
A Fully Time-domain Neural Model for Subband-based Speech Synthesizer

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper introduces a deep neural network model for subband-based speech
2 synthesizer. The model benefits from the short bandwidth of the subband signals
3 to reduce the complexity of the time-domain speech generator. We employed
4 the multi-level wavelet analysis/synthesis to decompose/reconstruct the signal to
5 subbands in time domain. Inspired from the WaveNet, a convolutional neural
6 network (CNN) model predicts subband speech signals fully in time domain. Due
7 to the short bandwidth of the subbands, a simple network architecture is enough to
8 train the simple patterns of the subbands accurately. In the ground truth experiments
9 with teacher forcing, the subband synthesizer outperforms the fullband model
10 significantly. In addition, by conditioning the model on the phoneme sequence
11 using a pronunciation dictionary, we have achieved the first fully time-domain
12 neural text-to-speech (TTS) system. The generated speech of the subband TTS
13 shows comparable quality as the fullband one with a slighter network architecture
14 for each subband.

15 1 Introduction

16 Text-to-speech (TTS) synthesizers have been vital assistants of disabled persons, new language
17 learners, as well as a wide range of human-computer interactions for decades. Researchers have
18 presented various techniques starting from concatenative synthesis [1], [2] to statistical parametric
19 speech synthesis [3]–[5], either based on hidden Markov model (HMM) or deep neural network
20 (DNN), and eventually end-to-end fully neural network based models [6], [7].

21 Recent speech synthesizers have employed giant neural networks and high configuration GPUs to
22 achieve remarkable success in more natural and fast speech generation. Of such models, WaveNet
23 [8] has achieved the most natural generated speech that significantly closes the gap with human. As
24 a deep generative network, WaveNet directly models the raw audio waveform, which has changed
25 the existing paradigms. The model is applicable for every audio such as speech and music. It
26 made a paradigm to absorb a tremendous amount of attention for sequential modeling [9], speech
27 enhancement [10], [11], and vocoder, which is the wave synthesizer from acoustic features [12]–[15].
28 Furthermore, the state-of-the-art TTS, Tacotron 2 [16] benefits from the WaveNet as the back-end
29 vocoder for transforming the spectrogram as acoustic features to the waveform.

30 Thanks to its convolutional structure, WaveNet benefits from parallel computing in train. However,
31 the generation is still a sequential sample-by-sample process. Thus, due to the very high temporal
32 resolution of speech signals (at least 16000 samples per second), the vanilla WaveNet suffers from the
33 long generation time. Therefore, fast [17] and parallel [18] models are introduced. The fast model is
34 an efficient implementation that removes redundant convolutional operations by caching them. While
35 the parallel model utilized a new method, named probability density distillation, which leads to the
36 speech synthesis faster than real-time.

Table 1: The list of symbols and notations used in this paper

Symbol	Description
$s(t)$	The fullband speech signal (\hat{s} is the estimation)
$s_l(t)$	The l^{th} subband obtained from the l^{th} level of the wavelet transform
c	Conditional features
h	Latent variables
x	Previous clean (generated) subband samples in train (test)
k	The dilation layer index, $k = 1, \dots, K$

Table 2: Applications of the model with different conditional features

Latent features	Application (if x is speech)
None	Speech-like wave generator
Speaker ID	Speech-like wave generator with the speaker’s voice
Acoustic features (like f_0 , MFCC)	Vocoder
Linguistic features	Text-to-Speech synthesizer

37 Unlike the huge network hired in the parallel model, some studies benefit from subband decomposition
 38 to reduce the complexity. Previously, a hybrid TTS [19] applied HMM-based and waveform-based
 39 synthesis for low and high frequencies, respectively. However, the TTS suffers from the drawbacks
 40 of the HMM-based models and the overall performance is not satisfying. In addition, a subband
 41 WaveNet vocoder [20] is presented using a frequency filterbank analysis. However, to have a TTS
 42 based on the subband vocoder, separate acoustic and linguistic models are required.

