method, an ablation study is conducted by dropping out each component respectively on AdaVITS-v1. To be specific, the effectiveness of linear attention, FLE, MCD, PPG predictor, and iSTFT decoder is analyzed, and the results can be found in Table 2. As can be seen, MCD and PPG predictor play important roles in obtaining high quality speech, while the addition of FLE, linear attention, and iSTFT decoder can effectively reduce the number of parameter or computation complexity without leading to a obvious impact on the results.

Table 2: Ablation study results. The MOS test are for the naturalness. w/o means without. In w/o Linear Att., the linear attention is replaced by the scaled-dot attention. In w/o ISTFT Dec, the ISTFT decoder is replaced by the decoder of HiFiGAN v2.

Model	Params	Com.	MOS
AdaVITS-v1	8.97	0.72	3.17 (± 0.12)
w/o Linear Att.	8.97	0.83	3.14 (± 0.13)
w/o FLE	11.88	0.72	3.18 (± 0.14)
w/o MCD	8.97	0.72	2.99 (± 0.14)
w/o PPG Predictor	8.97	0.72	2.83 (± 0.13)
w/o ISTFT Dec.	7.24	1.46	3.32 (± 0.13)
Recording	-	-	4.04 (± 0.10)

4. Conclusions

In this paper, a VITS-based lightweight adaptive TTS system, referred to as AdaVITS, is proposed to support speaker adaptation’s need for low cost. To effectively reduce the computational complexity led by the upsampling-based vocoder, an iSTFT-based wave construction decoder is proposed. In addition, NanoFlow is utilized to reduce the parameters of the prior encoder, and scaled-dot attention in FFT is replaced with linear attention to further reduce the computational complexity. To ensure the stability of the simplified model, PPG is used as frame-level linguistic features. Extensive experiments demonstrate the obvious superiority of AdaVITS in terms of the model size and computational complexity compared with other standard models. When compared to the model with similar parameters, the proposed AdaVITS achieves less computational complexity and better speech quality.

