

SPEECH SYNTHESIS AND CONTROL USING DIFFERENTIABLE DSP

Giorgio Fabbro^{1,2} Vladimir Golkov² Thomas Kemp¹ Daniel Cremers²

¹ Sony Europe B.V., Stuttgart ² Computer Vision Group, Technical University of Munich

ABSTRACT

Modern text-to-speech systems are able to produce natural and high-quality speech, but speech contains factors of variation (e.g. pitch, rhythm, loudness, timbre) that text alone cannot contain. In this work we move towards a speech synthesis system that can produce diverse speech renditions of a text by allowing (but not requiring) explicit control over the various factors of variation. We propose a new neural vocoder that offers control of such factors of variation. This is achieved by employing differentiable digital signal processing (DDSP) (previously used only for music rather than speech), which exposes these factors of variation. The results show that the proposed approach can produce natural speech with realistic timbre, and individual factors of variation can be freely controlled.

Index Terms— Neural vocoder, speech synthesis, digital signal processing, neural networks, deep learning

1. INTRODUCTION

A sentence can be pronounced in many different ways: fast or slow, happy or angry, with emphasis on certain words. We can say much more than what is written because we are able to control some fundamental aspects of the production of speech. In this work, we call these aspects *factors of variation*, since their variation affects the overall meaning of speech. These factors of variation include *pitch*, *rhythm*, *loudness*, and *timbre*. In this work, we move towards a speech synthesis system that can produce diverse speech renditions of a text by allowing explicit control over the various factors of variation.

In literature, control is achieved in different ways. Here we make an important distinction between models that *require* control and models that offer *optional* control. Models that require control expect additional inputs to generate speech, whereas models that offer optional control disentangle the data into various components and provide the possibility to modify them. Optional control is preferable, as not always affecting the generation is desired and not always we possess all the required inputs.

Digital signal processing (DSP) algorithms are central methods in audio engineering. Recently, differentiable DSP (DDSP) [1] has been introduced as a new method to generate audio with deep learning: DSP algorithms are used as parts of a neural network, ensuring end-to-end optimization. DDSP so far has only been applied to music. We use DDSP to generate speech. Since DDSP generators directly map the controllable variables to audio, we propose control where the audio is generated (i.e. in the neural vocoder), and not earlier in the pipeline as others do. This choice has several advantages: it can be used not only with a spectrogram generator in text-to-speech (TTS), but also to modify properties of existing audio clips. Moreover, our neural vocoder provides direct control to any spectrogram generator, without the effort and the difficulty of redesigning the latter to offer control. This work is intended as a

first step towards DDSP-based speech synthesis systems for a wide range of applications that are able to let the user affect aspects of speech that text alone cannot encode.

2. RELATED WORK

The goal of TTS is to convert a text sequence into an audio rendition of someone’s voice pronouncing that text. The typical TTS system employs two steps to achieve that goal. First, a model converts the text into an acoustic representation of speech, usually a mel spectrogram. We call such a model a (*mel*) *spectrogram generator*. Then, another model converts the mel spectrogram into the audio waveform. This model is usually referred to as a (*neural*) *vocoder*. In existing methods, control of factors of variation happens during the mel spectrogram generation. Our method, on the other hand, offers control in the neural vocoder. In the following we first outline some models that generate mel spectrograms while offering control. Then we outline the major neural vocoders in literature.

Various approaches have been proposed to affect the mel spectrogram generation from text. Those that *require* control use additional inputs to the networks as in Fig. 1a [4, 5, 2, 6, 7]. Those that offer control (but *do not require* it) aim at controlling latent variables while ensuring that they are interpretable [8, 9, 10, 11], or disentangle one interpretable control variable (often the duration of the utterance) and make it controllable as in Fig. 1b [3, 12, 13].