43 Similar to [20], the aim of this research is to break down the WaveNet architecture into smaller
 44 networks for each subband of the speech signal. The benefits of the subband model is the reduced
 45 computational complexity and the feasibility of training accurately for each subband due to their short
 46 bandwidth. In addition, The similar morphological structure of the dilated convolutions in WaveNet
 47 and the wavelet transform has inspired us to use the wavelet. Thus, the innovation is utilizing the
 48 wavelet analysis to decompose the time-domain speech signal $s(t)$ into subbands $s_l(t) (l = 1, \dots, L)$.
 49 Then, an integrated model generates each subband signal in parallel. The subband signal generator
 50 is based on the fast WaveNet [17]. Our wavelet decomposition seems to be more accurate for the
 51 reconstruction in time domain compared to the frequency domain filterbank used in [19] and [20].

52 Even though many recent studies utilized the WaveNet as a vocoder, we believe that converting the
 53 spectrogram information to waveform is an inverse spectrogram process and may not necessarily
 54 need such a huge architecture. Instead, our hypothesis is that the WaveNet is able to perform some
 55 parts of the TTS front stage, as well. In addition, a single integrated model is likely to be more
 56 stable than a multi-stage model [6], [21]. Hence, another contribution of this paper is that by simply
 57 conditioning the proposed model on the phoneme sequence and benefitting from an encoder, we have
 58 achieved the first fully time-domain neural TTS.

59 Table 1 reports the list of symbols and notations used in this paper. Section 2 describes the proposed
 60 subband speech synthesizer. Section 3 explains our experiments and results. Finally, conclusion
 61 comes in Section 4.

62 2 Proposed subband speech synthesizer

63 The aim of this paper is to reduce the complexity of the time-domain TTS by decomposing the
 64 fullband speech signal s into the subbands $s_l (l = 1, \dots, L)$ using the wavelet analysis. Benefiting
 65 from the parallel processing, our designed model estimates the subband signals based on conditional
 66 features. Due to the short bandwidth of the subbands, the structure of the subband generator can be
 67 much slighter than the fullband one. Our hypothesis is that estimations can be more accurate because
 68 subband generators are trained for the localized frequency patterns. When the subband signals are
 69 generated according to the corresponding conditional features using the localized TTS, then the
 70 wavelet synthesis reconstructs the fullband signal. Details of the wavelet transform is described in
 71 Subsection 2.1.

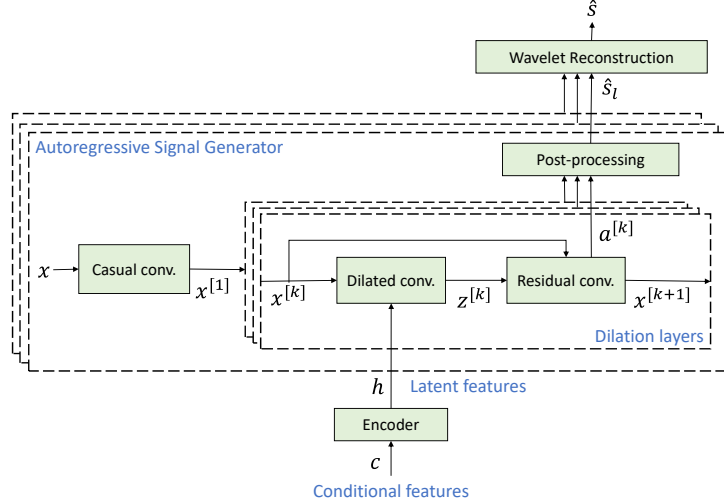


Figure 1: Schematic diagram of the proposed time-domain subband-based speech synthesizer. The model is trained to estimate subband signals s_l conditioning on the latent variable h extracted from c and the previous time samples of the subband signal x . Linguistic and acoustic features can feed to the model as the conditional features for TTS.

72 Figure 1 depicts the architecture of the proposed subband-based time-domain speech synthesizer. In
 73 the designed model, given the conditional features c , an encoder extracts the latent variables h for
 74 generating samples conditioning on them (detailed in Subsection 2.2). Table 2 explains applications
 75 of the model with different latent features. According to the table, if the conditional features are
 76 linguistic features such as character or phoneme sequence, then the latent features would be linguistic
 77 features to make the model as a TTS.

