

PARALLEL WAVEFORM SYNTHESIS BASED ON GENERATIVE ADVERSARIAL NETWORKS WITH VOICING-AWARE CONDITIONAL DISCRIMINATORS

Ryuichi Yamamoto¹, Eunwoo Song², Min-Jae Hwang³ and Jae-Min Kim²

¹LINE Corp., Tokyo, Japan
²NAVER Corp., Seongnam, Korea
³Search Solutions Inc., Seongnam, Korea

ABSTRACT

This paper proposes voicing-aware conditional discriminators for Parallel WaveGAN-based waveform synthesis systems. In this framework, we adopt a projection-based conditioning method that can significantly improve the discriminator’s performance. Furthermore, the conventional discriminator is separated into two waveform discriminators for modeling voiced and unvoiced speech. As each discriminator learns the distinctive characteristics of the harmonic and noise components, respectively, the adversarial training process becomes more efficient, allowing the generator to produce more realistic speech waveforms. Subjective test results demonstrate the superiority of the proposed method over the conventional Parallel WaveGAN and WaveNet systems. In particular, our speaker-independently trained model within a FastSpeech 2 based text-to-speech framework achieves the mean opinion scores of 4.20, 4.18, 4.21, and 4.31 for four Japanese speakers, respectively.

Index Terms— Text-to-speech, neural vocoder, generative adversarial networks, waveform synthesis

1. INTRODUCTION

Deep generative models in text-to-speech (TTS) frameworks have significantly improved the perceptual quality of synthetic speech signals [1, 2]. In particular, the autoregressive generative models, such as WaveNet, have shown superior quality over conventional parametric vocoders [3–7]. However, they suffer from slow generation due to their autoregressive nature and thus are limited in their applications to real-time scenarios.

To achieve real-time TTS systems, non-autoregressive waveform synthesis models have been proposed based on teacher-student frameworks [8,9], normalizing flows [10,11], or generative adversarial networks (GANs) [12, 13]. Specifically, in our previous work, we proposed the *Parallel WaveGAN* methods [14, 15], characterized by efficient training and fast inference while maintaining a quality that is competitive to the state-of-the-art Parallel WaveNet. However, as it is insufficient for a single discriminator to distinguish the complex nature of speech signal — e.g., voiced and unvoiced characteristics — the generated speech often suffers from unnatural artifacts. In addition, this problem becomes more severe when the training database has more diversity such as in the scenario of speaker-independent modeling.

To address the aforementioned problems, we propose voicing-aware conditional discriminators for Parallel WaveGAN. In this method, we adopt a projection-based conditioning framework that

incorporates acoustic features into the discriminators [16]. This enables the discriminator to classify the input speech well to be consistent with the given acoustic features. Furthermore, we introduce two separate voicing-aware discriminators that individually model the voiced and unvoiced speech, respectively. In detail, one discriminator is designed to have long receptive fields for capturing slowly varying harmonic components, which mainly represents the voiced speech; whereas the other has small receptive fields for capturing rapidly varying noise components of unvoiced speech. Because each discriminator learns the distinctive characteristics of the harmonic and noise components, respectively, the adversarial training process becomes more effective.

We investigate the performance of our proposed method by conducting perceptual listening tests in a TTS framework. Specifically, a speaker-independently trained Parallel WaveGAN with the FastSpeech 2 acoustic model significantly outperforms the conventional Parallel WaveGAN and similarly configured WaveNet systems, achieving mean opinion scores of 4.20, 4.18, 4.21, and 4.31 for four Japanese speakers, respectively.

2. RELATED WORK

There have been several attempts to improve the discriminator’s performance for GAN-based neural waveform synthesis systems. For instance, MelGAN [13] and VocGAN [17] employ multi-scale discriminators to learn the waveform structure on different time scales. GAN-TTS [18] adopts a blend of multiple conditional and unconditional discriminators based on multi-frequency random windows. Although these multi-resolution architectures are found to be effective for high perceptual quality, their methods tend to require complicated discriminators (e.g., hierarchically-nested joint conditional and unconditional discriminators [17]), which are more difficult to train. To keep the discriminator simple yet effective, our proposed method adopts two separate discriminators, which can explicitly focus on the distinctive voiced and unvoiced characteristics of speech.

