

# NEURAL SOURCE-FILTER-BASED WAVEFORM MODEL FOR STATISTICAL PARAMETRIC SPEECH SYNTHESIS

Xin Wang<sup>1</sup>, Shinji Takaki<sup>1</sup>, Junichi Yamagishi<sup>1\*</sup>

<sup>1</sup>National Institute of Informatics, Japan

wangxin@nii.ac.jp, takaki@nii.ac.jp, jyamagis@nii.ac.jp

## ABSTRACT

Neural waveform models such as the WaveNet are used in many recent text-to-speech systems, but the original WaveNet is quite slow in waveform generation because of its autoregressive (AR) structure. Although faster non-AR models were recently reported, they may be prohibitively complicated due to the use of a distilling training method and the blend of other disparate training criteria. This study proposes a non-AR neural source-filter waveform model that can be directly trained using spectrum-based training criteria and the stochastic gradient descent method. Given the input acoustic features, the proposed model first uses a source module to generate a sine-based excitation signal and then uses a filter module to transform the excitation signal into the output speech waveform. Our experiments demonstrated that the proposed model generated waveforms at least 100 times faster than the AR WaveNet and the quality of its synthetic speech is close to that of speech generated by the AR WaveNet. Ablation test results showed that both the sine-wave excitation signal and the spectrum-based training criteria were essential to the performance of the proposed model.

**Index Terms**— speech synthesis, neural network, waveform modeling

## 1. INTRODUCTION

Text-to-speech (TTS) synthesis, a technology that converts texts into speech waveforms, has been advanced by using end-to-end architectures [1] and neural-network-based waveform models [2, 3, 4]. Among those waveform models, the WaveNet [2] directly models the distributions of waveform sampling points and has demonstrated outstanding performance. The vocoder version of WaveNet [5], which converts the acoustic features into the waveform, also outperformed other vocoders for the pipeline TTS systems [6].

As an autoregressive (AR) model, the WaveNet is quite slow in waveform generation because it has to generate the waveform sampling points one by one. To improve the generation speed, the Parallel WaveNet [3] and the ClariNet [4] introduce a distilling method to transfer ‘knowledge’ from a teacher AR WaveNet to a student non-AR model that simultaneously generates all the waveform sampling points. However, the concatenation of two large models and the mix of distilling and other training criteria reduce the model interpretability and raise the implementation cost.

In this paper, we propose a neural source-filter waveform model that converts acoustic features into speech waveforms. Inspired by classical speech modeling methods [7, 8], we used a source

module to generate a sine-based excitation signal with a specified fundamental frequency (F0). We then used a dilated-convolution-based filter module to transform the sine-based excitation into the speech waveform. The proposed model was trained by minimizing spectral amplitude and phase distances, which can be efficiently implemented using discrete Fourier transforms (DFTs). Because the proposed model is a non-AR model, it generates waveforms much faster than the AR WaveNet. A large-scale listening test showed that the proposed model was close to the AR WaveNet in terms of the Mean opinion score (MOS) on the quality of synthetic speech. An ablation test showed that both the sine-wave excitation and the spectral amplitude distance were crucial to the proposed model.

The model structure and training criteria are explained in Section 2, after which the experiments are described in Section 3. Finally, this paper is summarized and concluded in Section 4.

## 2. PROPOSED MODEL AND TRAINING CRITERIA

### 2.1. Model structure

The proposed model (shown in Figure 1) converts an input acoustic feature sequence  $\mathbf{c}_{1:B}$  of length  $B$  into a speech waveform  $\hat{\mathbf{o}}_{1:T}$  of length  $T$ . It includes a source module that generates an excitation signal  $\mathbf{e}_{1:T}$ , a filter module that transforms  $\mathbf{e}_{1:T}$  into the speech waveform, and a condition module that processes the acoustic features for the source and filter modules. None of the modules takes the previously generated waveform sample as the input. The waveform is assumed to be real-valued, i.e.,  $\hat{o}_t \in \mathbb{R}, 0 < t \leq T$ .