78 The main part of the model is the autoregressive signal generators, shown by the outer dashed blocks
 79 in Figure 1. Each generator is in charge of modeling the probability distribution of each subband.
 80 The subband generator has similar structure as the WaveNet. Subsection 2.3 explains details of the
 81 autoregressive signal generator.

82 In the training phase, the loss is defined by summation of the subband losses, which is the cross-
 83 entropy of the estimated and target subband signal, as

$$loss = - \sum_{l=1}^L \mathbb{E}_{p_l} [\log q_l], \quad (1)$$

84 in which p_l and q_l are the probability distribution of $s_l(t)$ and $\hat{s}_l(t)$, respectively. Since it is a
 85 probabilistic model, the generation model estimates the t^{th} sample of each subband by sampling the
 86 learned probability distribution.

87 2.1 Subband decomposition/reconstruction

88 A set of analysis filters can decompose speech signal $s(t)$ into subbands $s_l(t)$, and their paired
 89 synthesis filters are able to reconstruct back the fullband signal. The proposed synthesizer utilizes
 90 multi-level orthogonal time-domain wavelet as follows,

$$\begin{cases} u_l(t) = u_{l-1}(t) * \varphi_l(t) \\ s_l(t) = u_{l-1}(t) * \psi_l(t) \end{cases} \quad (2)$$

91 where $\varphi_l(t)$ and $\psi_l(t)$ are Daubechies scaling (low-pass) and mother wavelet (high-pass) functions
 92 [22], respectively. Moreover, $l = 1, \dots, L$ refers to the wavelet level and $u_0(t) = s(t)$. The
 93 downsampling is omitted in every level of the wavelet transform because the downsampling widens
 94 the bandwidth, which needs more complex model for training. In addition, it decreases the size of the

95 dataset. Since there is no data like more data for the training, we ignored the downsampling after
96 each layer.

97 Reasons for selecting the wavelet transform rather than the short time Fourier transform (STFT)
98 filterbank are as follows. First, the wavelet transform is very robust for reconstruction [23]. Corruption
99 of the wavelet coefficients will only affect the reconstructed signal locally near the perturbed position,
100 while the STFT will spread out the error everywhere in time. Second, output of the Fourier analysis
101 filters are complex. Most of the spectrogram-based speech synthesizers ignore modeling the phase
102 spectrogram [6], while the Fourier synthesis filters are sensitive to phase errors. Therefore, compared
103 to the wavelet, the STFT models are unable to reconstruct the phase correctly. Third, the logarithmic
104 spectral resolution of the wavelet are more compatible with the nature of speech compared to the
105 uniform tiling of the spectrogram. Due to the nonlinear bandwidth divisions of the wavelet, high
106 frequencies (e.g. above 4 kHz for 16 kHz sampling rate) fall in one subband. Whereas, there are fine
107 divisions for the low frequencies.

108 Later in the experiments, we will see the signal-to-noise ratio (SNR) of the consecutive decomposition
109 and reconstruction is about 41 dB, in which the noise is hardly sensible by the human ear.

110 2.2 Conditional/latent features

111 A variety of conditional features can be fed to the model. Table 2 gives some examples. Of such
112 features, we use phoneme sequence produced by a text normalization and lexicon to have a TTS
113 model. The phoneme sequence speeds up the training [24]. As shown in Figure 1 by the encoder
114 block, a number of convolutional layers along time axis can extract the linguistic features implicitly.
115 The activation of the last layer, denoted by latent features h , is used for the generators. In fact, the
116 encoder plays the role of the linguistic model for TTS.

117 2.3 Subband autoregressive signal generator

118 The subband generator has a similar architecture as the WaveNet. Unlike the WaveNet, our au-
119 toregressive signal generator is in charge of generating subband signals. The model estimates the
120 posterior probability of each subband time-sample x_t conditioned on the previous samples, $x_{<t}$ and
121 some latent features h_t as $p(x_t|x_{<t}, h_t)$.

122 As shown in Figure 1, each generator contains:

- 123 • a causal convolution layer as the preprocessing,
- 124 • dilation area, which is illustrated by the inner dashed blocks in the figure, and
- 125 • post-processing.

