

Neural Text-to-Speech with a Modeling-by-Generation Excitation Vocoder

Eunwoo Song¹, Min-Jae Hwang², Ryuichi Yamamoto³, Jin-Seob Kim¹,
Ohsung Kwon¹, and Jae-Min Kim¹

¹NAVER Corp., Seongnam, Korea

²Search Solutions Inc., Seongnam, Korea

³LINE Corp., Tokyo, Japan

Abstract

This paper proposes a modeling-by-generation (MbG) excitation vocoder for a neural text-to-speech (TTS) system. Recently proposed neural excitation vocoders can realize qualified waveform generation by combining a vocal tract filter with a WaveNet-based glottal excitation generator. However, when these vocoders are used in a TTS system, the quality of synthesized speech is often degraded owing to a mismatch between training and synthesis steps. Specifically, the vocoder is separately trained from an acoustic model front-end. Therefore, estimation errors of the acoustic model are inevitably boosted throughout the synthesis process of the vocoder back-end. To address this problem, we propose to incorporate an MbG structure into the vocoder’s training process. In the proposed method, the excitation signal is extracted by the acoustic model’s generated spectral parameters, and the neural vocoder is then optimized not only to learn the target excitation’s distribution but also to compensate for the estimation errors occurring from the acoustic model. Furthermore, as the generated spectral parameters are shared in the training and synthesis steps, their mismatch conditions can be reduced effectively. The experimental results verify that the proposed system provides high-quality synthetic speech by achieving a mean opinion score of 4.57 within the TTS framework.

Index Terms: neural text-to-speech, WaveNet, ExcitNet, modeling-by-generation vocoder

1. Introduction

Generative models for raw speech waveform have significantly improved the quality of neural text-to-speech (TTS) systems [1, 2]. Specifically, by conditioning acoustic features to the network input, neural vocoding models such as WaveNet, WaveRNN, and WaveGlow successfully generate a time-sequence of speech signal [2–5]. More recently, neural excitation vocoders such as GlotNet, ExcitNet, LP-WaveNet and LPCNet [6–10] have exploited the advantages of linear prediction (LP)-based parametric vocoders. In this type of vocoder, an adaptive predictor is used to decouple the formant-related spectral structure from the input speech signal, and the probability distribution of its residual signal (i.e. the excitation signal) is then modeled by the vocoding network. As variation in the excitation signal is only constrained by vocal cord movement, the training and generation processes become much more efficient.

However, because the vocoding and acoustic models have been trained separately, it is not known whether or not combining them within the TTS framework would benefit synthesis quality. Furthermore, as parameters estimated from the acoustic model are used as a direct input of the vocoding model in the synthesis step, estimation errors of the acoustic features can be

propagated throughout the synthesis process. It is therefore crucial to model the interactions between the acoustic and vocoding elements during the training process in order to achieve the best complete performance of the TTS system.

In this paper, we propose a neural excitation model based on modeling-by-generation (MbG) in which the spectral parameters generated from the acoustic model are utilized in the neural vocoder’s training process. Specifically, the target excitation is defined as a combination of the prediction errors from the LP analysis and those from the acoustic model. The vocoding model is then optimized to learn the distribution of the target excitation while compensating for the errors from the acoustic model. It has been reported elsewhere that training the neural vocoder with generated acoustic parameters improves synthetic quality [11]. Although the MbG method is similar to this approach, there are also clear differences in that MbG aligns even the target excitation signal with the acoustic model’s generated spectral parameters.

We investigated the effectiveness of the proposed method by conducting subjective evaluation tasks. The MbG structure can be extended to any neural excitation vocoder that uses LP coefficients, but the focus here is on the WaveNet-based ExcitNet vocoder [7]. The experimental results show that a TTS system with the proposed MbG-ExcitNet vocoder provides significantly better perceptual quality than a similarly configured system with a conventional vocoder. In particular, our TTS framework achieves 4.57 mean opinion score (MOS).

2. Related work

The idea of using an MbG structure is not new. In a study of parametric *glottal vocoders*, Juvola et al. [12] first proposed the closed-loop extraction of glottal excitation from the generated spectral parameters, and our own previous work proposed the MbG structure to compensate for missing noise components in generated glottal signals [13]. However, it was not possible to fully utilize the effectiveness of the MbG training strategy because our experiments were only performed with simple deep learning models including stacked feed-forward and/or long short-term memory (LSTM) networks.