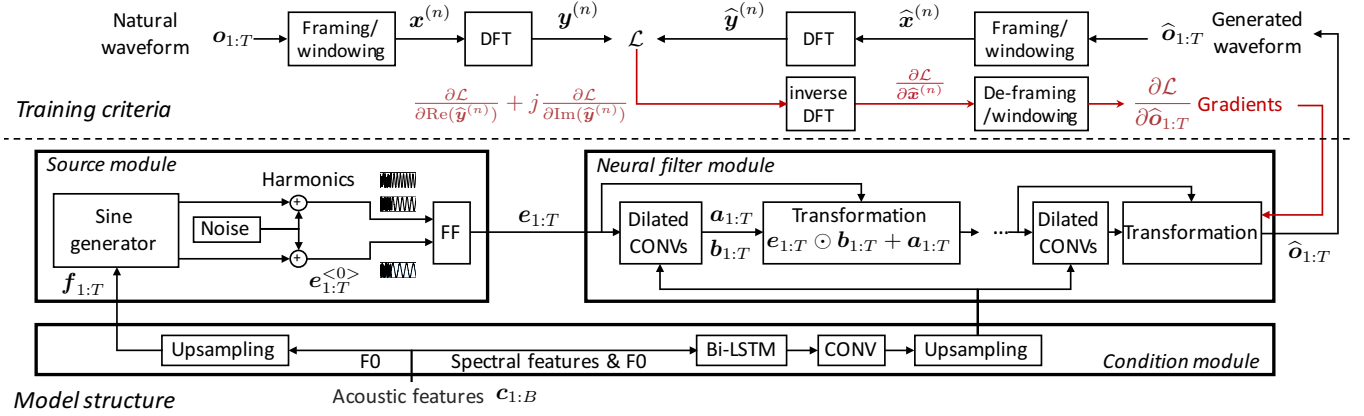
#### 2.1.1. Condition module

The condition module takes as input the acoustic feature sequence  $\mathbf{c}_{1:B} = \{\mathbf{c}_1, \dots, \mathbf{c}_B\}$ , where each  $\mathbf{c}_b = [f_b, \mathbf{s}_b^T]^T$  contains the F0  $f_b$  and the spectral features  $\mathbf{s}_b$  of the  $b$ -th speech frame. The condition module upsamples the F0 by duplicating  $f_b$  to every time step within the  $b$ -th frame and feeds the upsampled F0 sequence  $\mathbf{f}_{1:T}$  to the source module. Meanwhile, it processes  $\mathbf{c}_{1:B}$  using a bi-directional recurrent layer with long-short-term memory (LSTM) units [9] and a convolutional (CONV) layer, after which the processed features are upsampled and sent to the filter module. The LSTM and CONV were used so that the condition module was similar to that of the WaveNet-vocoder [10] in the experiment. They can be replaced with a feedforward layer in practice.

#### 2.1.2. Source module

Given the input F0 sequence  $\mathbf{f}_{1:T}$ , the source module generates a sine-based excitation signal  $\mathbf{e}_{1:T} = \{e_1, \dots, e_T\}$ , where  $e_t \in \mathbb{R}, \forall t \in \{1, \dots, T\}$ . Suppose the F0 value of the  $t$ -th time step

\*This work was partially supported by JST CREST Grant Number JPMJCR18A6, Japan and by MEXT KAKENHI Grant Numbers (16H06302, 16K16096, 17H04687, 18H04120, 18H04112, 18KT0051), Japan.



**Fig. 1.** Structure of proposed model.  $B$  and  $T$  denote lengths of input feature sequence and output waveform, respectively. FF, CONV, and Bi-LSTM denote feedforward, convolutional, and bi-directional recurrent layers, respectively. DFT denotes discrete Fourier transform.

is  $f_t \in \mathbb{R}_{\geq 0}$ , and  $f_t = 0$  denotes being unvoiced. By treating  $f_t$  as the instantaneous frequency [11], a signal  $e_{1:T}^{<0>}$  can be generated as

$$e_t^{<0>} = \begin{cases} \alpha \sin\left(\sum_{k=1}^t 2\pi \frac{f_k}{N_s} + \phi\right) + n_t, & \text{if } f_t > 0 \\ \frac{1}{3\sigma} n_t, & \text{if } f_t = 0 \end{cases}, \quad (1)$$

where  $n_t \sim \mathcal{N}(0, \sigma^2)$  is a Gaussian noise,  $\phi \in [-\pi, \pi]$  is a random initial phase, and  $N_s$  is equal to the waveform sampling rate.

Although we can directly set  $e_{1:T} = e_{1:T}^{<0>}$ , we tried two additional tricks. First, a ‘best’ phase  $\phi^*$  for  $e_{1:T}^{<0>}$  can be determined in the training stage by maximizing the correlation between  $e_{1:T}^{<0>}$  and the natural waveform  $o_{1:T}$ . During generation,  $\phi$  is randomly generated. The second method is to generate harmonics by increasing  $f_k$  in Equation (1) and use a feedforward (FF) layer to merge the harmonics and  $e_{1:T}^{<0>}$  into  $e_{1:T}$ . In this paper we use 7 harmonics and set  $\sigma = 0.003$  and  $\alpha = 0.1$ .