126 As an input of the generator, x refers to the previous clean samples of each subband s_l for training.
127 Similarly, in generation or test phase, x is previously generated samples \hat{s}_l . The causal convolution
128 layer is used to make sure that the model does not violate the order and therefore the generation is
129 based on the previous time samples. Later, stacks of K dilation layers in the dilations area perform
130 dilated convolutions, residual connections, and skip connections. Note that the superscripts in Figure
131 1 show the layer index ($k = 1, \dots, K$). Convolutions with holes, as the dilated convolution layers,
132 process the input in a fine to coarse scale with fewer weight parameters in the sufficient receptive
133 field size. However, the residual and skip connection layers help avoiding the gradient vanishing
134 problem. In addition, the output of the skip connection layers $a^{[k]}$ contains various latent feature of
135 the input in different scales.

136 The morphological structure of the dilated convolutions resemble the wavelet transform. In fact, with
137 a specific set of weights, the first dilation layer can resemble the first level of the wavelet transform.
138 Hence, the first layer mostly models the high frequency features, likewise, the higher dilations for
139 the low frequencies. Thus, a stack of r repeats of $1, 2, 4, \dots, 2^r$ dilations for modeling the fullband
140 signal could be equivalent to r repeats of one dilation layer for each wavelet subband. Therefore,
141 in our experiments with subband signals, the number of dilation layers K is much smaller than the
142 original fullband WaveNet.

143 As the last block in each autoregressive subband generator, the post-processing performs two conse-
144 quent convolutional layers on summation of $a^{[k]}$ s, which are activations of the skip connection layers.

145 Because the signal is represented as one-hot vector, the post-processing ends with a softmax layer
 146 to increase the probability of the maximum value compared to others and to have a summation of
 147 probabilities equal to one.

148 3 Experiments

149 We used the TTS benchmark dataset LJ Speech¹ consisting of 13,100 short audio clips uttered by a
 150 female speaker, varying in length from 1 to 10 seconds, recorded in 16kHz sampling rate. We kept
 151 around 11 minutes of the speech signals (100 utterances) for test, which was not included in the train.
 152 The training set lengths more than 23 hours after the silence removal using voice activity detector
 153 (VAD).

154 3.1 Parameter settings

155 The subband decomposition is performed by Daubechies wavelet db10 for eight levels (L=8). Sub-
 156 band amplitude normalization is unavoidable because of the quantization in generator.

157 We found the Carnegie Mellon university pronouncing dictionary (CMUdict)² as a good choice for
 158 the lexicon including three levels of stress. The input phoneme sequence has 70 dimensions. The
 159 encoder contains three convolutional layers with filter width equals 5 and 256 channels. The HTK³
 160 aligns the phoneme sequence with the speech samples using forced-alignment. We have replaced the
 161 monophone with the triphone sequence but not that much change in results. In addition, we have tried
 162 summation of the activations of each layer in encoder as the latent feature but the results were worse.

163 The dilations of each generator are 1, 2, 4, 8, and 16. The channel size for dilation, residual, and
 164 skip-connection were set to 256. Adam optimizer [25] is used for training with the learning rate
 165 initiating from 10^{-3} and decaying every 50k iteration by a factor of 0.5.

166 3.2 Evaluation metrics

167 The evaluation metrics are signal-to-noise ratio in time domain and logarithmic spectral distortion
 168 (SD) which are defined as follows:

$$SNR_{[dB]} = 10 \log_{10} \frac{\sum_{t=1}^T s(t)^2}{\left| \sum_{t=1}^T s(t)^2 - \sum_{t=1}^T \hat{s}(t)^2 \right|} \quad (3)$$

$$SD_{[dB]} = \frac{1}{T} \sum_{t=1}^T \sqrt{\frac{1}{F} \sum_{f=1}^F \left[20 \log_{10} \frac{|S(f,t)|}{|\hat{S}(f,t)|} \right]^2}, \quad (4)$$