5. References

- [1] J. Kim, J. Kong, and J. Son, "Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech," in *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, ser. Proceedings of Machine Learning Research, M. Meila and T. Zhang, Eds., vol. 139. PMLR, 2021, pp. 5530–5540.
- [2] Y. Wang, R. J. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, Q. V. Le, Y. Agiomyriannakis, R. Clark, and R. A. Saurous, "Tacotron: Towards end-to-end speech synthesis," in *Interspeech 2017, 18th Annual Conference of the International Speech Communication Association, Stockholm, Sweden, August 20-24, 2017*, F. Lacerda, Ed. ISCA, 2017, pp. 4006–4010.
- [3] Y. Ren, Y. Ruan, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T. Liu, "Fastspeech: Fast, robust and controllable text to speech," in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. B. Fox, and R. Garnett, Eds., 2019, pp. 3165–3174.
- [4] K. Kumar, R. Kumar, T. de Boissiere, L. Gestin, W. Z. Teoh, J. Sotelo, A. de Brébisson, Y. Bengio, and A. C. Courville, "Melgan: Generative adversarial networks for conditional waveform synthesis," in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. B. Fox, and R. Garnett, Eds., 2019, pp. 14 881–14 892.
- [5] J. Su, Z. Jin, and A. Finkelstein, "Hifi-gan: High-fidelity denoising and dereverberation based on speech deep features in adversarial networks," in *Interspeech 2020, 21st Annual Conference of the International Speech Communication Association, Virtual Event, Shanghai, China, 25-29 October 2020*, H. Meng, B. Xu, and T. F. Zheng, Eds. ISCA, 2020, pp. 4506–4510.
- [6] S. Ö. Arik, J. Chen, K. Peng, W. Ping, and Y. Zhou, "Neural voice cloning with a few samples," in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, S. Bengio, H. M. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., 2018, pp. 10 040–10 050.
- [7] Y. Chen, Y. M. Assael, B. Shillingford, D. Budden, S. E. Reed, H. Zen, Q. Wang, L. C. Cobo, A. Trask, B. Laurie, Ç. Gülgehre, A. van den Oord, O. Vinyals, and N. de Freitas, "Sample efficient adaptive text-to-speech," in *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- [8] M. Chen, X. Tan, B. Li, Y. Liu, T. Qin, S. Zhao, and T. Liu, "Adaspeech: Adaptive text to speech for custom voice," *CoRR*, vol. abs/2103.00993, 2021.
- [9] S. Liu, D. Su, and D. Yu, "Meta-voice: Fast few-shot style transfer for expressive voice cloning using meta learning," *CoRR*, vol. abs/2111.07218, 2021.
- [10] J. Cong, S. Yang, L. Xie, G. Yu, and G. Wan, "Data efficient voice cloning from noisy samples with domain adversarial training," in *Interspeech 2020, 21st Annual Conference of the International Speech Communication Association, Virtual Event, Shanghai, China, 25-29 October 2020*, H. Meng, B. Xu, and T. F. Zheng, Eds. ISCA, 2020, pp. 811–815.
- [11] Y. Ren, C. Hu, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T. Liu, "Fastspeech 2: Fast and high-quality end-to-end text to speech," in *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [12] G. E. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *CoRR*, vol. abs/1503.02531, 2015.
- [13] R. Luo, X. Tan, R. Wang, T. Qin, J. Li, S. Zhao, E. Chen, and T. Liu, "Lightspeech: Lightweight and fast text to speech with neural architecture search," in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2021, Toronto, ON, Canada, June 6-11, 2021*. IEEE, 2021, pp. 5699–5703.
- [14] T. Wang, J. Tao, R. Fu, J. Yi, Z. Wen, and R. Zhong, "Spoken content and voice factorization for few-shot speaker adaptation," H. Meng, B. Xu, and T. F. Zheng, Eds. ISCA, 2020, pp. 796–800.
- [15] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret, "Transformers are rnns: Fast autoregressive transformers with linear attention," in *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 5156–5165.
- [16] Y. Zhang, J. Cong, H. Xue, L. Xie, P. Zhu, and M. Bi, "Visinger: Variational inference with adversarial learning for end-to-end singing voice synthesis," *CoRR*, vol. abs/2110.08813, 2021.
- [17] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. W. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," in *The 9th ISCA Speech Synthesis Workshop, Sunnyvale, CA, USA, 13-15 September 2016*. ISCA, 2016, p. 125.
- [18] S. Lee, S. Kim, and S. Yoon, "Nanoflow: Scalable normalizing flows with sublinear parameter complexity," in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., 2020.
- [19] W. Jang, D. Lim, and J. Yoon, "Universal melgan: A robust neural vocoder for high-fidelity waveform generation in multiple domains," *CoRR*, vol. abs/2011.09631, 2020.
- [20] Y. Hu, Y. Liu, S. Lv, M. Xing, S. Zhang, Y. Fu, J. Wu, B. Zhang, and L. Xie, "DCCRN: deep complex convolution recurrent network for phase-aware speech enhancement," in *Interspeech 2020, 21st Annual Conference of the International Speech Communication Association, Virtual Event, Shanghai, China, 25-29 October 2020*, H. Meng, B. Xu, and T. F. Zheng, Eds. ISCA, 2020, pp. 2472–2476.
- [21] H. Zen, V. Dang, R. Clark, Y. Zhang, R. J. Weiss, Y. Jia, Z. Chen, and Y. Wu, "Libritts: A corpus derived from librispeech for text-to-speech," in *Interspeech 2019, 20th Annual Conference of the International Speech Communication Association, Graz, Austria, 15-19 September 2019*, G. Kubin and Z. Kacic, Eds. ISCA, 2019, pp. 1526–1530.
- [22] C. Veaux, J. Yamagishi, K. MacDonald *et al.*, "Superseded-cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit," 2017.
- [23] Z. Yao, D. Wu, X. Wang, B. Zhang, F. Yu, C. Yang, Z. Peng, X. Chen, L. Xie, and X. Lei, "Wenet: Production oriented streaming and non-streaming end-to-end speech recognition toolkit," in *Interspeech 2021, 22nd Annual Conference of the International Speech Communication Association, Brno, Czechia, 30 August - 3 September 2021*, H. Hermansky, H. Cernocký, L. Burget, L. Lamel, O. Scharenborg, and P. Motlíček, Eds. ISCA, 2021, pp. 4054–4058.
- [24] G. Yang, S. Yang, K. Liu, P. Fang, W. Chen, and L. Xie, "Multi-band melgan: Faster waveform generation for high-quality text-to-speech," in *2021 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2021, pp. 492–498.
- [25] J. Cong, S. Yang, L. Xie, and D. Su, "Glow-wavegan: Learning speech representations from gan-based variational auto-encoder for high fidelity flow-based speech synthesis," in *Interspeech 2021, 22nd Annual Conference of the International Speech Communication Association, Brno, Czechia, 30 August - 3 September 2021*, H. Hermansky, H. Cernocký, L. Burget, L. Lamel, O. Scharenborg, and P. Motlíček, Eds. ISCA, 2021, pp. 2182–2186.