### 2.1.3. Neural filter module

Given the excitation signal  $e_{1:T}$  from the source module and the processed acoustic features from the condition module, the filter module modulates  $e_{1:T}$  using multiple stages of dilated convolution and affine transformations similar to those in ClariNet [4]. For example, the first stage takes  $e_{1:T}$  and the processed acoustic features as input and produces two signals  $a_{1:T}$  and  $b_{1:T}$  using dilated convolution. The  $e_{1:T}$  is then transformed using  $e_{1:T} \odot b_{1:T} + a_{1:T}$ , where  $\odot$  denotes element-wise multiplication. The transformed signal is further processed in the following stages, and the output of the final stage is used as generated waveform  $\hat{o}_{1:T}$ .

The dilated convolution blocks are similar to those in Parallel WaveNet [3]. Specifically, each block contains multiple dilated convolution layers with a filter size of 3. The outputs of the convolution layers are merged with the features from the condition module through gated activation functions [3]. After that, the merged features are transformed into  $a_{1:T}$  and  $b_{1:T}$ . To make sure that  $b_{1:T}$  is positive,  $b_{1:T}$  is parameterized as  $b_{1:T} = \exp(\tilde{b}_{1:T})$ .

Unlike ClariNet or Parallel WaveNet, the proposed model does not use the distilling method. It is unnecessary to compute the mean and standard deviation of the transformed signal. Neither is it necessary to form the convolution and transformation blocks as an inverse autoregressive flow [12].

## 2.2. Training criteria in frequency domain

Because speech perception heavily relies on acoustic cues in the frequency domain, we define training criteria that minimize the spectral amplitude and phase distances, which can be implemented using DFTs. Given these criteria, the proposed model is trained using the stochastic gradient descent (SGD) method.

### 2.2.1. Spectral amplitude distance

Following the convention of short-time Fourier analysis, we conduct waveform framing and windowing before producing the spectrum of each frame. For the generated waveform  $\hat{o}_{1:T}$ , we use  $\hat{x}^{(n)} = [\hat{x}_1^{(n)}, \dots, \hat{x}_M^{(n)}]^\top \in \mathbb{R}^M$  to denote the  $n$ -th waveform frame of length  $M$ . We then use  $\hat{y}^{(n)} = [\hat{y}_1^{(n)}, \dots, \hat{y}_K^{(n)}]^\top \in \mathbb{C}^K$  to denote the spectrum of  $\hat{x}^{(n)}$  calculated using  $K$ -point DFT. We similarly define  $x^{(n)}$  and  $y^{(n)}$  for the natural waveform  $o_{1:T}$ .

Suppose the waveform is sliced into  $N$  frames. Then the log spectral amplitude distance is defined as follows:

$$\mathcal{L}_s = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \left[ \log \frac{\text{Re}(y_k^{(n)})^2 + \text{Im}(y_k^{(n)})^2}{\text{Re}(\hat{y}_k^{(n)})^2 + \text{Im}(\hat{y}_k^{(n)})^2} \right]^2, \quad (2)$$

where  $\text{Re}(\cdot)$  and  $\text{Im}(\cdot)$  denote the real and imaginary parts of a complex number, respectively.