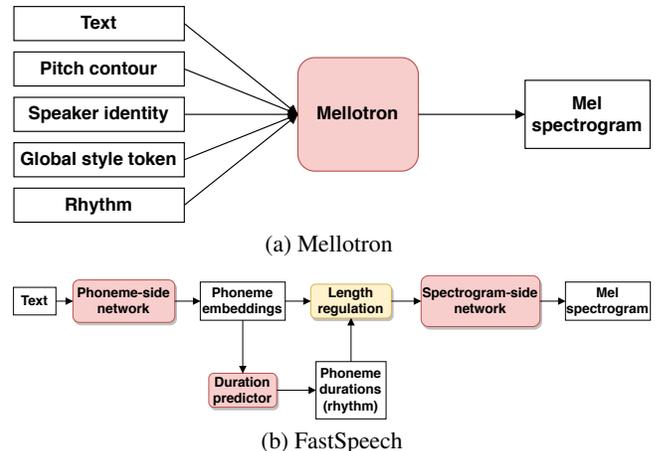


Fig. 1: Comparison of two existing models that use factors of variation such as rhythm to affect the audio generation. Red blocks are learned models, yellow blocks are fixed operations. (a) Mellotron [2] does not infer the factors of variation, it always requires them as input. (b) FastSpeech [3] infers phoneme durations (i.e. rhythm) and offers (but *does not require*) their control. Offering control without requiring it makes a model applicable in a wider variety of settings.

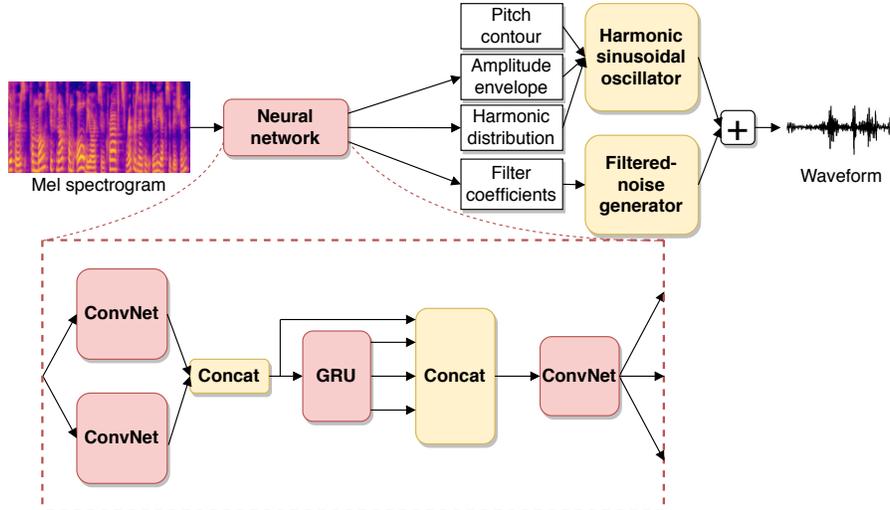


Fig. 2: The proposed DDSP-based neural vocoder. The neural network decomposes the spectrogram into control variables (amplitude envelope, harmonic distribution, and filter coefficients) for the DSP generators: this allows us to control them, if needed. We use the ground-truth pitch contour during training instead of letting the neural network infer it, in order to study in isolation the network’s capabilities of modeling the other three control variables. The oscillator produces vowels, the noise generator produces consonants.

Existing neural vocoders use various principles such as autoregressive networks [14], flows [15, 16, 17, 18, 19], adversarial learning [20, 21, 22, 23], and source-filtering [24]. These principles do not expose the factors of variation as interpretable variables and thus do not allow their control. We propose using a different principle to generate speech, namely DDSP [1], that so far was used only for music synthesis.

3. METHODS

Our objective is to have a model that offers (but *does not require*) control of factors of variation, so that we can modify them if needed. A simple way of achieving this is to include in the model architecture a fixed (i.e. not learned) operation that requires such variable (and make sure that no other operation will further perturb its outcome). For instance, the length regulation algorithm in Figure 1b requires phonemes duration as input: in this way, the authors can disentangle the speech rhythm from the data and make it available to the user for control. Another approach is to use DSP algorithms such as oscillators, since they generate audio signals using variables that we would like to control. Following [1], we insert differentiable DSP algorithms to synthesize audio (in our case speech): our model architecture has then fixed operations that use our controllable variables.

