

FRAMEWISE WAVEGAN: HIGH SPEED ADVERSARIAL VOCODER IN TIME DOMAIN WITH VERY LOW COMPUTATIONAL COMPLEXITY

Ahmed Mustafa* Jean-Marc Valin* Jan Bütthe* Paris Smaragdis*[†] Mike Goodwin*

*Amazon Web Services, Palo Alto, USA

[†]University of Illinois at Urbana-Champaign

{ahdmust, jmvalin, jbuethe, parsmara, mmg}@amazon.com

ABSTRACT

GAN vocoders are currently one of the state-of-the-art methods for building high-quality neural waveform generative models. However, most of their architectures require dozens of billion floating-point operations per second (GFLOPS) to generate speech waveforms in samplewise manner. This makes GAN vocoders still challenging to run on normal CPUs without accelerators or parallel computers. In this work, we propose a new architecture for GAN vocoders that mainly depends on recurrent and fully-connected networks to directly generate the time domain signal in framewise manner. This results in considerable reduction of the computational cost and enables very fast generation on both GPUs and low-complexity CPUs. Experimental results show that our Framewise WaveGAN vocoder achieves significantly higher quality than auto-regressive maximum-likelihood vocoders such as LPCNet at a very low complexity of 1.2 GFLOPS. This makes GAN vocoders more practical on edge and low-power devices.

Index Terms— GAN vocoder, LPCNet, TTS, Speech synthesis, Neural speech coding

1. INTRODUCTION

The task of neural vocoding has seen rapid progress in recent years. This came with lots of applications that depend on neural vocoders for rendering high quality speech; such as TTS [1], speech enhancement [2] and neural speech coding [3, 4, 5]. Since the advent of Wavenet [6], deep generative models have become the standard way for building high quality neural vocoders, with clear outperformance over classical vocoding methods. By modelling raw waveform data distributions in the time domain, deep generative models can capture intricate representational details of the speech signal, which are not easy to track by fully-parametric methods. This results in generating high quality speech signals; but it also incurs huge amount of computations that scales rapidly with the target sampling rate. In case of WaveNet, this high computational cost becomes more challenging due to the auto-regressive (AR) sampling that makes the overall generation very slow. Later models such as WaveRNN [7] tried to address one part of this problem by using a sparse RNN that runs auto-regressively to learn the data distribution. LPCNet [8] further improved WaveRNN by using linear prediction to build lighter architectures while maintaining or even enhancing the signal quality; and currently it is one of the lowest complexity generative models available for raw audio. Further AR models [9, 10] follow the same approach but generate more than one sample at a time to speed up the inference process.

Parallel generative models propose a different way for dealing with the high computational cost. These models avoid auto-

regressive generation to parallelize computations via GPUs or CPU accelerators. There are many types of parallel vocoders which are similar in their architecture but differ in the way they are trained. GAN vocoders [11, 12, 13] are more competitive in parallel generation due to their light-weight models, as compared to maximum-likelihood (ML) vocoders like WaveGlow [14], and their one-shot generation, as opposed to diffusion vocoders [15, 16] that usually require iterative denoising steps to achieve high quality signal generation. However, despite their high quality and very fast signal generation, GAN vocoders still require multiples of Giga Floating-Point Operations per Second (GFLOPS), especially when using WaveNet-based architectures. This amount of computation is challenging for low complexity devices with limited power and parallelization capabilities.

Our work aims to reduce the large computational amount of current GAN vocoders. Specifically, we report the following contributions:

- We introduce Framewise WaveGAN, a neural vocoder generating wideband speech signals in the time domain frame-by-frame instead of sample-by-sample.
- We show how to train the model reliably and stably by training in perceptual domain using a combination of multi-resolution STFT regression and adversarial losses.
- We demonstrate the quality of the proposed model that outperforms LPCNet at a very low complexity of 1.2 GFLOPS according to subjective and objective scores.

2. FRAMEWISE WAVEGAN

Most GAN vocoders can be categorized into *WaveNet-based* and *latent-based* models. In WaveNet-based models, all network parameters are used at the same rate as the target signal. That is, if we seek to generate speech signal at sampling rate f_s , then every parameter in the model will do f_s multiplications plus f_s accumulations as the total number of floating-point operations per second (FLOPS). For light-weight vocoders such as Parallel WaveGAN [17] with 1.44 M parameters, generating one second of speech at 16 KHz would require 46 GFLOPS.