Although  $\mathcal{L}_s$  is defined on complex-valued spectra, the gradient  $\frac{\partial \mathcal{L}_s}{\partial \hat{o}_{1:T}} \in \mathbb{R}^T$  for SGD training can be efficiently calculated. Let us consider the  $n$ -th frame and compose a complex-valued vector  $g^{(n)} = \frac{\partial \mathcal{L}_s}{\partial \text{Re}(\hat{y}^{(n)})} + j \frac{\partial \mathcal{L}_s}{\partial \text{Im}(\hat{y}^{(n)})} \in \mathbb{C}^K$ , where the  $k$ -th element is  $g_k^{(n)} = \frac{\partial \mathcal{L}_s}{\partial \text{Re}(\hat{y}_k^{(n)})} + j \frac{\partial \mathcal{L}_s}{\partial \text{Im}(\hat{y}_k^{(n)})} \in \mathbb{C}$ . It can be shown that, as long as  $g^{(n)}$  is Hermitian symmetric, the inverse DFT of  $g^{(n)}$  is equal to  $\frac{\partial \mathcal{L}_s}{\partial \hat{x}^{(n)}} = [\frac{\partial \mathcal{L}_s}{\partial \hat{x}_1^{(n)}}, \dots, \frac{\partial \mathcal{L}_s}{\partial \hat{x}_m^{(n)}}, \frac{\partial \mathcal{L}_s}{\partial \hat{x}_M^{(n)}}] \in \mathbb{R}^{M \times 1}$ . Using the same method,  $\frac{\partial \mathcal{L}_s}{\partial \hat{x}^{(n)}}$  for  $n \in \{1, \dots, N\}$  can be computed in parallel. Given  $\{\frac{\partial \mathcal{L}_s}{\partial \hat{x}^{(1)}}, \dots, \frac{\partial \mathcal{L}_s}{\partial \hat{x}^{(N)}}\}$ , the value of each  $\frac{\partial \mathcal{L}_s}{\partial \hat{o}_t}$  in  $\frac{\partial \mathcal{L}_s}{\partial \hat{o}_{1:T}}$  can be easily accumulated since the relationship between  $\hat{o}_t$  and each  $\hat{x}_m^{(n)}$  has been determined by the framing and windowing operations.

<sup>1</sup>In the implementation using fast Fourier transform,  $\hat{x}^{(n)}$  of length  $M$  is zero-padded to length  $K$  before DFT. Accordingly, the inverse DFT of  $g^{(n)}$  also gives the gradients w.r.t. the zero-padded part, which should be discarded (see <https://arxiv.org/abs/1810.11946>).

**Table 1.** Three framing and DFT configurations for  $\mathcal{L}_s$  and  $\mathcal{L}_p$ 

	$\mathcal{L}_{s1} \& \mathcal{L}_{p1}$	$\mathcal{L}_{s2} \& \mathcal{L}_{p2}$	$\mathcal{L}_{s3} \& \mathcal{L}_{p3}$
DFT bins $K$	512	128	2048
Frame length $M$	320 (20 ms)	80 (5 ms)	1920 (120 ms)
Frame shift	80 (5 ms)	40 (2.5 ms)	640 (40 ms)

Note: all configurations use Hann window.

In fact,  $\frac{\partial \mathcal{L}_s}{\partial \mathbf{o}_{1:T}} \in \mathbb{R}^T$  can be calculated in the same manner no matter how we set the framing and DFT configuration, i.e., the values of  $N$ ,  $M$ , and  $K$ . Furthermore, multiple  $\mathcal{L}_s$ s with different configurations can be computed, and the gradients  $\frac{\partial \mathcal{L}_s}{\partial \mathbf{o}_{1:T}}$  can be simply summed up. For example, using the three  $\mathcal{L}_s$ s in Table 1 was found to be essential to the proposed model (see Section 3.3).

The Hermitian symmetry of  $\mathbf{g}^{(n)}$  is satisfied if  $\mathcal{L}_s$  is carefully defined. For example,  $\mathcal{L}_s$  can be the square error or Kullback-Leibler divergence (KLD) of the spectral amplitudes [13, 14]. The phase distance defined below also satisfies the requirement.

### 2.2.2. Phase distance

Given the spectra, a phase distance [15] is computed as

$$\mathcal{L}_p = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \left| 1 - \exp(j(\hat{\theta}_k^{(n)} - \theta_k^{(n)})) \right|^2$$

$$= \sum_{n=1}^N \sum_{k=1}^K \left[ 1 - \frac{\text{Re}(\hat{y}_k^{(n)})\text{Re}(y_k^{(n)}) + \text{Im}(\hat{y}_k^{(n)})\text{Im}(y_k^{(n)})}{|\hat{y}_k^{(n)}||y_k^{(n)}|} \right], \quad (3)$$

where  $\hat{\theta}_k^{(n)}$  and  $\theta_k^{(n)}$  are the phases of  $\hat{y}_k^{(n)}$  and  $y_k^{(n)}$ , respectively. The gradient  $\frac{\partial \mathcal{L}_p}{\partial \mathbf{o}_{1:T}}$  can be computed by the same procedure as  $\frac{\partial \mathcal{L}_s}{\partial \mathbf{o}_{1:T}}$ . Multiple  $\mathcal{L}_p$ s and  $\mathcal{L}_s$ s with different framing and DFT configurations can be added up as the ultimate training criterion  $\mathcal{L}$ . For different  $\mathcal{L}_{*s}$ , additional DFT/iDFT and framing/windowing blocks should be added to the model in Figure 1.