Our aim here was to extend the usage of the MbG structure to recently proposed neural excitation models (e.g. ExcitNet) with autoregressive acoustic models (e.g. Tacotron) [11, 14, 15]. As the accuracy of acoustic models has been significantly improved, it is now possible to extract stable excitation signals from the generated spectral parameters. Furthermore, the ExcitNet vocoder directly models the time-domain excitation sequence which enables straightforward application of the MbG structure to the training process. As a result, the entire model can be stably and easily trained while the perceptual quality of the synthesized speech is significantly improved.

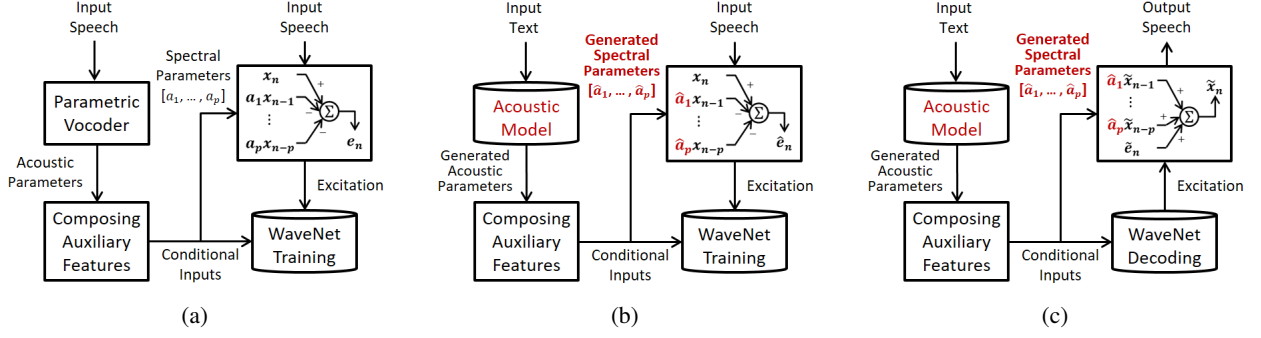


Figure 1: An ExcitNet vocoder for a TTS system: (a) conventional training; (b) proposed MbG training; and (c) synthesis methods.

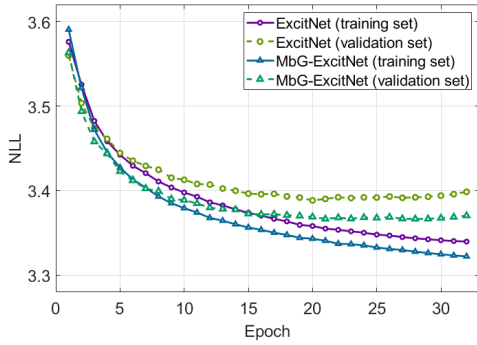


Figure 2: Negative log-likelihood (NLL) obtained during the training process with respect to the plain ExcitNet and MbG-based ExcitNet (MbG-ExcitNet) training methods.

3. ExcitNet TTS systems

3.1. ExcitNet vocoders

The basic WaveNet framework is an autoregressive network which generates a probability distribution of discrete speech symbols from a fixed number of past samples [16]. The ExcitNet vocoder is an advanced version of this network which takes advantages of both the LP vocoder and the WaveNet structure. In an ExcitNet framework, an LP-based adaptive predictor is used to decouple the spectral formant structure from the input speech signal (Fig. 1a). The WaveNet model is then used to train the distribution of the prediction residuals (i.e. excitation) as follows:

$$p(\mathbf{e}|\mathbf{h}) = \prod_{n=1}^N p(e_n|e_1, \dots, e_{n-1}, \mathbf{h}), \quad (1)$$

$$e_n = x_n - \sum_{k=1}^p \alpha_k x_{n-k}, \quad (2)$$

where x_n and e_n denote the n^{th} sample of speech and excitation, respectively; α_k denotes the k^{th} LP coefficient with the order p ; \mathbf{h} denotes the conditional inputs composed of acoustic parameters.