Latent-based models alleviate this computational burden by building the network as a stack of upsampling layers; where each layer operates at different resolution starting from the acoustic feature rate until reaching the target sampling rate of the signal. This can achieve moderate computational cost in ranges of 10 GFLOPS [18] and 8 GFLOPS [19], but with sophisticated architecture setups to maintain high quality.

Frame-wise WaveGAN is a first step towards running GAN vocoders in the time domain at the acoustic feature rate, without having to use upsampling layers which are the main source of high complexity. This is achieved by making the model generate one frame at a time. Normally, in WaveNet-based and latent-based models, the feature representations starting from the first layer until the last one are organized as tensors of $[Batch_dim, Channel_dim, Temporal_dim]$; with $Temporal_dim$ equal to the target signal resolution at the output layer. In Frame-wise WaveGAN, all feature representations are organized as $[Batch_dim, Sequence_dim, Frame_dim]$, where $Sequence_dim$ is equal everywhere to the acoustic feature resolution that is commonly much smaller than the signal one; and $Frame_dim$ holds the representation of the target frame that is being generated. The final waveform is obtained by simply *flattening* the generated frames at the model output. This leads to significant computational saving even with models of large memory footprint.

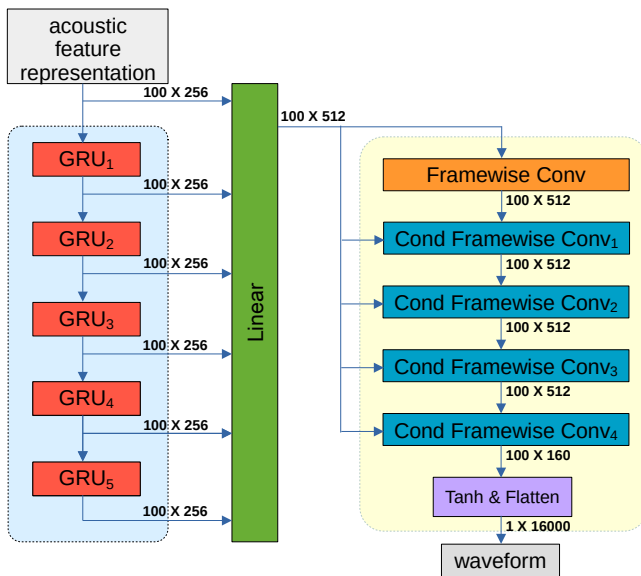


Fig. 1: Frame-wise WaveGAN generator architecture. The numbers show $[Sequence_dim, Frame_dim]$ of the output representation from each layer to generate 1 sec of speech waveform at sampling rate of 16 kHz, using conditioning acoustic features at 100 Hz.

2.1. Architecture Overview

Fig. 1 illustrates the architecture of the proposed model. It mainly consists of two stacks of recurrent and fully-connected layers. The recurrent stack has 5 GRUs [20] to model long-term dependencies of the signal. All GRU outputs are concatenated with the conditioning (i.e., acoustic feature) representation and converted into lower-dimensional latent representation through a simple fully-connected layer. This representation is then utilized by the fully-connected stack that operates in frame-wise convolutional manner to model the short-term dependencies of the signal.

2.2. Frame-wise Convolution

The term *frame-wise convolution* refers to a kernel whose elements are frames instead of samples. We implement this by making the fully-connected layer receive at frame index i a concatenation of

k frames at indices $\{i - k + 1, \dots, i\}$ from the input tensor, where k is the kernel size. The rest of the operation is same as normal convolution. There is also *conditional frame-wise convolution* that only differs from frame-wise convolution in concatenating an external feature frame (i.e., conditioning vector) to the layer input. In our model, we use 1 frame-wise convolution layer that receives the latent representation from the previous stack, with a kernel size of 3 frames, stride = dilation = 1 frame; and padding in non-causal manner (i.e., 1 look-ahead frame). Hence, if the input tensor to this layer has $Frame_dim$ of e.g., 512, then the fully-connected network should have $3 * 512 = 1536$ input dimensions. In addition, there are 4 conditional frame-wise convolution layers coming afterwards with a kernel size of 2 frames which are concatenated with 1 conditioning frame provided by the same latent representation obtained from the previous stack; with same stride, dilation and padding applied in causal manner. That's why the fully-connected network for this conditional layer has the same dimensionality as the non-conditional one. All of these frame-wise convolution operations are running in a single-channel sense; i.e., there is only one fully-connected network per layer. We used this way of implementation instead of traditional multi-channel convolution layers to ease the efficient implementation of the model, especially when applying sparsification methods to these layers.