169 where $S(f,t)$ and $\hat{S}(f,t)$ are spectrograms of the target signal and the generated signal, respectively.
 170 The spectrograms are calculated by 16 ms frame length, 1 ms shift and Hanning window. In addition,
 171 because the human auditory perception is based on the Mel spectrogram representation, we considered
 172 Mel spectral distortion (MSD) as the third quantitative metric for the objective evaluation. The MSD
 173 is calculated similar to the SD, replacing the linear spectrogram with the 40-filters Mel spectrogram,
 174 which is obtained by 25 ms window length, and 5 ms shift. In addition, we calculated the SNR in the
 175 linear spectrogram domain. We did not mention the spectrogram SNR results because with two digits
 176 precision they are the same as the time domain ones.

177 The generation is time consuming in the proposed model because the speech is synthesized sample-
 178 by-sample and sequentially. Therefore, in addition to the above-mentioned metrics, we have measured
 179 the training and the synthesis time. The next subsection will explain results and discussions.

¹<https://keithito.com/LJ-Speech-Dataset/>

²<http://www.speech.cs.cmu.edu/cgi-bin/cmudict>

³<http://htk.eng.cam.ac.uk/>

Table 3: Evaluation results (mean \pm 95% CI) for 100 test set utterances

		SNR [dB]	SD [dB]	MSD [dB]
Decomposition-Reconstruction		41.5 \pm 1.14	0.61 \pm 0.01	0.08 \pm 0.02
Ground truth	Subband ($K = 5$)	23.5 \pm 0.31	4.3 \pm 0.02	2.5 \pm 0.01
	Fullband ($K = 24$)	18.8 \pm 0.47	8.1 \pm 0.03	5.5 \pm 0.04
Synthesis	Subband ($K = 5$)	4.0 \pm 0.88	13.3 \pm 0.01	10.0 \pm 0.10
	Fullband ($K = 24$)	5.2 \pm 0.93	15.2 \pm 0.10	11.8 \pm 0.11

Table 4: Average required time (minutes) for Generating 1 second of speech

CPU configuration	Fullband	Subband	
		sequential	parallel
Intel(R) Core(TM) i7, 2.93 GHz, 8 cores	1.67	6.8	2.08
Intel(R) Xeon(R), 2.4 GHz, 32 cores	2.09	5.36	1.87

180 3.3 Results and discussions

181 First experiment investigated the effect of the wavelet analysis/synthesis on the quality of speech
 182 without engaging any neural network model. The average results on 100 test set utterances with 95%
 183 confidence interval are reported in the first row of Table 3 as the extreme case for evaluations. For
 184 SNR, higher value shows more accurate model; whereas for both SD and MSD lower value means
 185 better performance. As shown in Table 3, the subband decomposition/reconstruction results provides
 186 near perfect performance.

187 Moreover, we compared the subband with the fullband speech synthesizer. The fullband term means
 188 that the model prediction \hat{s} is the speech signal in its full frequency range. Hence, there is no
 189 subband decomposition. Therefore, one complex signal generator models the probability distribution.
 190 Basically, the two models are exactly the same, except in the fullband TTS, $K = 24$ dilation layers
 191 are defined with 4 stacks of 1, 2, 4, \dots , 32 dilations in our experiments; while in the subband TTS,
 192 the number of dilation layers is much lower than the fullband ($K = 5$). Fast WaveNet algorithm
 193 [17] is utilized for the synthesis of both models. We have examined the fullband model without the
 194 encoder, which is in fact the original WaveNet conditioning on phoneme sequence; but the results
 195 were worse since the features were not enough for the training.

196 We compare the two models by conditioning on the phoneme sequence as the conditional features in
 197 two cases: *ground truth* and *synthesis*. The ground truth means feeding the previous clean samples
 198 to the model and evaluating the accuracy of the prediction of the next sample. As depicted in Table 3,
 199 the subband model performs significantly better than the fullband one in ground truth. For synthesis,
 200 the results are somehow comparable. In fact, the results of synthesis are not satisfying for both
 201 subband and fullband models, which is probably due to the lack of acoustical conditioning features.