## 3. EXPERIMENTS

### 3.1. Corpus and features

This study used the same Japanese speech corpus and data division recipe as our previous study [16]. This corpus [17] contains neutral reading speech uttered by a female speaker. Both validation and test sets contain 480 randomly selected utterances. Among the 48-hour training data, 9,000 randomly selected utterances (15 hours) were used as the training set in this study. For the ablation test in Section 3.3, the training set was further reduced to 3,000 utterances (5 hours). Acoustic features, including 60 dimensions of Mel-generalized cepstral coefficients (MGCs) [18] and 1 dimension of F0, were extracted from the 48 kHz waveforms at a frame shift of 5 ms using WORLD [19]. The natural waveforms were then downsampled to 16 kHz for model training and the listening test.

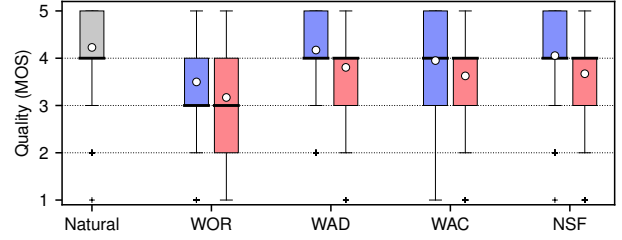
### 3.2. Comparison of proposed model, WaveNet, and WORLD

The first experiment compared the four models listed in Table 2<sup>2</sup>. The WAD model, which was trained in our previous study [6],

<sup>2</sup>The models were implemented using a modified CURRENNT toolkit [20] on a single P100 Nvidia GPU card. Codes, recipes, and generated speech can be found on <https://nii-yamagishilab.github.io>.

**Table 2.** Models for comparison test in Section 3.2

WOR	WORLD vocoder
WAD	WaveNet-vocoder for 10-bit discrete $\mu$ -law waveform
WAC	WaveNet-vocoder using Gaussian dist. for raw waveform
NSF	Proposed model for raw waveform



**Fig. 2.** MOS scores of natural speech, synthetic speech given natural acoustic features (blue), and synthetic speech given acoustic features generated from acoustic models (red). White dots are mean values.

**Table 3.** Average number of waveform points generated in 1 s

WAD	NSF (memory-save mode)	NSF (normal mode)
0.19k	20k	227k

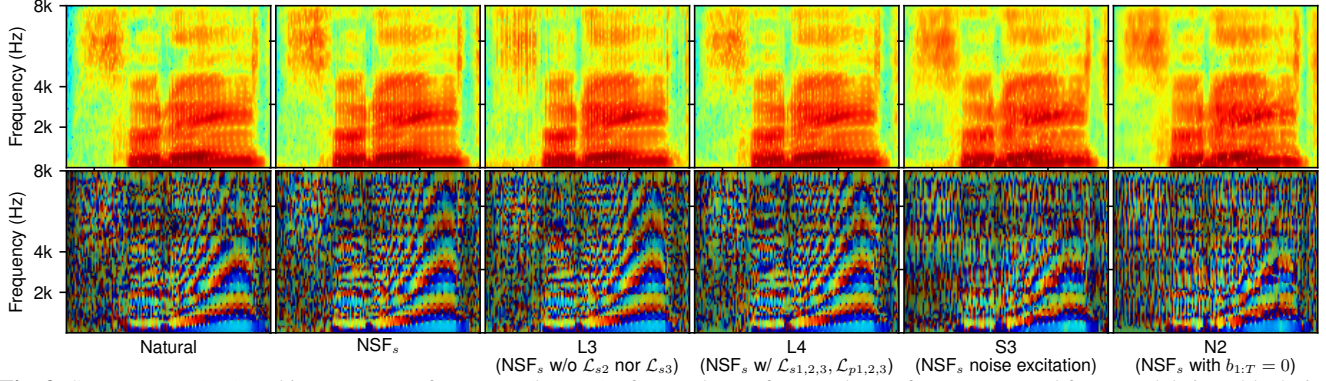
contained a condition module, a post-processing module, and 40 dilated CONV blocks, where the  $k$ -th CONV block had a dilation size of  $2^{\text{modulo}(k,10)}$ . WAC was similar to WAD but used a Gaussian distribution to model the raw waveform at the output layer [4].