2.3. Activation Layers

For all layers in the recurrent and frame-wise convolution stacks, we use Gated Linear Unit (GLU) [21] to activate their feature representations:

$$GLU(X) = X \otimes \sigma(FC(X)), \quad (1)$$

where FC is a simple fully-connected network to learn the sigmoid gate and it has the same output dimension as X , \otimes denotes element-wise multiplication. We also disabled the bias for all layers in the model; which was helpful for faster convergence with lower reconstruction artifacts.

2.4. Conditioning Acoustic Features

We use LPCNet acoustic features [8] to condition our vocoder model. They consist of 18 Bark-Frequency Cepstral Coefficients (BFCCs), a pitch period and a pitch correlation; which are extracted by 20 ms overlapping windows at 100 Hz frame rate from a 16 kHz speech waveform. The model generates one 10 ms frame per conditioning vector. The pitch period is fed to an embedding layer of 256 levels and 128 dimensions, while the BFCCs with the pitch correlation are fed to a simple causal convolution layer of 128 output channels and kernel size of 3. The outputs from these two layers are then concatenated and fed to another causal convolution layer of 256 output channels, kernel size of 3 and *Leaky_ReLU* activation to obtain the acoustic feature representation that is used for frame-wise generation, as shown in figure 1. LPCNet features are calculated on 10 ms frames with 5 ms look-ahead; and Frame-wise WaveGAN requires one feature frame look-ahead. Hence, the total delay is 10 ms for framing plus 15 ms look-ahead, which sums to 25 ms.

3. TRAINING PROCEDURE

3.1. Training in the Perceptual Domain

Speech signals are characterized by their high dynamic range as they go wider in bandwidth. When we apply a simple pre-emphasis filter before training, we find the vocoder able to learn high frequency

components faster than training in the normal signal domain. We reinforce this benefit by additionally using perceptual filtering, as detailed in AMR-WB [22], so that the vocoder can learn high frequency content even faster. The perceptual weighting filter is defined by the following transfer function:

$$W(z) = \frac{A(z/\gamma_1)}{(1 - \gamma_2 z^{-1})}, \quad (2)$$

where $A(z)$ is the linear prediction (LPC) filter whose coefficients are computed from BFCCs, $\gamma_1 = 0.92$ and $\gamma_2 = 0.85$. This filtering increases the spectral flatness of signals during the training, which enables clearly faster convergence. Moreover, when applying inverse filtering to obtain the final signal, the noise of reconstruction artifacts is shaped by $W^{-1}(z)P^{-1}(z)$, where $P^{-1}(z)$ is the de-emphasis applied at end of the synthesis. The computational cost of this perceptual filtering is also quite cheap and still keeps the overall complexity low.

3.2. Spectral Pre-training

We first pre-train the model using the spectral reconstruction loss \mathcal{L}_{aux} defined by Equation (6) in [17]. It is a combination of spectral magnitude and convergence losses obtained by different STFT resolutions. We use all power-of-two FFT sizes between 64 and 2048 (6 sizes), with same values for window sizes and 75% window overlap. For the spectral magnitude loss \mathcal{L}_{mag} defined by Equation (5) in [17], we apply *sqrt* instead of *log* as a non-linearity, which was found better for early convergence. The spectral pre-training gives a metallic-sounding signal with over-smoothed high frequency content, which is a good prior signal to start adversarial training for achieving realistic signal reconstruction.