3. METHOD

3.1. Parallel WaveGAN

Parallel WaveGAN is a non-autoregressive WaveNet model that generates a time-domain speech waveform from the corresponding conditional acoustic parameters [14]. Specifically, the conventional Parallel WaveGAN consists of a non-causal WaveNet generator, G , and a single convolutional neural network (CNN) dis-

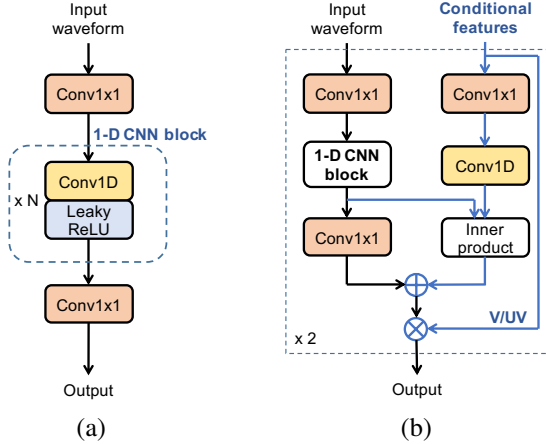


Fig. 1. Block diagram of (a) a conventional and (b) the proposed discriminators. Note that in the proposed method, two separate discriminators with different dilation factors of the 1-D convolutional neural network (CNN) blocks were used for modeling the voiced and unvoiced segments, respectively.

criminator, D . Based on GANs [19], the generator learns a distribution of realistic waveforms by trying to deceive the discriminator into recognizing the generated samples as *real*. Moreover, the discriminator is trained to correctly classify the generated sample as *fake* while classifying the ground truth as *real*. By combining adversarial training with an auxiliary multi-resolution short-time Fourier transform (STFT) loss function, Parallel WaveGAN learns the time-frequency characteristics of realistic speech efficiently.

3.2. Proposed Parallel WaveGAN with voicing-aware conditional discriminators

Fig. 1 depicts an overview of the Parallel WaveGAN’s discriminator. Compared with the conventional method (Fig. 1(a)), there are two main improvements in the proposed method (Fig. 1(b)), as follows. First, we adopt a projection-based conditioning method where the acoustic features are incorporated into the discriminator as conditional inputs [16]. This helps the discriminator to better classify the input signals to be consistent with the given conditional features. Second, we replace the conventional discriminator with voicing-aware ones using a voiced and unvoiced binary flag (V/UV). Considering the fact that the voiced and unvoiced segments of speech signals have distinctive characteristics, the two separate discriminators independently operate to capture each segment, respectively. Note that voicing masks are used to make each discriminator see only the region of its interest.

The voiced segment can be characterized by slowly evolving harmonic components. To control these components, we design the first discriminator with a dilated CNN [3]. Note that the use of dilated convolution allows the discriminator to increase the size of the receptive field while keeping a small number of parameters. With a sufficient size of the receptive field, the discriminator not only covers long-term variations of the harmonic component, but also penalizes any unwanted aperiodic noise components in the voiced regions. On the other hand, the second discriminator is composed of a non-dilated CNN with a small receptive field (i.e., a dilation factor of 1 in the 1-D CNN block). Because the char-

acteristics of the noise component vary rapidly, employing a short window is advantageous to focus on the detailed high-frequency structure of speech.