In the speech synthesis step (Fig. 1c), the acoustic parameters of the given input text are generated by a pre-trained acoustic model. These parameters are then used as conditional inputs for the WaveNet model to generate the corresponding time se-

quence of the excitation signal. Finally, the speech signal is reconstructed by passing the generated excitation signal through the LP synthesis filter.

3.2. MbG-structured ExcitNet vocoders

To further improve the quality of the synthesized speech, we propose the incorporation of an MbG structure into the training process of the ExcitNet vocoder. As illustrated in Fig. 1a, conventional vocoding models are trained separately from the acoustic model, even though the generated acoustic parameters, which contain estimation errors, are used as direct conditional inputs (Fig. 1c). This inevitably causes quality degradation of the synthesized speech as the estimation errors from the acoustic model are boosted non-linearly throughout the synthesis process in the vocoder back-end.

Fig. 1b shows the proposed MbG training method which uses closed-loop extraction¹ of the excitation signal. To minimize the mismatch between the training and the generation processes, the LP coefficients in the training step are replaced with those generated by the pre-trained acoustic model as follows:

$$\hat{e}_n = x_n - \sum_{k=1}^p \hat{\alpha}_k x_{n-k}, \quad (3)$$

where $\{\hat{\alpha}_1, \dots, \hat{\alpha}_p\}$ denotes the generated LP coefficients. By combining equations (2) and (3), the excitation sequence can be represented as follows:

$$\hat{e}_n = e_n + e_n^{am}, \quad (4)$$

where e_n^{am} denotes an *intermediate prediction* defined as follows:

$$e_n^{am} = \sum_{k=1}^p (\alpha_k - \hat{\alpha}_k) x_{n-k}. \quad (5)$$

Using the excitation signal (i.e. \hat{e}_n) as the training target means that it becomes possible to guide the model to learn the distributions of the true excitation signal (i.e. e_n) as well as compensate for the acoustic model's estimation errors (i.e. e_n^{am}). Furthermore, because the training and synthesis processes share the same LP coefficients, it is also possible to minimize any mismatch.

¹This extraction method has been adopted in analysis-by-synthesis speech coding frameworks [17,18] where the encoder and decoder share the same quantized filter parameters for minimizing their mismatch conditions.

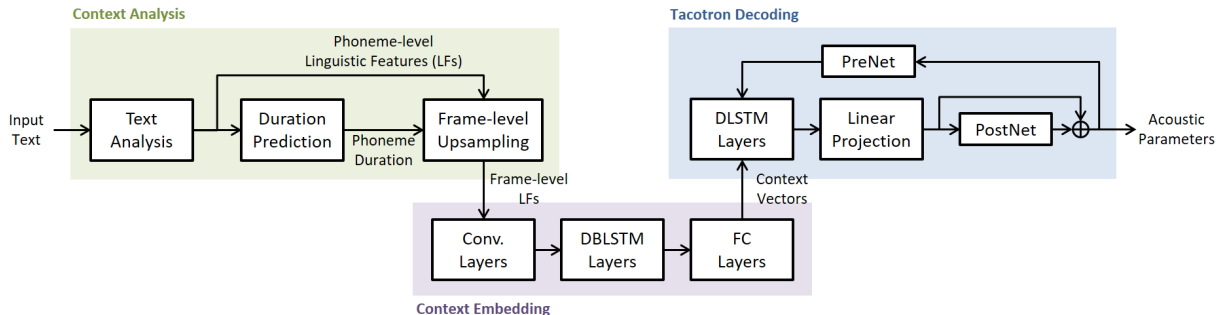


Figure 3: Acoustic model consisting of three sub-modules: context analysis, context embedding, and Tacotron decoding.

The merits of the proposed method are presented in Fig. 2 which shows the negative log-likelihood obtained from the training and validation sets. The proposed MbG-ExcitNet model enables a reduction in both training and validation errors as compared to a plain ExcitNet approach. It is therefore expected that the proposed method will provide more accurate training and generation results, to be further discussed in the following section.