3.3. Spectral Adversarial Training

Using time-domain discriminators is a major challenge in adversarial training of the proposed model. We tried different time-domain discriminator architectures [23, 12, 11] that all failed to achieve stable training. Instead, we use multi-resolution spectrogram discriminators, which achieve much better training behavior and reliably increase the fidelity of generated signals. We follow the same spectrogram discriminator architecture defined in [24] and use 6 models running on spectrograms of the same STFT resolutions used for spectral pre-training; with *sqrt* used as a non-linearity. The adversarial training uses least-square loss as a metric for evaluating discriminator outputs, with the same formulation given by Equations (1) and (2) in [11]. The spectral reconstruction loss is kept to regularize the adversarial training. Hence, the final generator objective is

$$\min_G \left(\mathbb{E}_z \left[\sum_{k=1}^6 (D_k(G(s)) - 1)^2 \right] + \mathcal{L}_{\text{aux}}(G) \right), \quad (3)$$

where s represents the conditioning features (e.g., LPCNet features). Weight normalization [25] is applied to all convolution layers of the discriminators (D_k) and all fully-connected layers of the generator (G).

4. EXPERIMENTS AND RESULTS

4.1. Experimental Setup

We train the proposed Framewise WaveGAN model for a universal vocoding task, where the synthesis is speaker and language independent. To achieve this, we use a training speech corpus of 205

hours sampled at 16 kHz and obtained from a combination of TTS datasets [26, 27, 28, 29, 30, 31, 32, 33, 34]. The training data contains more than 900 speakers in more than 34 languages and dialects. Pre-emphasis with a factor of 0.85 is applied to the speech signals when extracting the input LPCNet features. The complete synthesis is performed by running the vocoder model and then applying inverse perceptual filtering followed by de-emphasis. The training runs on an NVIDIA Tesla V100 GPU. We use batch size of 32 and each sample in the batch corresponds to 1 sec for both features and speech data tensors. The spectral pre-training is carried out for 1 M steps (i.e., ~ 50 epochs), with $lr_g = 10^{-4}$; and then the adversarial training runs for another 1 M steps, with $lr_d = 2 * 10^{-4}$ and lr_g is reduced to $5 * 10^{-5}$. AdamW [35] is used for both generator and discriminator optimizers, with $\beta = \{0.8, 0.99\}$.

4.2. Complexity

In Framewise WaveGAN, all model parameters are used at the same rate of the conditioning features, which is 100 Hz. Thus, each parameter contributes to one multiply and one accumulation (i.e., 2 floating-point operations) per generation step. The total computational complexity to generate one second is then given as follows:

$$C = N * 2 * S, \quad (4)$$

where C denotes the number of floating-point operations per second (FLOPS), N is the total count of model parameters and S is the number of generation steps to create one second. Framewise WaveGAN according to the architecture defined in sec 2 has 7.8 M parameters. This gives a total complexity of 1.5 GFLOPS, including 7.3 MFLOPS for calculating *tanh* and *sigmoid* functions. We further reduce this complexity by sparsifying the dense model, as done in [7, 8], with weight densities of 0.6 for all GRUs and 0.65 for all fully-connected layers except the last three, which are kept dense. This decreases the number of active parameters to 5.9 M and makes the total complexity 1.2 GFLOPS. We also checked the parallel generation speed for the Pytorch implementation of Framewise WaveGAN at 1.2 GFLOPS in real-time factor (RTF). The model runs 20x and 75x faster than real-time on CPU (Intel Xeon Platinum 8175M 2.50GHz) and GPU (Nvidia Tesla V100), respectively.

4.3. Quality

We evaluate the model on the PTDB-TUG speech corpus [36] and the NTT Multi-Lingual Speech Database for Telephonometry. From PTDB-TUG, we use all English speakers (10 male, 10 female) and randomly pick 200 concatenated pairs of sentences. For the NTT database, we select the American English and British English speakers (8 male, 8 female), which account for a total of 192 samples (12 samples per speaker). The training material did not include any data from the datasets used in testing. For comparison, we use the latest published low-complexity LPCNet [37]. In addition, we also evaluate the reference speech as an upper bound on quality, and we include the Speex 4kb/s wideband vocoder [38] as an anchor. We evaluate the models by the mean opinion score (MOS) results, obtained according to the crowd-sourcing methodology described in P.808 [39]. To check how the quality scales with complexity, we test two LPCNet models at {3, 1.2} GFLOPS against three Framewise WaveGAN vocoders: the dense and sparse baseline models at {1.5, 1.2} GFLOPS; and a higher-complexity model of 3 GFLOPS, whose architecture includes GRU size of 320 and 10 framewise convolution layers.