3.1. Network Architecture

Our neural vocoder is depicted in Figure 2. The model takes as input the mel spectrogram of a sentence and uses a recurrent neural network to decompose it into a set of control variables. These variables are input to a harmonic oscillator and a filtered-noise generator. These two final blocks output two audio sequences, which get summed to produce the final speech signal. The oscillator and the noise generator are differentiable (as in [1]), therefore we can apply a loss to the output signal and expect the gradients to flow back to the weights of the neural network. Note that the usage of an oscillator to

model the voiced parts of speech (vowels) and a filtered-noise generator to model consonant sounds makes our model similar to spectral modeling synthesis [25] and to neural source-filter models [24].

3.1.1. Trainable Layers

For simplicity, we use the same network architecture as the decoder in DDSP [1]. The decoder first applies to each of its inputs a 1D ConvNet with layer normalization, leaky ReLUs, and filter size 1 (equivalent to applying a multilayer perceptron to each time frame), then concatenates the outputs of the ConvNets over the channels dimension and feeds them to a gated recurrent unit (GRU) that processes them over the time dimension. The output of the GRU is then concatenated with its input and fed to another ConvNet. Its output is then split over the channels dimension to create the desired control variables. Our network is different from the DDSP decoder [1] only in the input stage: their inputs are intermediate variables (pitch, loudness, and a latent variable) each processed by a separate ConvNet, whereas our input is only the mel spectrogram. Therefore, we reduced the number of parallel ConvNets in the input stage. Still, in [1] the authors could have used a single ConvNet and feed it with the three inputs concatenated along the channels dimension, but they used three separate ConvNets in parallel. This indicates that the network performs better if it processes different information (for example pitch and loudness) in separate ConvNets (without cross-talk) before merging the extracted features. Similarly, we use two ConvNets in parallel. In this way the network can learn to extract and process different parts of information from the spectrogram separately before merging high-level features. Hyperparameters are listed in Section 3.2.

In order to make the optimization easier in the first experiments and to focus on learning to infer a specific subset of control variables, we provide to the harmonic oscillator the pitch contour extracted from the original audio, so that the efforts of the neural network are focused on loudness, timbre, and rhythm. This does not affect the capabilities of the neural vocoder in controlling pitch.

We let the network work at an intermediate temporal resolution

that is higher than in the mel spectrogram but lower than in the audio. To this end, the mel spectrogram is upsampled by a factor of 8.5 using bilinear interpolation before going into the network. Lower upsampling factors yielded lower audio quality, whereas higher ones yielded similar quality but slower computation. The network output is further upsampled using bilinear interpolation to reach the temporal resolution of the audio.

3.1.2. Harmonic Sinusoidal Oscillator

The time-upsampled outputs of the neural network are used to control two DSP generators: a harmonic sinusoidal oscillator and a filtered-noise generator. The harmonic sinusoidal oscillator is controlled by the fundamental frequency (pitch contour) $f_1(n) \in \mathbb{R}^+$ for each time instant n , an amplitude envelope $A(n) \in \mathbb{R}^+$ for each n , and a distribution $c_k(n)$ over harmonics that for each time step contains the weight to be applied to each harmonic (the weights sum up to 1 for every time step). The distribution c_k characterizes the timbre and the vowel. Following [1], the amplitude for each harmonic is $A_k(n) = A(n)c_k(n)$.

The oscillator generates a superposition $y(n)$ of harmonically related sinusoidal signals, i.e. signals whose frequency $f_k(n)$ is an integer multiple of the fundamental frequency $f_1(n)$, i.e. $f_k(n) = kf_1(n)$ with $k \in \mathbb{Z}^+$. Since digital signals are bandlimited, we only consider a finite number H of harmonics, i.e. $1 < k < H$. The output of the oscillator is:

$$y(n) = \sum_{k=1}^H A_k(n) \sin(\phi_k(n)), \quad (1)$$

where $\phi_k(n) = 2\pi \sum_{i=0}^n f_k(i)$ is the so-called instantaneous phase for time instant n and harmonic k . The only hyperparameter that must be chosen for the harmonic oscillator is the number H of harmonics it generates. The dataset we used was recorded at a sampling frequency $f_s = 22\,050$ Hz, so the Nyquist frequency is $f_{\text{Nyq}} = f_s/2 = 11\,025$ Hz. A generic female speech signal has a fundamental frequency with lower bound $f_1^{(\min)} = 165$ Hz [26] (we model female voice because our dataset contains only female speech). Therefore, in order to effectively model all harmonic content that our signal potentially has, we need $H = f_{\text{Nyq}}/f_1^{(\min)} \approx 67$ harmonics.