3.3. Training objectives

To train the proposed models, we adopt the least-squares GANs thanks to their training stability [20]. The training objectives for the discriminators and the generator are defined as follows:

$$\min_D \mathbb{E}_{\mathbf{z}, \mathbf{h}} [(1 - D(\mathbf{x}, \mathbf{h}))^2] + \mathbb{E}_{\mathbf{z}, \mathbf{h}} [D(G(\mathbf{z}, \mathbf{h}), \mathbf{h})^2], \forall D \in \{D^v, D^{uv}\} \quad (1)$$

$$\min_G \mathbb{E}_{\mathbf{x}, \mathbf{z}, \mathbf{h}} [L_{\text{mr_stft}}(\mathbf{x}, G(\mathbf{z}, \mathbf{h}))] + \frac{1}{2} \lambda_{\text{adv}} \mathbb{E}_{\mathbf{z}, \mathbf{h}} \left[\sum_{D \in \{D^v, D^{uv}\}} (1 - D(G(\mathbf{z}, \mathbf{h}), \mathbf{h}))^2 \right], \quad (2)$$

where \mathbf{z} , \mathbf{h} , and \mathbf{x} denote the Gaussian noise, conditional acoustic features, and the target speech waveform, respectively; D^v and D^{uv} are the voiced and unvoiced discriminators, respectively; and λ_{adv} represents a hyperparameter that balances the two adversarial losses and the multi-resolution STFT loss defined as follows:

$$L_{\text{mr_stft}}(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{M} \sum_{m=1}^M L_{\text{stft}}^{(m)}(\mathbf{x}, \hat{\mathbf{x}}), \quad (3)$$

where $\hat{\mathbf{x}}$ represents the generated waveform; and $L_{\text{stft}}^{(m)}(\mathbf{x}, \hat{\mathbf{x}})$ denotes the m^{th} STFT loss, represented by the sum of *spectral convergence* and *log STFT magnitude losses*¹ [21]. Note that unlike the original Parallel WaveGAN [14], the discriminator is now divided into distinctive voiced and unvoiced parts, and the generator is designed to deceive both of them.

4. EXPERIMENTS

4.1. Experimental setup

4.1.1. Data and feature configurations

The experiments used four phonetically and prosodically rich speech corpora recorded by two female (F1, F2) and two male (M1, M2) Japanese professional speakers. The speech signals were sampled at 24 kHz, and each sample was quantized by 16 bits. Each corpus included 5,000 utterances, among which 4,500, 250, and 250 samples were used for training, validation, and evaluation, respectively. The training data size for each speaker was between 5.5 and 5.9 hours.

The acoustic features were extracted using an improved time-frequency trajectory excitation vocoder at the analysis intervals of 5 ms [22] and included 40-dimensional line spectral frequencies, the fundamental frequency, the energy, the binary V/UV flag, a 32-dimensional slowly evolving waveform, and a 4-dimensional rapidly evolving waveform, all of which constituted a 79-dimensional feature vector². The acoustic features were then

¹The detailed setups for designing the multi-resolution STFT loss are the same as those in the original Parallel WaveGAN [14].

²We have also tried mel-spectrograms as acoustic features but found that the vocoder parameters were more effective to avoid buzzy synthetic speech.

Table 1. The dilation factors and receptive fields in the 1-D CNN blocks of the voicing-aware discriminators.

Discriminator	Dilation factors	Receptive field
D^v	[1, 2, 4, 8, 16, 32]	127
D^{uv}	[1, 1, 1, 1, 1, 1]	13

normalized to have zero mean and unit variance using the statistics of the training data.

4.1.2. Model details

The proposed Parallel WaveGAN consists of a WaveNet-based generator and voiced and unvoiced discriminators. The generator comprises 30 layers of dilated residual 1-D convolution blocks with three exponentially increasing dilation cycles [14]. The number of residual and skip channels was set to 64, and the convolution filter size was set to 5. The size of the receptive field for the generator was 12,277. The discriminators for the voiced and unvoiced regions were each composed of a 1-D CNN block and 1-by-1 convolution layers. Each 1-D CNN block contains six convolution layers interleaved with leaky ReLU activation. The number of channels and kernel size in the 1-D CNN blocks were set to 64 and 3, respectively. The dilation factors and receptive fields³ of the 1-D CNN blocks for the voicing-aware discriminator are summarized in Table 1. For conditional input, a 1-D convolution with the kernel size of the discriminator’s receptive field was used before the inner product projection. The number of channels was 64, the same as in the 1-D CNN block.