4. Experiments

4.1. Experimental setup

4.1.1. Database

The experiments used a phonetically and prosodically balanced speech corpus recorded by a Korean female professional speaker. The speech signals were sampled at 24 kHz with 16 bit quantization. In total, 4,408 utterances (7.9 hours) were used for training, 230 utterances (0.4 hours) were used for validation, and a further 120 utterances (0.2 hours) were used for testing. The acoustic features were extracted using the improved time-frequency trajectory excitation vocoder at analysis intervals of 5 ms [19], and these features included 40-dimensional line spectral frequencies (LSFs), fundamental frequency (F0), energy, voicing flag (v/uv), 32-dimensional slowly evolving waveform (SEW), and 4-dimensional rapidly evolving waveform (REW), all of which constituted a 79-dimensional feature vector.

4.1.2. Acoustic model

Although there are many state-of-the-art acoustic models available, including Tacotron and Transformer [11, 14, 20], we opted to pursue a Tacotron model with phoneme alignment approach [15] because of its fast and stable generation and competitive synthetic quality. Fig. 3 is a block diagram of the acoustic model which consists of three sub-modules, namely context analysis, context embedding, and Tacotron decoding.

In the context analysis module, the phoneme-level linguistic feature vectors were extracted from the input text. These were composed of 330 binary features for categorical linguistic contexts and 24 features for numerical linguistic contexts. Having input these features, the corresponding phoneme duration was estimated through three fully connected (FC) layers with 1,024, 512, 256 units followed by a unidirectional LSTM network with 128 memory blocks. Based on this estimated duration, the phoneme-level linguistic features were then upsampled to frame-level adding two numerical vectors of phoneme duration and its relative position.

In context embedding, the linguistic features were transformed into high-level context vectors. The module here consisted of three convolution layers with a 101 kernel and 512 channels per layer, a bi-directional LSTM network with 512 memory blocks, and an FC layer with 512 units.

We used a Tacotron 2 decoder network to generate the output acoustic features [11]. First, the previously generated acoustic features were fed into two FC layers with 256 units (i.e. the PreNet), and those features and the vectors from the context embedding module were then passed through two uni-directional LSTM layers with 1,024 memory blocks followed by two projection layers. Finally, to improve generation accuracy, five convolution layers with 5×1 kernels and 512 channels per layer were used as a post-processing network (i.e. the PostNet) to add the residual elements of the generated acoustic features.

Before training, the input and output features were normalized to have zero mean and unit variance. The weights were initialized using *Xavier* initialization and *Adam* optimization was used [21, 22]. The learning rate was scheduled to be decayed from 0.001 to 0.0001 via a decaying rate of 0.33 per 100,000 steps.

4.1.3. Vocoding model

The architecture of the proposed MbG-ExcitNet comprised three convolutional blocks, each with 10 convolution layers with dilations of 1, 2, 4, and so on, up to 512. The numbers of dilated causal convolution channels and 1×1 convolutions in the residual block were both set to 512, and the number of 1×1 convolution channels between the skip connection and the softmax layer was set to 256.

Before training, the LSFs in the training set were generated by the pre-trained acoustic model and were used to compose the conditional inputs together with auxiliary parameters extracted from the input speech, such as F0, energy, v/uv, SEW, and REW. It is possible to use auxiliary parameters also generated by the pre-trained acoustic model, but we recommend to use ground-truth observations to avoid generating unstable speech segments. The conditional inputs were normalized to have zero mean and unit variance and were duplicated from frame to sample to match the length of the input speech signals [3]. The corresponding excitation signals were obtained by passing the input speech signals through the LP analysis filters formed by the generated LSFs. They were then normalized to range between -1.0 and 1.0 followed by 8-bit -law encoding.

The weights were initialized using *Xavier* initialization and *Adam* optimization was used. The learning rate was set to 0.0001, and the batch size was set to 30,000 (1.25 sec).

Table 1: TTS naturalness MOS results with 95% confidence intervals with respect to the different vocoding models: the best MOS scores are in bold.

Index	System	MOS
Test 1	WaveNet	3.23±0.11
Test 2	ExcitNet	4.43±0.08
Test 3	G-WaveNet	3.36±0.11
Test 4	G-ExcitNet	3.29±0.12
Test 5	MbG-ExcitNet (ours)	4.57±0.07
Test 6	Raw	4.66±0.07

4.1.4. TTS system

In the synthesis step, all of the acoustic feature vectors were predicted by the acoustic model with the given input text. By inputting these features, the MbG-ExcitNet vocoder generated a discrete symbol of the quantized excitation signal, and its dynamic was recovered via μ -law expansion. Finally, the speech signal was reconstructed by applying the LP synthesis filter to the generated excitation signal.