3.1.3. Filtered-Noise Generator

To add non-harmonic components (consonants) to our synthesis process, we use a module that generates white noise and filters it with a linear time-varying filter bank [1]. We control this *filtered-noise generator* by letting the neural network output parameters for the time-varying filter bank. For speed, the filtering operation occurs in the frequency domain: the time-upsampled outputs of the neural network are elementwise multiplied with white noise in the frequency domain. The signal is then converted to time-domain using the overlap-add method to account for the overlap between adjacent time frames. We can choose the frequency resolution of the filter by changing the number M of frequency-domain coefficients (network output channels). By increasing M , we decrease the amount of frequencies that each coefficient corresponds to, therefore increasing the filter resolution. We found that $M = 101$ gave us good results, while keeping the computational load low.

3.1.4. Training Objective

Our loss is computed using spectrograms with different time resolutions, similarly to the multi-scale spectrogram loss from DDSPP [1].

Since we propose to adapt DDSPP to speech, we define the loss via mel spectrograms, which are widely used for speech applications, instead of short-time Fourier transforms (STFT). A mel spectrogram is computed by mapping the result of STFT to the mel scale. This implies that in our case the frequency resolution will remain always constant for STFT results computed with different time-frequency resolutions, as we map all of them to the same mel scale. Therefore, our loss is

$$\mathcal{L} = \sum_{i \in R} \left\| S_i^{\text{mel}} - \hat{S}_i^{\text{mel}} \right\|_1 + \alpha \left\| \log S_i^{\text{mel}} - \log \hat{S}_i^{\text{mel}} \right\|_1, \quad (2)$$

where $R = \{2048, 1024, 512, 256, 128, 64\}$ is the set of time-frequency resolutions measured in the number of samples, S^{mel} is the mel spectrogram of the ground-truth audio and \hat{S}^{mel} is the mel spectrogram of the generated audio. In all the experiments, we used mel spectrograms with 80 frequency bins. The time overlap of the audio data between neighbouring frames in each spectrogram is 75%.

3.2. Experimental Setup

We trained on the LJSpeech dataset [27], which contains 13100 short audio clips of one female speaker reading non-fictional passages. The clips vary in length between 1 and 10 seconds and have an average length of 6.5 seconds. The total length of the dataset is approximately 24 hours. The audio sampling frequency is $f_s = 22\,050$ Hz and the encoding is 16-bit PCM WAV. No pre-processing was applied to the original audio and the input mel spectrograms were obtained from the ground truth audio by applying the STFT where each frame is the result of an FFT applied to 1024 samples of audio; each frame overlaps with the previous one and the next one by 75%. They have 80 frequency bins that cover the interval [0 Hz; 8000 Hz]. We provided our model with ground truth pitch contours extracted with the YIN algorithm [28] using the same parameters as for the mel spectrograms. We used the first 12822 clips as the training set, and the remaining 278 clips (i.e. the files named LJ050-*.wav) as the test set. For validation, we randomly selected a portion of the training set. As baselines, we used the official implementation of WaveGlow [16] and our implementation of WaveNet [14] in NNabla¹.

We used the Adam optimizer [29] with learning rate 10^{-4} , $\beta_1 = 0.9$ and $\beta_2 = 0.999$. We applied an exponential learning rate decay of 0.98 every 10^4 iterations. We stopped training after $4 \cdot 10^5$ iterations as we could not notice any more improvement.

Our best configuration uses the following hyperparameters: batch size 8, batch length $6 \cdot 10^4$ samples, and $M = 101$. The GRU layer in the architecture has 512 units. Each ConvNet has 3 layers with filter size 1 and 512 filters in the hidden layers. We also tried a smaller architecture (1 layer in each ConvNet with 256 filters in hidden layers), but results were worse.