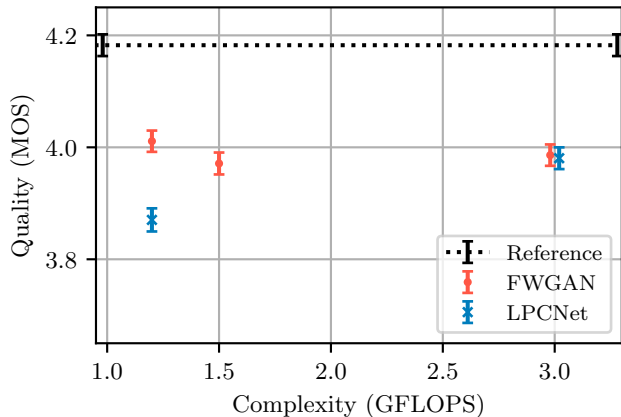


Fig. 2: Results from the MOS quality evaluation with 95% confidence interval for LPCNet and Framewise WaveGAN (FWGAN) at different complexities.

Fig. 2 demonstrates the MOS test results performed by 30 listeners using Amazon Mechanical Turk. Framewise WaveGAN (FWGAN) at 1.2 GFLOPS achieves significantly higher quality than LPCNet at the same complexity¹, even slightly outperforming (statistically significant) LPCNet operating at 3 GFLOPS. On the other hand, the higher-complexity FWGANs cannot outperform the 1.2 GFLOPS model. This is mainly due to the weaker discriminator behavior in adversarial training against generators with higher number of parameters. This may motivate further future work to come up with better discriminator models that enable quality scaling with higher generator complexities. We also found Framewise WaveGAN has clearly better pitch consistency than the LPCNet model, as shown in Table 1; where PMAE is the pitch mean absolute error between ground truth and vocoder samples, VDE is the voicing decision error (i.e., the percentage of frames with incorrect voicing decision). We use YAAPT pitch tracker [40] for extracting the voicing information. Although we do not target singing voice generation in this work, this finding may encourage using Framewise WaveGAN for building low complexity expressive speech synthesizers.

Table 1: Objective evaluation of voicing features, lower is better.

Model	PMAE	VDE
LPCNet 3 GFLOPS	5.5865	0.0168
LPCNet 1.2 GFLOPS	6.0965	0.0177
FWGAN 1.2 GFLOPS	5.0632	0.0163
FWGAN 1.5 GFLOPS	5.3502	0.0175
FWGAN 3 GFLOPS	5.4733	0.0169

¹See our demo samples at the following url: https://ahmed-fau.github.io/fwgan_demo/

5. CONCLUSION

This work introduced Framewise WaveGAN, an adversarial vocoder for wideband speech synthesis with very low complexity, close to 1 GFLOPS, and a low delay of 25 ms. The model architecture comprises recurrent and fully-connected layers that capture long- and short-term dependencies of the speech signal, with ability to generate the waveform frame-by-frame faster than real-time. Quality evaluations show that Framewise WaveGAN significantly outperforms state-of-the-art auto-regressive vocoders such as LPCNet. Potential future work includes better adversarial training methods to boost this kind of framewise vocoder in achieving high-fidelity speech generation at low delay and complexity.

6. ACKNOWLEDGMENTS

We would like to thank Umüt Isik, Timothy B. Terriberry, Karim Helwani and Erfan Soltanmohammadi for their kind support and fruitful tips.