202 Table 4 reports the average required time for synthesizing one second of speech in terms of minutes
 203 on two different machines with 8 and 32 cores. The required time of the subband TTS is reported in
 204 two cases: *sequential* and *parallel*. For the earlier experiment, the speech signal is decomposed into
 205 subband signals; and they are kept in the original sampling rate, which is 16 kHz. Thus, the samples
 206 are *redundant*. Obviously, without parallelization, the synthesis time of the redundant samples
 207 should be 8 times more than the fullband because there are 8 subband signals in the experiments.
 208 Nevertheless, since the complexity of the signal generators in the subband model is less than the
 209 fullband one, it is 4 and 2.5 times slower for the first and the second machine, respectively. For the
 210 last experiment, the subband signals are downsampled by a factor of 2. Hence, the subband signals
 211 are not redundant any more. Even though the parallelization and the downsampling speed up the
 212 synthesis, but it is still not that much far from the fullband model.

213 Both models need less than a day (around 18 hours) for training up to an admissible output quality on
 214 a Titan X GPU. Such a fast training is because of their fully CNN architecture, which is much better
 215 than the RNN-based TTS, e.g. Tacotron [6]. It is reported that an implementation of the Tacotron

216 takes 12 days (877K iterations) on a GTX 1080 Ti⁴. Note that the number of iterations is still much
217 less than the original Tacotron reported by Google (2M iterations) [6].

218 4 Conclusion

219 We proposed a subband time-domain TTS system inspiring from the WaveNet. The main differences
220 of our TTS with the WaveNet are twofold: first, rather than a complex deep neural network for
221 modeling the probability distribution of the speech signal, we designed separate (but integrated)
222 networks for each subband signal, which has much simple architecture and could estimate the
223 probability distributions of the subband signals accurately. Second, the original WaveNet TTS
224 benefits from pre-trained linguistic and acoustic feature extraction models; while an encoder in our
225 system extracts the latent features from the phoneme sequence input in a nearly end-to-end way,
226 which is more preferred.

227 The force alignment should be replaced by an attention mechanism for automatic aligning to have a
228 fully end-to-end model. Still enriching the conditional features by acoustic features beside the current
229 linguistic features is unavoidable. As another future work, we are trying to utilize the current dilated
230 architecture to extract acoustic features in a top-down way to improve the quality of both fullband
231 and subband models.

232 Acknowledgments

233 This work was supported by Institute for Information & communications Technology Promotion
234 (IITP) grant funded by the Korea government (MSIT) [2016-0-00562(R0124-16-0002), Emotional
235 Intelligence Technology to Infer Human Emotion and Carry on Dialogue Accordingly].

236 References

- 237 [1] A. J. Hunt and A. W. Black, "Unit selection in a concatenative speech synthesis system using a large speech
238 database," in Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 1996,
239 vol. 1, pp. 373–376.
- 240 [2] N. S. Kim and S. S. Park, "Discriminative training for concatenative speech synthesis," IEEE Signal Process.
241 Lett., vol. 11, no. 1, pp. 40–43, 2004.
- 242 [3] H. Zen, K. Tokuda, and A. W. Black, "Statistical parametric speech synthesis," Speech Commun., vol. 51,
243 no. 11, pp. 1039–1064, 2009.
- 244 [4] Y.-J. Hu and Z.-H. Ling, "DBN-based spectral feature representation for statistical parametric speech
245 synthesis," IEEE Signal Process. Lett., vol. 23, no. 3, pp. 321–325, 2016.
- 246 [5] Z.-C. Liu, Z.-H. Ling, and L.-R. Dai, "Statistical Parametric Speech Synthesis Using Generalized Distillation
247 Framework," IEEE Signal Process. Lett., vol. 25, no. 5, pp. 695–699, 2018.
- 248 [6] Y. Wang et al., "Tacotron: A fully end-to-end text-to-speech synthesis model," in Proc. Interspeech, 2017, pp.
249 4006–4010.
- 250 [7] H. Tachibana, K. Uenoyama, and S. Aihara, "Efficiently trainable text-to-speech system based on deep
251 convolutional networks with guided attention," in Proc. IEEE International Conference on Acoustics, Speech
252 and Signal Processing (ICASSP), 2018, pp. 4784–4788.
- 253 [8] A. Van Den Oord et al., "Wavenet: A generative model for raw audio," CoRR, vol. abs/1609.0, 2016.
- 254 [9] G. Lai, B. Li, G. Zheng, and Y. Yang, "Stochastic WaveNet: A Generative Latent Variable Model for
255 Sequential Data," arXiv Prepr. arXiv1806.06116, 2018.
- 256 [10] K. Qian, Y. Zhang, S. Chang, X. Yang, D. Florêncio, and M. Hasegawa-Johnson, "Speech enhancement
257 using Bayesian wavenet," in Proc. Interspeech, 2017, pp. 2013–2017.
- 258 [11] D. Rethage, J. Pons, and X. Serra, "A Wavenet for speech denoising," in Proc. IEEE International Conference
259 on Acoustics, Speech and Signal Processing (ICASSP), 2018, pp. 5069–5073.
- 260 [12] A. Tamamori, T. Hayashi, K. Kobayashi, K. Takeda, and T. Toda, "Speaker-dependent WaveNet vocoder,"
261 in Proc. Interspeech, 2017, vol. 2017–August, pp. 1118–1122.