4. RESULTS AND DISCUSSION

4.1. Audio Clips

Typical clips generated by our neural vocoder are available on the project's GitHub page². Our speech is fully intelligible and we model accurately the speaker's timbre. Some consonants still sound slightly artificial and some slight noise from the noise generator can be heard even when vowels are pronounced.

¹www.github.com/sony/nnabla

²<https://thesmith1.github.io/DDSPeech/>

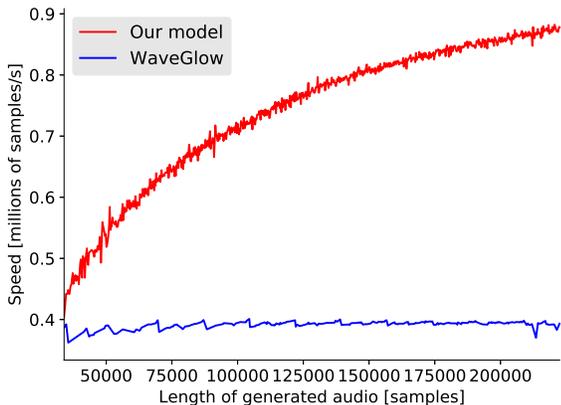


Fig. 3: Synthesis speed for WaveGlow (which is fast and high-quality) and our model. Our model can synthesize audio at almost $40\times$ real-time for long audio sequences.

On the same webpage, we also provide some clips to showcase the control capabilities of our neural vocoder. In particular, we show how we can freely change the pitch of the utterance, either by modification or by substitution. Moreover, in some cases we perceive traces of the original pitch (e.g. in the clip where the pitch is modified to be constant). These traces appear to be encoded in the timbre.

4.2. Synthesis Speed

Figure 3 shows a comparison of the inference speed for WaveGlow and our model, given different lengths of the audio to be synthesized. We avoid the inclusion of WaveNet, since it is a sequential approach and therefore very slow. The speed of WaveGlow is overall constant with the output size, whereas our model generates more samples per second the more samples the audio clip has, outperforming WaveGlow for clips longer than 2 seconds. We attribute this different scaling behavior of the two models to the GPU using different optimizations for different architectures. To confirm this, we repeated the measurement on the CPU (using a single core) and established that the speed (in samples per second) of our model as a function of clip length is constant on the CPU.

4.3. Training Time

Our DDSP-based neural vocoder converged in 2.5 days on one GPU. This is relatively fast compared to other popular neural vocoders such as WaveNet, which required 4–5 days, or WaveGlow, which required 10 days and more than one GPU.

4.4. Model Size

Table 1 shows the sizes of various neural vocoders, expressed as the number of learnable parameters. Our neural vocoder is relatively small: without compression it needs less than 20 MB of memory space. Such a small size makes it possible to utilize our model on edge devices. Unfortunately, the majority of such devices lacks the presence of a GPU, therefore the inference speed for now would be very low. On the other hand, once that issue is solved, our neural vocoder will represent a viable solution also for embedded systems.

Model	Size
WaveNet	6.8 million
WaveGlow	87.8 million
Our model	4.9 million

Table 1: Numbers of parameters of different neural vocoders. Our model is much smaller than WaveGlow and is comparable to WaveNet. Moreover, our model offers control of interpretable factors of variation of speech, such as pitch.

4.5. Listening Test

We performed a MUSHRA listening test [30, 31] with seven clips from the test set, which offers statistical significance for small numbers of participants. To reject unreliable listeners who give high scores to bad audio as per the MUSHRA protocol, we created bad audio by randomly shifting in time each frequency band of ground-truth audio by up to 5 spectrogram frames, adding the elementwise product of the original signal with white noise, and low-pass filtering. Of the 37 people who participated in our listening test, 8 were filtered out because they either assigned a score below 90 to a hidden copy of the reference, or they assigned a score higher than 30 to the lower anchor. MUSHRA scores for our neural vocoder had a median of 40 (out of 100), with a relatively high variance, because, even though the speaker’s timbre is well modeled, not all consonants sound realistic and it is easy to tell it apart from the real audio clip. The Griffin–Lim algorithm [32] had similar scores to our neural vocoder, even though its outputs have different characteristics, as consonants did not sound artificial but the speaker’s timbre is contaminated by the imperfect estimation of the phase. This indicates that a few-dimensional quality score cannot reflect many dimensions of audio quality, especially in the case of moderate scores. In other words, moderate scores do not tell us which of the dimensions of quality are worth improving. WaveNet and WaveGlow yield results that are good along all dimensions of quality, and hence also received a high one-dimensional score, namely a median MUSHRA score of 90 and 84, respectively.