The proposed NSF contained 5 stages of dilated CONV and transformation, each stage including 10 convolutional layers with a dilation size of  $2^{\text{modulo}(k,10)}$  and a filter size of 3. Its condition module was the same as that of WAD and WAC. NSF was trained using  $\mathcal{L} = \mathcal{L}_{s1} + \mathcal{L}_{s2} + \mathcal{L}_{s3}$ , and the configuration of each  $\mathcal{L}_{s*}$  is listed in Table 1. The phase distance  $\mathcal{L}_{p*}$  was not used in this test.

Each model generated waveforms using natural and generated acoustic features, where the generated acoustic features were produced by the acoustic models in our previous study [6]. The generated and natural waveforms were then evaluated by paid native Japanese speakers. In each evaluation round the evaluator listened to one speech waveform in each screen and rated the speech quality on a 1-to-5 MOS scale. The evaluator can take at most 10 evaluation rounds and can replay the sample during evaluation. The waveforms in an evaluation round were for the same text and were played in a random order. Note that the waveforms generated from NSF and WAC were converted to 16-bit PCM format before evaluation.

A total of 245 evaluators conducted 1444 valid evaluation rounds in all, and the results are plotted in Figure 2. Two-sided Mann-Whitney tests showed that the difference between any pair of models is statistically significant ( $p < 0.01$ ) except NSF and WAC when the two models used generated acoustic features. In general, NSF outperformed WOR and WAC but performed slightly worse than WAD. The gap of the mean MOS scores between NSF and WAD was about 0.12, given either natural or generated acoustic features. A possible reason for this result may be the difference between the non-AR and AR model structures, which is similar to the difference between the finite and infinite impulse response filters. WAC performed worse than WAD because some syllables were perceived to be trembling in pitch, which may be caused by the random sampling generation method. WAD alleviated this artifact by using a one-best generation method in voiced regions [6].

After the MOS test, we compared the waveform generation speed of NSF and WAD. The implementation of NSF has a normal



**Fig. 3.** Spectrogram (top) and instantaneous frequency (bottom) of natural waveform and waveforms generated from models in Table 4 given natural acoustic features in test set (utterance AOZORAR\_03372.T01). Figures are plotted using 5 ms frame length and 2.5 ms frame shift.

**Table 4.** Models for ablation test (Section 3.3)

NSF <sub>s</sub>	NSF trained on 5-hour data
L1	NSF <sub>s</sub> without using $\mathcal{L}_{s3}$ (i.e., $\mathcal{L} = \mathcal{L}_{s1} + \mathcal{L}_{s2}$ )
L2	NSF <sub>s</sub> without using $\mathcal{L}_{s2}$ (i.e., $\mathcal{L} = \mathcal{L}_{s1} + \mathcal{L}_{s3}$ )
L3	NSF <sub>s</sub> without using $\mathcal{L}_{s2}$ nor $\mathcal{L}_{s3}$ (i.e., $\mathcal{L} = \mathcal{L}_{s1}$ )
L4	NSF <sub>s</sub> using $\mathcal{L} = \mathcal{L}_{s1} + \mathcal{L}_{s2} + \mathcal{L}_{s3} + \mathcal{L}_{p1} + \mathcal{L}_{p2} + \mathcal{L}_{p3}$
L5	NSF <sub>s</sub> using KLD of spectral amplitudes
S1	NSF <sub>s</sub> without harmonics
S2	NSF <sub>s</sub> without harmonics or ‘best’ phase $\phi^*$
S3	NSF <sub>s</sub> only using noise as excitation
N1	NSF <sub>s</sub> with $b_{1:T} = 1$ in filter’s transformation layers
N2	NSF <sub>s</sub> with $b_{1:T} = 0$ in filter’s transformation layers

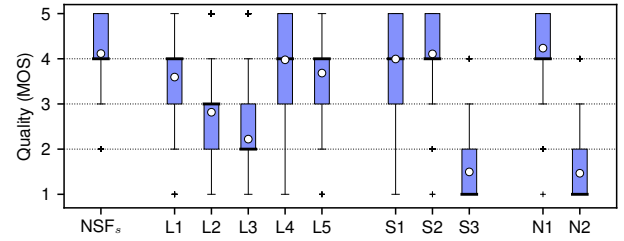
generation mode and a memory-save one. The normal mode allocates all the required GPU memory once but cannot generate waveforms longer than 6 seconds because of the insufficient memory space in a single GPU card. The memory-save mode can generate long waveforms because it releases and allocates the memory layer by layer, but the repeated memory operations are time consuming.