At the training stage, the multi-resolution STFT loss was computed by the sum of three different STFT losses, as described in Parallel WaveGAN [14]. The hyperparameter λ_{adv} in equation (2) was chosen to be 4.0. The models were trained for 400K steps with a RAdam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 1e^{-6}$ [23]. The discriminators were fixed for the first 100K steps, and the models were jointly trained afterwards. The mini-batch size was set to eight, and the length of each audio clip was set to 24K time samples (1.0 second). The initial learning rate was set to 0.0001 for both the generator and discriminators. The learning rate was reduced by half for every 200K steps. We trained Parallel WaveGAN models with two NVIDIA Tesla V100 GPUs, which took about 44 and 58 hours for conventional and proposed Parallel WaveGAN systems, respectively.

To validate our discriminator design choices, we investigated six Parallel WaveGAN systems with different discriminator configurations, as described in Table 2. Note that all the Parallel WaveGAN systems used the same generator architecture and training configurations as described above; they only differed in the discriminator settings. Therefore, the proposed method retained the original Parallel WaveGAN’s fast inference speed.

As a baseline system, we used the autoregressive Gaussian WaveNet [9], which consists of 24 layers of dilated residual convolution blocks with exponentially increasing four dilation cycles. The number of residual and skip channels was set to 128, and the filter size was set to 3. The model was trained for 1M steps with a RAdam optimizer. The learning rate was set to 0.001 and reduced by half for every 200K steps. The mini-batch size was set to eight,

³The receptive field for the discriminator of the voiced regions was kept not too large because a larger receptive field poses training difficulty.

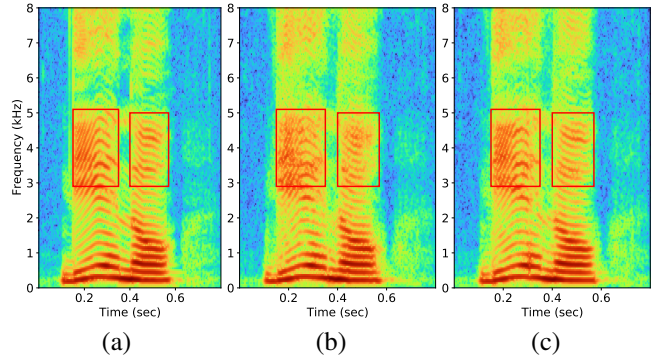


Fig. 2. Spectrograms of (a) natural speech, (b) generated speech from the conventional Parallel WaveGAN (S2), and (c) generated speech from the proposed Parallel WaveGAN (S7). As demonstrated in rectangle areas, our proposed method is able to model spectral harmonics more accurately.

and the length of each audio clip was set to 12 K time samples (0.5 seconds). The log-scale parameters of Gaussian were clipped at -9.0 during training to reduce noisy artifacts [9].

Across all the neural vocoders, the input auxiliary features were up-sampled by nearest neighbor interpolation followed by 1-D convolutions so that the time resolution of the auxiliary features matched the sampling rate of the speech waveforms [12, 24]. For the models using conditional discriminators, the up-sampled features were used as the conditional input. Note that the binary V/UV flag used in the voicing-aware discriminator was up-sampled to the sample level by repetition as an exception. All the vocoder models were trained in a speaker-independent manner by putting all the speaker’s data together.

4.2. Evaluation

We performed mean opinion score (MOS)⁴ tests to investigate the effectiveness of our proposed method. Seventeen native Japanese speakers were asked to make quality judgments about the synthesized speech samples using the following five possible responses: 1 = Bad; 2 = Poor; 3 = Fair; 4 = Good; and 5 = Excellent. In total, 20 utterances were randomly selected from the evaluation set and were then synthesized using the different models.