5. CONCLUSIONS AND FUTURE WORK

This paper proposed using DDSP operators within neural vocoders. This offers the possibility of controlling the pitch contour and other factors of variation of an utterance. Even though the quality of the synthesized speech can be improved, the control capabilities that the model offers open up new research directions to study further how speech is generated and new opportunities for many application fields.

Future research directions include the usage of a different synthesis operator to make the model faster on the CPU: one example could be to use a wavetable oscillator instead of a harmonic one. Moreover, taking into account the dependency between pitch and timbre will improve synthesis quality. Finally, disentangling the factors of variation from the mel spectrogram enables the usage of other statistical models together with our neural vocoder: such models would operate on the factors of variation and could provide natural-sounding variations to generate alternative versions of the same utterance.

This work was supported by the Munich Center for Machine Learning (Grant No. 01IS18036B) and the BMBF project MLwin.

References

- [1] J. H. Engel, L. Hantrakul, C. Gu, and A. Roberts, “DDSP: differentiable digital signal processing,” in *8th International Conference on Learning Representations, ICLR*. 2020, OpenReview.net.
- [2] R. Valle, J. Li, R. Prenger, and B. Catanzaro, “Mellotron: Multispeaker expressive voice synthesis by conditioning on rhythm, pitch and global style tokens,” in *2020 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP*. 2020, pp. 6189–6193, IEEE.
- [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, NeurIPS*, H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, Eds., 2019, pp. 3165–3174.
- [4] Y. Wang, D. Stanton, Y. Zhang, R. J. Skerry-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, ICML*. 2018, vol. 80, pp. 5167–5176, PMLR.
- [5] S. Shechtman, C. Rabinovitz, A. Sorin, Z. Kons, and R. Hoory, “Controllable sequence-to-sequence neural TTS with LPC-Net backend for real-time speech synthesis on CPU,” *arXiv preprint arXiv:2002.10708*, 2020.
- [6] J. Shen, Y. Jia, M. Chrzanowski, Y. Zhang, I. Elias, H. Zen, and Y. Wu, “Non-attentive tacotron: Robust and controllable neural TTS synthesis including unsupervised duration modeling,” *CoRR*, vol. abs/2010.04301, 2020.
- [7] J. Bae, H. Bae, Y. Joo, J. Lee, G. Lee, and H. Cho, “Speaking speed control of end-to-end speech synthesis using sentence-level conditioning,” *CoRR*, vol. abs/2007.15281, 2020.
- [8] T. Raitio, R. Rasipuram, and D. Castellani, “Controllable neural text-to-speech synthesis using intuitive prosodic features,” *CoRR*, vol. abs/2009.06775, 2020.
- [9] R. Valle, K. J. Shih, R. Prenger, and B. Catanzaro, “Flowtron: an autoregressive flow-based generative network for text-to-speech synthesis,” *CoRR*, vol. abs/2005.05957, 2020.
- [10] W. Hsu, Y. Zhang, R. J. Weiss, H. Zen, Y. Wu, Y. Wang, Y. Cao, Y. Jia, Z. Chen, J. Shen, P. Nguyen, and R. Pang, “Hierarchical generative modeling for controllable speech synthesis,” in *Proceedings of the 35th International Conference on Machine Learning, ICML*. 2019, OpenReview.net.
- [11] R. Habib, S. Mariooryad, M. Shannon, E. Battenberg, R. J. Skerry-Ryan, D. Stanton, D. Kao, and T. Bagby, “Semi-supervised generative modeling for controllable speech synthesis,” in *8th International Conference on Learning Representations, ICLR*. 2020, OpenReview.net.
- [12] 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,” *CoRR*, vol. abs/2006.04558, 2020.
- [13] J. Kim, S. Kim, J. Kong, and S. Yoon, “Glow-TTS: A generative flow for text-to-speech via monotonic alignment search,” *CoRR*, vol. abs/2005.11129, 2020.
- [14] 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*. 2016, p. 125, ISCA.
- [15] 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, N. Casagrande, D. Grewe, S. Noury, S. Dieleman, E. Elsen, N. Kalchbrenner, H. Zen, A. Graves, H. King, T. Walters, D. Belov, and D. Hassabis, “Parallel WaveNet: Fast high-fidelity speech synthesis,” in *Proceedings of the 35th International Conference on Machine Learning, ICML*. 2018, vol. 80, pp. 3915–3923, PMLR.
- [16] R. Prenger, R. Valle, and B. Catanzaro, “WaveGlow: A flow-based generative network for speech synthesis,” in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP*. 2019, pp. 3617–3621, IEEE.
- [17] S. Kim, S. Lee, J. Song, J. Kim, and S. Yoon, “FloWaveNet: A generative flow for raw audio,” in *Proceedings of the 36th International Conference on Machine Learning, ICML*, K. Chaudhuri and R. Salakhutdinov, Eds. 2019, vol. 97 of *Proceedings of Machine Learning Research*, pp. 3370–3378, PMLR.
- [18] B. Zhai, T. Gao, F. Xue, D. Rothchild, B. Wu, J. E. Gonzalez, and K. Keutzer, “SqueezeWave: Extremely lightweight vocoders for on-device speech synthesis,” *CoRR*, vol. abs/2001.05685, 2020.
- [19] W. Ping, K. Peng, K. Zhao, and Z. Song, “WaveFlow: A compact flow-based model for raw audio,” *CoRR*, vol. abs/1912.01219, 2019.
- [20] C. Donahue, J. J. McAuley, and M. S. Puckette, “Adversarial audio synthesis,” in *7th International Conference on Learning Representations, ICLR*. 2019, OpenReview.net.
- [21] M. Binkowski, J. Donahue, S. Dieleman, A. Clark, E. Elsen, N. Casagrande, L. C. Cobo, and K. Simonyan, “High fidelity speech synthesis with adversarial networks,” in *8th International Conference on Learning Representations, ICLR*. 2020, OpenReview.net.
- [22] K. Kumar, R. Kumar, T. de Boissiere, L. Gestein, 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, NeurIPS*, H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, Eds., 2019, pp. 14881–14892.
- [23] 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 International Conference on Acoustics, Speech and Signal Processing, ICASSP*. 2020, pp. 6199–6203, IEEE.