7. REFERENCES

- [1] X. Tan, T. Qin, F. Soong, and T.-Y. Liu, “A survey on neural speech synthesis,” *arXiv preprint arXiv:2106.15561*, 2021.
- [2] J. Su, Z. Jin, and A. Finkelstein, “Hifi-gan-2: Studio-quality speech enhancement via generative adversarial networks conditioned on acoustic features,” in *2021 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 2021, pp. 166–170.
- [3] A. Mustafa, J. Büthe, S. Korse, K. Gupta, G. Fuchs, and N. Pia, “A streamwise gan vocoder for wideband speech coding at very low bit rate,” in *2021 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2021, pp. 66–70.
- [4] N. Zeghidour and et al., “Soundstream: An end-to-end neural audio codec,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 495–507, 2021.
- [5] N. Pia and et al., “Nesc: Robust neural end-2-end speech coding with gans,” *arXiv preprint arXiv:2207.03282*, 2022.
- [6] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, et al., “WaveNet: A Generative Model for Raw Audio,” *arXiv:1609.03499*, 2016.
- [7] N. Kalchbrenner, E. Elsen, K. Simonyan, S. Noury, et al., “Efficient neural audio synthesis,” *arXiv:1802.08435*, 2018.
- [8] J. Valin and J. Skoglund, “LPCNET: Improving Neural Speech Synthesis through Linear Prediction,” in *IEEE (ICASSP)*, 2019, pp. 5891–5895.
- [9] R. Vippera, S. Park, K. Choo, S. Ishtiaq, K. Min, S. Bhattacharya, A. Mehrotra, A. G. C. P. Ramos, and N. D. Lane, “Bunched LPCNet: Vocoder for low-cost neural text-to-speech systems,” *arXiv preprint arXiv:2008.04574*, 2020.
- [10] Q. Tian, Z. Zhang, H. Lu, L.-H. Chen, and S. Liu, “Featherwave: An efficient high-fidelity neural vocoder with multi-band linear prediction,” *arXiv preprint arXiv:2005.05551*, 2020.
- [11] J. Kong, J. Kim, and J. Bae, “Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis,” *arXiv preprint arXiv:2010.05646*, 2020.