Table 2 shows the MOS test results with respect to different neural vocoders. The findings can be analyzed as follows. (1) Among all speakers, the Parallel WaveGAN system using the conditional discriminator (S3) obtained a better score than the baseline WaveNet (S1) and Parallel WaveGAN (S2) systems. The results confirmed the effectiveness of incorporating acoustic features into the discriminator through conditional information. (2) The systems using poorly configured voicing-aware discriminators (S4, S5, and S6) performed worse than S3. More specifically, they performed even worse than the baseline Parallel WaveGAN (S2) in some cases (e.g., comparing S2 and S5). Notably, the synthetic male voices contained buzzy noise, which resulted in significantly lower scores than the baseline Parallel WaveGAN. (3) Finally, the system with the proposed voicing-aware discriminators (S7) obtained the best scores and consistently outperformed the

⁴Audio samples are available at the following URL: <https://r9y9.github.io/demos/projects/icassp2021/>

Table 2. MOS test results with 95% confidence intervals in analysis/synthesis: The speech samples were generated using the acoustic features extracted from the recorded speech. PWG denotes Parallel WaveGAN for short. Note that systems S2 and S3 used D^v as the primary discriminator. All the models were trained in a speaker-independent manner.

System	Model	Voiced segments	Unvoiced segments	Discriminator conditioning	MOS			
					F1	F2	M1	M2
S1	WaveNet	-	-	-	3.64±0.12	3.83±0.11	3.33±0.12	3.13±0.11
S2	PWG	-	-	-	3.61±0.11	3.55±0.11	3.57±0.12	3.61±0.11
S3	PWG-cGAN-D	-	-	Yes	4.04±0.10	3.95±0.10	3.91±0.11	3.97±0.10
S4	PWG-V/UV-D	D^v	D^v	Yes	3.60±0.12	3.59±0.11	3.34±0.11	3.48±0.11
S5	PWG-V/UV-D	D^{uv}	D^v	Yes	3.67±0.11	3.48±0.11	3.29±0.12	3.38±0.11
S6	PWG-V/UV-D	D^{uv}	D^{uv}	Yes	3.77±0.11	3.88±0.10	3.57±0.11	3.34±0.11
S7	PWG-V/UV-D (proposed)	D^v	D^{uv}	Yes	4.11±0.10	4.05±0.10	4.04±0.10	4.08±0.10
R1	Recordings	-	-	-	4.63±0.08	4.67±0.07	4.61±0.08	4.64±0.08

Table 3. MOS test results with 95% confidence intervals: Acoustic features generated from the FastSpeech 2 acoustic model were used to compose the input auxiliary features.

System	Model	MOS			
		F1	F2	M1	M2
S1	FastSpeech 2 + WaveNet	3.90±0.11	3.81±0.10	3.43±0.11	3.09±0.10
S2	FastSpeech 2 + PWG	3.76±0.11	3.62±0.11	3.63±0.11	3.78±0.10
S3	FastSpeech 2 + PWG-cGAN-D	4.02±0.10	4.03±0.10	4.16±0.10	4.06±0.10
S7	FastSpeech 2 + PWG-V/UV-D (proposed)	4.20±0.10	4.18±0.09	4.21±0.09	4.31±0.09
R1	Recordings	4.63±0.08	4.67±0.07	4.61±0.08	4.64±0.08

other systems (from S1 to S6). The results proved the importance of the discriminator design and the effectiveness of our proposed approach. The benefits of our method can also be confirmed in spectrogram visualization; as shown in Fig.2, the proposed method was able to better reconstruct the spectral harmonics.

4.3. Text-to-speech

To further verify the effectiveness of the proposed method within the TTS framework, we combined the proposed Parallel WaveGAN with a FastSpeech 2 based acoustic model [25]. This model was configured based on the setup of FastSpeech 2, except for a few modifications as follows: we changed the variance predictor module to operate on phoneme-level rather than frame-level [26]. Manually annotated phoneme alignment was used instead of performing forced alignment. The model used accent information as an external input to better model pitch accents of Japanese [27].