4.2. Evaluations

To evaluate the perceptual quality of the proposed system, naturalness MOS tests were performed² by asking 13 native Korean speakers to make quality judgments about the synthesized speech samples using the following five responses: 1 = Bad; 2 = Poor; 3 = Fair; 4 = Good; and 5 = Excellent. In total, 20 utterances were randomly selected from the test set and were synthesized using the different generation models. In particular, the speech samples synthesized by the below conventional vocoding methods were evaluated together to confirm performance differences:

- **WaveNet**: Plain WaveNet vocoder [3]
- **ExcitNet**: Plain ExcitNet vocoder [7]
- **G-WaveNet**: WaveNet vocoder trained with generated acoustic parameters [11]
- **G-ExcitNet**: ExcitNet vocoder trained with generated acoustic parameters

The G-ExcitNet vocoder was configured similarly to the proposed MbG-ExcitNet, but its target excitation was extracted from the ground-truth spectral parameters.

Table 1 presents the MOS test results for the TTS systems with respect to the different vocoding models, and the analysis can be summarized as follows: First, when training vocoding models using ground-truth acoustic parameters, ExcitNet performed better than WaveNet (Tests 1 and 2). This implies that ExcitNet’s adaptive spectral filter is beneficial to reconstruct a more accurate speech signal [7]. Second, training the model with generated parameters provided better perceptual quality than using the ground-truth approach in WaveNet (Tests 1 and 3), but vice versa in ExcitNet (Tests 2 and 4). This result confirms that target excitation should be replaced by considering the acoustic model’s estimation errors in excitation-based methods. Lastly, the proposed MbG-ExcitNet performed best across the different vocoders (Tests 5 and the others). Because the

²Generated audio samples are available at the following URL: <https://sewplay.github.io/demos/mbg-excitenet>

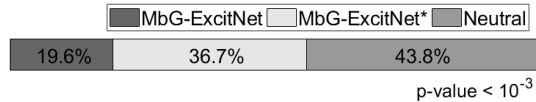


Figure 4: A/B/X preference comparison of MbG-ExcitNet and its initialization-refined version, MbG-ExcitNet*.

MbG training strategy guided the vocoding model to compensate for errors from the acoustic model, it was possible to significantly improve synthesis accuracy. Consequently, the TTS system with the proposed MbG-ExcitNet vocoder achieved 4.57 MOS.

To further verify the effectiveness of the proposed method, we designed additional experiments to refine the initialization of MbG-ExcitNet’s model weights. Since the MbG training process utilizes generated spectral parameters and the corresponding excitation signals as the input conditions and target outputs, respectively, it may be difficult to capture the speech signals’ original characteristics. We therefore adopted a transfer learning method [23] through which the MbG-ExcitNet was initialized by the plain ExcitNet model whose own weights were optimized by ground-truth speech spectra and excitations. All weights were then fine-tuned by the MbG framework. As a result, it was possible to guide the entire training process to learn the characteristics of both the original and the generated speech segments.

Fig. 4 depicts the results of an A/B/X preference test between the proposed MbG-ExcitNet and this initialization-refined version (MbG-ExcitNet*). The setup for this test was the same as for the MOS assessment except that listeners were asked to rate the quality preference of the synthesized speech samples. The results confirm that the initialization-refined system provided better perceptual quality than the originally proposed MbG-ExcitNet. This confirms that adopting a transfer learning method is advantageous to generating more natural speech signal in an MbG-structured TTS system.

5. Conclusions

This paper has proposed a high-quality neural TTS system that incorporates an MbG structure into the ExcitNet vocoder. The MbG-ExcitNet back-end was optimized to learn excitation output distributions while simultaneously compensating for estimation errors from the acoustic model front-end. As such, the proposed method was effective in minimizing the mismatch between the acoustic model and the vocoder. The experimental results verified that a TTS system with the proposed MbG-ExcitNet vocoder performed significantly better than conventional systems with similarly configured WaveNet vocoders. Future research should include extending the framework into speech synthesis systems based on WaveRNN and/or WaveGlow vocoders.