We evaluated NSF using both modes on a smaller test set, in which each of the 80 generated test utterances was around 5 seconds long. As the results in Table 3 show, NSF is much faster than WAD. Note that WAD allocates and re-uses a small size of GPU memory, which needs no repeated memory operation. WAD is slow mainly because of the AR generation process. Of course, both WAD and NSF can be improved if our toolkit is further optimized. Particularly, if the memory operation can be sped up, the memory-save mode of NSF will be much faster.

### 3.3. Ablation test on proposed model

This experiment was an ablation test on NSF. Specifically, the 11 variants of NSF listed in Table 4 were trained using the 5-hour training set. For a fair comparison, NSF was re-trained using the 5-hour data, and this variant is referred to as NSF<sub>s</sub>. The speech waveforms were generated given the natural acoustic features and rated in 1444 evaluation rounds by the same group of evaluators in Section 3.2. This test excluded natural waveform for evaluation.

The results are plotted in Figure 4. The difference between NSF<sub>s</sub> and any other model except S2 was statistically significant ( $p < 0.01$ ). Comparison among NSF<sub>s</sub>, L1, L2, and L3 shows that using multiple  $\mathcal{L}_s$ s listed in Table 1 is beneficial. For L3 that used only  $\mathcal{L}_{s1}$ , the generated waveform points clustered around one peak in each frame, and the waveform suffered from a pulse-train noise. This can be observed from L3 of Figure 3, whose spectrogram in the



**Fig. 4.** MOS scores of synthetic samples from NSF<sub>s</sub> and its variants given natural acoustic features. White dots are mean MOS scores.

high frequency band shows more clearly vertical strips than other models. Accordingly, this artifact can be alleviated by adding  $\mathcal{L}_{s2}$  with a frame length of 5 ms for model training, which explained the improvement in L1. Using phase distance (L4) didn’t improve the speech quality even though the value of the phase distance was consistently decreased on both training and validation data.

The good result of S2 indicates that a single sine-wave excitation with a random initial phase also works. Without the sine-wave excitation, S3 generated waveforms that were intelligible but lacked stable harmonic structure. N1 slightly outperformed NSF<sub>s</sub> while N2 produced unstable harmonic structures. Because the transformation in N1 is equivalent to skip-connection [21], the result indicates that the skip-connection may help the model training.

## 4. CONCLUSION

In this paper, we proposed a neural waveform model with separated source and filter modules. The source module produces a sine-wave excitation signal with a specified F0, and the filter module uses dilated convolution to transform the excitation into a waveform. Our experiment demonstrated that the sine-wave excitation was essential for generating waveforms with harmonic structures. We also found that multiple spectral-based training criteria and the transformation in the filter module contributed to the performance of the proposed model. Compared with the AR WaveNet, the proposed model generated speech with a similar quality at a much faster speed.

The proposed model can be improved in many aspects. For example, it is possible to simplify the dilated convolution blocks. It is also possible to try classical speech modeling methods, including glottal waveform excitations [22, 23], two-bands or multi-bands approaches [24, 25] on waveforms. When applying the model to convert linguistic features into the waveform, we observed the over-smoothing affect in the high-frequency band and will investigate the issue in the future work.