For evaluation, we trained four speaker-dependent acoustic models. At the training stage, a dynamic batch size with an average of 24 samples was used for making a mini-batch [28, 29], and the models were trained for 200K iterations. In the synthesis step, the input phoneme and accent sequences were converted to the corresponding acoustic parameters by the FastSpeech 2 model. By inputting the resulting acoustic parameters, the vocoder models generated the time-domain speech signals.

To evaluate the quality of the generated speech samples, we performed naturalness MOS tests. The test setups were the same as those described in section 4.2, except that we excluded the systems with poorly designed discriminators (S4, S5, and S6 in Table 2, respectively)

The results of the MOS tests are shown in Table 3. Similar to the analysis/synthesis results, the proposed system with

voicing-aware conditional discriminators (S7) achieved the best scores among all speakers. In particular, the proposed method significantly outperformed the baseline WaveNet (S1) and Parallel WaveGAN (S2) systems, and even the improved Parallel WaveGAN with conditional discriminator (S3). Note that the MOS of the TTS samples tended to be higher than that of analysis/synthesis samples. This result was because the unwanted artifacts produced by the analysis/synthesis process were statistically excluded during the generation process. Most listeners preferred consistent results of the predicted duration and accent than those of the ground-truth.

5. CONCLUSION

We proposed voicing-aware conditional discriminators for a Parallel WaveGAN-based TTS system. Our framework incorporated a projection-based conditioning method into the discriminator and divided it into two separate discriminators. By controlling the voiced and unvoiced speech segments independently, the performance of each discriminator was significantly improved, which allowed the generator to produce more natural speech waveforms. The experimental results demonstrated the superiority of our proposed method over the conventional Parallel WaveGAN and similarly configured WaveNet systems. Future work includes improving the performance of the generator by utilizing the voicing information of speech.

6. ACKNOWLEDGEMENTS

This work was supported by Clova Voice, NAVER Corp., Seongnam, Korea.