6. Acknowledgements

The authors would like to thank Hyungseob Lim, Kyunguen Byun, Seyun Um, and Suhyeon Oh at DSP&AI Lab., Yonsei University, Seoul, Korea, for their support.

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] 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.
- [3] A. Tamamori, T. Hayashi, K. Kobayashi, K. Takeda, and T. Toda, “Speaker-dependent WaveNet vocoder,” in *Proc. INTERSPEECH*, 2017, pp. 1118–1122.
- [4] N. Kalchbrenner, E. Elsen, K. Simonyan, S. Noury, N. Casagrande, E. Lockhart, F. Stimberg, A. van den Oord, S. Dieleman, and K. Kavukcuoglu, “Efficient neural audio synthesis,” in *Proc. ICML*, 2018, pp. 2410–2419.
- [5] R. Prenger, R. Valle, and B. Catanzaro, “WaveGlow: A flow-based generative network for speech synthesis,” in *Proc. ICASSP*, 2019, pp. 3617–3621.
- [6] L. Juvela, V. Tsiaras, B. Bollepalli, M. Airaksinen, J. Yamagishi, and P. Alku, “Speaker-independent raw waveform model for glottal excitation,” in *Proc. INTERSPEECH*, 2018, pp. 2012–2016.
- [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] M.-J. Hwang, F. Soong, E. Song, X. Wang, H. Kang, and H.-G. Kang, “LP-WaveNet: Linear prediction-based WaveNet speech synthesis,” *arXiv preprint arXiv:1811.11913*, 2018.
- [9] J.-M. Valin and J. Skoglund, “LPCNet: Improving neural speech synthesis through linear prediction,” in *Proc. ICASSP*, 2019, pp. 5891–5895.
- [10] M.-J. Hwang, E. Song, R. Yamamoto, F. Soong, and H.-G. Kang, “Improving LPCNet-based text-to-speech with linear prediction-structured mixture density network,” in *Proc. ICASSP*, 2020, pp. 7219–7223.
- [11] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Skerry-Ryan *et al.*, “Natural TTS synthesis by conditioning WaveNet on Mel spectrogram predictions,” in *Proc. ICASSP*, 2018, pp. 4779–4783.
- [12] L. Juvela, B. Bollepalli, J. Yamagishi, P. Alku *et al.*, “Reducing mismatch in training of DNN-based glottal excitation models in a statistical parametric text-to-speech system,” in *Proc. INTERSPEECH*, 2017, pp. 1368–1372.
- [13] M.-J. Hwang, E. Song, K. Byun, and H.-G. Kang, “Modeling-by-generation-structured noise compensation algorithm for glottal vocoding speech synthesis system,” in *Proc. ICASSP*, 2018, pp. 5669–5673.
- [14] Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, Q. Le, Y. Agiomyrgiannakis, R. Clark, and R. A. Saurous, “Tacotron: Towards end-to-end speech synthesis,” in *Proc. INTERSPEECH*, 2017, pp. 4006–4010.
- [15] T. Okamoto, T. Toda, Y. Shiga, and H. Kawai, “Tacotron-based acoustic model using phoneme alignment for practical neural text-to-speech systems,” in *Proc. ASRU*, 2019, pp. 214–221.
- [16] 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,” *CoRR abs/1609.03499*, 2016.
- [17] B. Atal and M. Schroeder, “Predictive coding of speech signals and subjective error criteria,” *IEEE Trans. Acoust., Speech Signal Process.*, vol. 27, no. 3, pp. 247–254, 1979.
- [18] T. Drugman, P. Alku, A. Alwan, and B. Yegnanarayana, “Glottal source processing: From analysis to applications,” *Comput. Speech Lang.*, vol. 28, no. 5, pp. 1117–1138, 2014.
- [19] 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.
- [20] 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.
- [21] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *Proc. AISTATS*, 2010, pp. 249–256.
- [22] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *CoRR*, vol. abs/1412.6980, 2014. [Online]. Available: <http://arxiv.org/abs/1412.6980>
- [23] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, “Learning and transferring mid-level image representations using convolutional neural networks,” in *Proc. CVPR*, 2014, pp. 1717–1724.