- [12] A. Mustafa, N. Pia, and G. Fuchs, "StyleMelGAN: An efficient high-fidelity adversarial vocoder with temporal adaptive normalization," in *ICASSP 2021-2021*. IEEE, 2021, pp. 6034–6038.
- [13] S.-H. Lee, J.-H. Kim, K.-E. Lee, and S.-W. Lee, "Fre-gan 2: Fast and efficient frequency-consistent audio synthesis," in *ICASSP 2022 - 2022*, 2022, pp. 6192–6196.
- [14] R. Prenger, R. Valle, and B. Catanzaro, "Waveglow: A flow-based generative network for speech synthesis," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 3617–3621.
- [15] Z. Kong, W. Ping, J. Huang, K. Zhao, and B. Catanzaro, "Diffwave: A versatile diffusion model for audio synthesis," *arXiv preprint arXiv:2009.09761*, 2020.
- [16] Y. Koizumi, H. Zen, K. Yatabe, N. Chen, and M. Bacchiani, "Specgrad: Diffusion probabilistic model based neural vocoder with adaptive noise spectral shaping," *arXiv preprint arXiv:2203.16749*, 2022.
- [17] R. Yamamoto, E. Song, and J. Kim, "Parallel Wavegan: A Fast Waveform Generation Model Based on Generative Adversarial Networks with Multi-Resolution Spectrogram," in *IEEE-ICASSP*, 2020, pp. 6199–6203.
- [18] A. Mustafa, J. Büthe, S. Korse, K. Gupta, G. Fuchs, and N. Pia, "A streamwise gan vocoder for wideband speech coding at very low bit rate," in *2021 WASPAA*, 2021, pp. 66–70.
- [19] Z. Liu and Y. Qian, "Basis-melgan: Efficient neural vocoder based on audio decomposition," *arXiv preprint arXiv:2106.13419*, 2021.
- [20] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [21] Y. N. Dauphin, A. Fan, M. Auli, and D. Grangier, "Language modeling with gated convolutional networks," in *International conference on machine learning*. PMLR, 2017, pp. 933–941.
- [22] B. Bessette, R. Salami, and et al., "The adaptive multirate wideband speech codec (amr-wb)," *IEEE Transactions on Speech and Audio Processing*, vol. 10, no. 8, pp. 620–636, 2002.
- [23] K. Kumar, R. Kumar, de T. Boissiere, L. Gestein, et al., "MelGAN: Generative Adversarial Networks for Conditional Waveform Synthesis," in *Advances in NeurIPS 32*, pp. 14910–14921. 2019.
- [24] W. Jang, D. Lim, J. Yoon, B. Kim, and J. Kim, "UnivNet: A Neural Vocoder with Multi-Resolution Spectrogram Discriminators for High-Fidelity Waveform Generation," in *Proc. Interspeech 2021*, 2021, pp. 2207–2211.
- [25] T. Salimans and D. P. Kingma, "Weight normalization: A simple reparameterization to accelerate training of deep neural networks," in *Advances in NeurIPS*, 2016, pp. 901–909.
- [26] I. Demirsahin, O. Kjartansson, A. Gutkin, and C. Rivera, "Open-source Multi-speaker Corpora of the English Accents in the British Isles," in *Proc. LREC*, 2020.
- [27] O. Kjartansson, A. Gutkin, A. Butryna, I. Demirsahin, and C. Rivera, "Open-Source High Quality Speech Datasets for Basque, Catalan and Galician," in *Proc. SLTU and CCURL*, 2020.
- [28] A. Guevara-Rukoz, I. Demirsahin, F. He, S.-H. C. Chu, S. Sarin, K. Pipatsrisawat, A. Gutkin, A. Butryna, and O. Kjartansson, "Crowdsourcing Latin American Spanish for Low-Resource Text-to-Speech," in *Proc. LREC*, 2020.
- [29] F. He, S.-H. C. Chu, O. Kjartansson, C. Rivera, A. Katanova, A. Gutkin, I. Demirsahin, C. Johnny, M. Jansche, S. Sarin, and K. Pipatsrisawat, "Open-source Multi-speaker Speech Corpora for Building Gujarati, Kannada, Malayalam, Marathi, Tamil and Telugu Speech Synthesis Systems," in *Proc. LREC*, 2020.
- [30] K. Sodimana, K. Pipatsrisawat, L. Ha, M. Jansche, O. Kjartansson, P. De Silva, and S. Sarin, "A Step-by-Step Process for Building TTS Voices Using Open Source Data and Framework for Bangla, Javanese, Khmer, Nepali, Sinhala, and Sundanese," in *Proc. SLTU*, 2018.
- [31] Y. M. Oo, T. Wattanavekin, C. Li, P. De Silva, S. Sarin, K. Pipatsrisawat, M. Jansche, O. Kjartansson, and A. Gutkin, "Burmese Speech Corpus, Finite-State Text Normalization and Pronunciation Grammars with an Application to Text-to-Speech," in *Proc. LREC*, 2020.
- [32] D. van Niekerk, C. van Heerden, M. Davel, N. Kleynhans, O. Kjartansson, M. Jansche, and L. Ha, "Rapid development of TTS corpora for four South African languages," in *Proc. INTERSPEECH*, 2017.
- [33] A. Gutkin, I. Demirşahin, O. Kjartansson, C. Rivera, and K. Túbòsún, "Developing an Open-Source Corpus of Yoruba Speech," in *Proc. INTERSPEECH*, 2020.
- [34] E. Bakhturina, V. Lavrukhin, B. Ginsburg, and Y. Zhang, "Hi-Fi Multi-Speaker English TTS Dataset," *arXiv preprint arXiv:2104.01497*, 2021.
- [35] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," *arXiv preprint arXiv:1711.05101*, 2017.
- [36] G. Pirker, M. Wohlmayr, S. Petrik, and F. Pernkopf, "A pitch tracking corpus with evaluation on multipitch tracking scenario.," in *Proc. INTERSPEECH*, 2011, pp. 1509–1512.
- [37] K. Subramani, J.-M. Valin, U. Isik, P. Smaragdis, and A. Krishnaswamy, "End-to-end LPCNet: A neural vocoder with fully-differentiable lpc estimation," *arXiv preprint arXiv:2202.11301*, 2022.
- [38] J.-M. Valin, "The speex codec manual," 2007.
- [39] ITU-T, *Recommendation P.808: Subjective evaluation of speech quality with a crowdsourcing approach*, 2018.
- [40] K. Kasi and S. A. Zahorian, "Yet another algorithm for pitch tracking," in *2002 IEEE International Conference on Acoustics, Speech, and Signal Processing*, 2002, vol. 1, pp. I–361–I–364.