7. REFERENCES

- [1] H Zen, A Senior, and M Schuster, “Statistical parametric speech synthesis using deep neural networks,” in *Proc. ICASSP*, 2013, pp. 7962–7966.
- [2] X Wang, J Lorenzo-Trueba, S Takaki, L Juvela, and J Yamagishi, “A comparison of recent waveform generation and acoustic modeling methods for neural-network-based speech synthesis,” in *Proc. ICASSP*, 2018, pp. 4804–4808.
- [3] A van den 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.
- [4] N Kalchbrenner, E Elsen, K Simonyan, S Noury, N Casagrande, E Lockhart, F Stimberg, A. v. d Oord, S Dieleman, and K Kavukcuoglu, “Efficient neural audio synthesis,” in *Proc. ICML*, 2018, pp. 2410–2419.
- [5] A Tamamori, T Hayashi, K Kobayashi, K Takeda, and T Toda, “Speaker-dependent WaveNet vocoder,” in *Proc. INTERSPEECH*, 2017, pp. 1118–1122.
- [6] T Hayashi, A Tamamori, K Kobayashi, K Takeda, and T Toda, “An investigation of multi-speaker training for WaveNet vocoder,” in *Proc. ASRU*, 2017, pp. 712–718.
- [7] E Song, K Byun, and H.-G Kang, “ExcitNet vocoder: A neural excitation model for parametric speech synthesis systems,” in *Proc. EUSIPCO*, 2019, pp. 1179–1183.
- [8] A van den Oord, Y Li, I Babuschkin, K Simonyan, O Vinyals, K Kavukcuoglu, G van den Driessche, E Lockhart, L. C Cobo, F Stimberg, et al., “Parallel WaveNet: Fast high-fidelity speech synthesis,” in *Proc. ICML*, 2018, pp. 3915–3923.
- [9] W Ping, K Peng, and J Chen, “ClariNet: Parallel wave generation in end-to-end text-to-speech,” in *Proc. ICLR*, 2019.
- [10] R Prenger, R Valle, and B Catanzaro, “WaveGlow: A flow-based generative network for speech synthesis,” in *Proc. ICASSP*, 2019, pp. 3617–3621.
- [11] S Kim, S Lee, J Song, J Kim, and S Yoon, “FloWaveNet : A generative flow for raw audio,” in *Proc. ICML*, 2019, pp. 3370–3378.
- [12] R Yamamoto, E Song, and J.-M Kim, “Probability density distillation with generative adversarial networks for high-quality parallel waveform generation,” in *Proc. INTERSPEECH*, 2019, pp. 699–703.
- [13] 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 *Proc. NeurIPS*, 2019, pp. 14881–14892.
- [14] R Yamamoto, E Song, and J.-M Kim, “Parallel WaveGAN: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram,” in *Proc. ICASSP*, 2020, pp. 6199–6203.
- [15] E Song, R Yamamoto, M.-J Hwang, J.-S Kim, O Kwon, and J.-M Kim, “Improved Parallel WaveGAN vocoder with perceptually weighted spectrogram loss,” in *Proc. SLT*, 2021, pp. 470–476.
- [16] T Miyato and M Koyama, “cGANs with projection discriminator,” in *Proc. ICLR*, 2018.
- [17] J Yang, J Lee, Y Kim, H.-Y Cho, and I Kim, “VocGAN: A high-fidelity real-time vocoder with a hierarchically-nested adversarial network,” in *Proc. INTERSPEECH*, 2020, pp. 200–204.
- [18] M Bińkowski, J Donahue, S Dieleman, A Clark, E Elsen, N Casagrande, L. C Cobo, and K Simonyan, “High fidelity speech synthesis with adversarial networks,” *Proc. ICLR*, 2020.
- [19] I Goodfellow, J Pouget-Abadie, M Mirza, B Xu, D Warde-Farley, S Ozair, A Courville, and Y Bengio, “Generative adversarial nets,” in *Proc. NIPS*, 2014, pp. 2672–2680.
- [20] X Mao, Q Li, H Xie, R. Y Lau, Z Wang, and S Paul Smolley, “Least squares generative adversarial networks,” in *Proc. ICCV*, 2017, pp. 2794–2802.
- [21] S. Ö Arık, H Jun, and G Diamos, “Fast spectrogram inversion using multi-head convolutional neural networks,” *IEEE Signal Process. Letters*, vol. 26, no. 1, pp. 94–98, 2019.
- [22] E Song, F. K Soong, and H.-G Kang, “Effective spectral and excitation modeling techniques for LSTM-RNN-based speech synthesis systems,” *IEEE/ACM Trans. Audio, Speech, and Lang. Process.*, vol. 25, no. 11, pp. 2152–2161, 2017.
- [23] L Liu, H Jiang, P He, W Chen, X Liu, J Gao, and J Han, “On the variance of the adaptive learning rate and beyond,” in *Proc. ICLR*, 2020.
- [24] A Odena, V Dumoulin, and C Olah, “Deconvolution and checkerboard artifacts,” *Distill*, 2016.
- [25] Y Ren, C Hu, T Qin, S Zhao, Z Zhao, and T.-Y Liu, “Fast-Speech 2: Fast and high-quality end-to-end text-to-speech,” in *Proc. ICLR (in press)*, 2021.
- [26] A Łańcucki, “FastPitch: Parallel text-to-speech with pitch prediction,” in *Proc. ICASSP (in press)*, 2021.
- [27] 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*, 2019, pp. 6905–6909.
- [28] N Li, S Liu, Y Liu, S Zhao, M Liu, and M. T Zhou, “Neural speech synthesis with Transformer network,” in *Proc. AAAI*, 2019, pp. 6706–6713.
- [29] T Hayashi, R Yamamoto, K Inoue, T Yoshimura, S Watanabe, T Toda, K Takeda, Y Zhang, and X Tan, “ESPnet-TTS: Unified, reproducible, and integratable open source end-to-end text-to-speech toolkit,” in *Proc. ICASSP*, 2020, pp. 7654–7658.