⁴<https://github.com/keithito/tacotron>

- 262 [13] T. Hayashi, A. Tamamori, K. Kobayashi, K. Takeda, and T. Toda, "An investigation of multi-speaker training
263 for wavenet vocoder," in Proc. Automatic Speech Recognition and Understanding Workshop (ASRU), 2017, pp.
264 712–718.
- 265 [14] W. Ping et al., "Deep voice 3: Scaling text-to-speech with convolutional sequence learning," in Proc. 6th
266 International Conference on Learning Representations (ICLR), 2018, vol. 79, no. 14, pp. 1094–1099.
- 267 [15] T. Yoshimura, K. Hashimoto, K. Oura, Y. Nankaku, and K. Tokuda, "Mel-cepstrum-based quantization
268 noise shaping applied to neural-network-based speech waveform synthesis," IEEE/ACM Trans. Audio, Speech,
269 Lang. Process., vol. 26, no. 7, pp. 1173–1180, 2018.
- 270 [16] J. Shen et al., "Natural TTS synthesis by conditioning WaveNet on Mel spectrogram predictions," in Proc.
271 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2018, pp. 4779–4783.
- 272 [17] T. Le Paine et al., "Fast wavenet generation algorithm," CoRR, vol. abs/1611.0, 2016.
- 273 [18] A. Van Den Oord et al., "Parallel WaveNet: Fast high-fidelity speech synthesis," CoRR, vol. abs/1711.1,
274 2017.
- 275 [19] T. Inoue, S. Hara, and M. Abe, "A hybrid text-to-speech based on sub-band approach," in Proc. Signal and
276 Information Processing Association Annual Summit and Conference (APSIPA), 2014 Asia-Pacific, 2014, pp.
277 1–4.
- 278 [20] T. Okamoto, K. Tachibana, T. Toda, Y. Shiga, and H. Kawai, "An investigation of subband WaveNet vocoder
279 covering entire audible frequency range with limited acoustic features," in Proc. IEEE International Conference
280 on Acoustics, Speech and Signal Processing (ICASSP), 2018, pp. 5654–5658.
- 281 [21] Y. Lee, A. Rabiee, and S.-Y. Lee, "Emotional End-to-End Neural Speech Synthesizer," in Workshop
282 Machine Learning for Audio Signal Processing at NIPS (ML4Audio@NIPS17), 2017.
- 283 [22] I. Daubechies, Ten lectures on wavelets, vol. 61. Siam, 1992.
- 284 [23] M. Farge, "Wavelet transforms and their applications to turbulence," Annu. Rev. Fluid Mech., vol. 24, no.
285 1, pp. 395–458, 1992.
- 286 [24] Y. Wang et al., "Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech
287 synthesis," in Proc. International Conference on Machine Learning (ICML), 2018.
- 288 [25] D. P. Kingma and J. L. Ba, "Adam: A method for stochastic optimization," in Proc. of the 3rd International
289 Conference on Learning Representations (ICLR), 2015.