## 5. REFERENCES

- [1] Jonathan Shen, Ruoming Pang, Ron J Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, et al., “Natural TTS synthesis by conditioning WaveNet on Mel spectrogram predictions,” in *Proc. ICASSP*, 2018, pp. 4779–4783.
- [2] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu, “WaveNet: A generative model for raw audio,” *arXiv preprint arXiv:1609.03499*, 2016.
- [3] Aaron van den Oord, Yazhe Li, Igor Babuschkin, Karen Simonyan, Oriol Vinyals, Koray Kavukcuoglu, George van den Driessche, Edward Lockhart, Luis Cobo, Florian Stimberg, Norman Casagrande, Dominik Grewe, Seb Noury, Sander Dieleman, Erich Elsen, Nal Kalchbrenner, Heiga Zen, Alex Graves, Helen King, Tom Walters, Dan Belov, and Demis Hassabis, “Parallel WaveNet: Fast high-fidelity speech synthesis,” in *Proc. ICML*, 2018, pp. 3918–3926.
- [4] Wei Ping, Kainan Peng, and Jitong Chen, “Clarinet: Parallel wave generation in end-to-end text-to-speech,” in *Proc. ICLP*, 2019.
- [5] Akira Tamamori, Tomoki Hayashi, Kazuhiro Kobayashi, Kazuya Takeda, and Tomoki Toda, “Speaker-dependent WaveNet vocoder,” in *Proc. Interspeech*, 2017, pp. 1118–1122.
- [6] Xin Wang, Jaime Lorenzo-Trueba, Shinji Takaki, Lauri Juvela, and Junichi Yamagishi, “A comparison of recent waveform generation and acoustic modeling methods for neural-network-based speech synthesis,” in *Proc. ICASSP*, 2018, pp. 4804–4808.
- [7] Per Hedelin, “A tone oriented voice excited vocoder,” in *Proc. ICASSP*. IEEE, 1981, vol. 6, pp. 205–208.
- [8] Robert McAulay and Thomas Quatieri, “Speech analysis/synthesis based on a sinusoidal representation,” *IEEE Trans. on Acoustics, Speech, and Signal Processing*, vol. 34, no. 4, pp. 744–754, 1986.
- [9] Alex Graves, *Supervised Sequence Labelling with Recurrent Neural Networks*, Ph.D. thesis, Technische Universität München, 2008.
- [10] Xin Wang, Shinji Takaki, and Junichi Yamagishi, “Investigation of WaveNet for text-to-speech synthesis,” Tech. Rep. 6, SIG Technical Reports, feb 2018.
- [11] John R Carson and Thornton C Fry, “Variable frequency electric circuit theory with application to the theory of frequency-modulation,” *Bell System Technical Journal*, vol. 16, no. 4, pp. 513–540, 1937.
- [12] Diederik P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling, “Improved variational inference with inverse autoregressive flow,” in *Proc. NIPS*, 2016, pp. 4743–4751.
- [13] Daniel D Lee and H Sebastian Seung, “Algorithms for non-negative matrix factorization,” in *Proc. NIPS*, 2001, pp. 556–562.
- [14] Shinji Takaki, Hirokazu Kameoka, and Junichi Yamagishi, “Direct modeling of frequency spectra and waveform generation based on phase recovery for DNN-based speech synthesis,” in *Proc. Interspeech*, 2017, pp. 1128–1132.
- [15] Shinji Takaki, Toru Nakashika, Xin Wang, and Junichi Yamagishi, “STFT spectral loss for training a neural speech waveform model,” in *Proc. ICASSP*, 2019, p. (accepted).
- [16] Hieu-Thi Luong, Xin Wang, Junichi Yamagishi, and Nobuyuki Nishizawa, “Investigating accuracy of pitch-accent annotations in neural-network-based speech synthesis and denoising effects,” in *Proc. Interspeech*, 2018, pp. 37–41.
- [17] Hisashi Kawai, Tomoki Toda, Jinfu Ni, Minoru Tsuzaki, and Keiichi Tokuda, “XIMERA: A new TTS from ATR based on corpus-based technologies,” in *Proc. SSW5*, 2004, pp. 179–184.
- [18] Keiichi Tokuda, Takao Kobayashi, Takashi Masuko, and Satoshi Imai, “Mel-generalized cepstral analysis a unified approach,” in *Proc. ICSLP*, 1994, pp. 1043–1046.
- [19] Masanori Morise, Fumiya Yokomori, and Kenji Ozawa, “WORLD: A vocoder-based high-quality speech synthesis system for real-time applications,” *IEICE Trans. on Information and Systems*, vol. 99, no. 7, pp. 1877–1884, 2016.
- [20] Felix Weninger, Johannes Bergmann, and Björn Schuller, “Introducing CURRENT: The Munich open-source CUDA recurrent neural network toolkit,” *The Journal of Machine Learning Research*, vol. 16, no. 1, pp. 547–551, 2015.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *Proc. CVPR*, 2016, pp. 770–778.
- [22] Gunnar Fant, Johan Liljencrants, and Qi-guang Lin, “A four-parameter model of glottal flow,” *STL-QPSR*, vol. 4, no. 1985, pp. 1–13, 1985.
- [23] Lauri Juvela, Bajibabu Bollepalli, Manu Airaksinen, and Paavo Alku, “High-pitched excitation generation for glottal vocoding in statistical parametric speech synthesis using a deep neural network,” in *Proc. ICASSP*, 2016, pp. 5120–5124.
- [24] John Makhoul, R Viswanathan, Richard Schwartz, and AWF Huggins, “A mixed-source model for speech compression and synthesis,” *The Journal of the Acoustical Society of America*, vol. 64, no. 6, pp. 1577–1581, 1978.
- [25] D. W. Griffin and J. S. Lim, “Multiband excitation vocoder,” *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 36, no. 8, pp. 1223–1235, Aug 1988.