- [24] X. Wang, S. Takaki, and J. Yamagishi, “Neural source-filter waveform models for statistical parametric speech synthesis,” *IEEE ACM Trans. Audio Speech Lang. Process.*, vol. 28, pp. 402–415, 2020.
- [25] X. Serra and J. O. Smith, “Spectral modeling synthesis: A sound analysis/synthesis based on a deterministic plus stochastic decomposition,” *Computer Music Journal*, vol. 14, pp. 12–24, 1990, SMS.
- [26] I. R. Titze, *Principles of Voice Production*, Prentice Hall, 1994.
- [27] K. Ito, “The LJ speech dataset,” <https://keithito.com/LJ-Speech-Dataset/>, 2017.
- [28] A. Cheveigné and H. Kawahara, “YIN, a fundamental frequency estimator for speech and music,” *The Journal of the Acoustical Society of America*, vol. 111, pp. 1917–30, 05 2002.
- [29] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *3rd International Conference on Learning Representations, ICLR*, Y. Bengio and Y. LeCun, Eds., 2015.
- [30] International Telecommunication Union, *Method for the subjective assessment of intermediate quality levels of coding systems*, Geneva, 2015.
- [31] M. Schoeffler, S. Bartoschek, F.-R. Stöter, M. Roess, S. Westphal, B. Edler, and J. Herre, “webMUSHRA — a comprehensive framework for web-based listening tests,” *Journal of Open Research Software*, vol. 6, 02 2018.
- [32] D. Griffin and J. Lim, “Signal estimation from modified short-time fourier transform,” *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 32, no. 2, pp. 236–243, April 1984.