

SIMD-SIZE AWARE WEIGHT REGULARIZATION FOR FAST NEURAL VOCODING ON CPU

Hiroki Kanagawa and Yusuke Ijima

NTT Corporation

ABSTRACT

This paper proposes weight regularization for a faster neural vocoder. Pruning time-consuming DNN modules is a promising way to realize a real-time vocoder on a CPU (*e.g.* WaveRNN, LPCNet). Regularization that encourages sparsity is also effective in avoiding the quality degradation created by pruning. However, the orders of weight matrices must be contiguous in SIMD size for fast vocoding. To ensure this order, we propose explicit SIMD size aware regularization. Our proposed method reshapes a weight matrix into a tensor so that the weights are aligned by group size in advance, and then computes the group Lasso-like regularization loss. Experiments on 70% sparse subband WaveRNN show that pruning in conventional Lasso and column-wise group Lasso degrades the synthetic speech’s naturalness. The vocoder with proposed regularization 1) achieves comparable naturalness to that without pruning and 2) performs meaningfully faster than other conventional vocoders using regularization.

Index Terms— speech synthesis, neural vocoder, regularization, SIMD, group pruning

1. INTRODUCTION

The neural vocoder represented by WaveNet [1] has dramatically improved the quality of text-to-speech (TTS) synthesis. While WaveNet can generate high-quality speech waveforms directly from conditioning features via a large causal convolution-based autoregressive model, its huge computational cost and autoregressive (AR) architecture prevent fast vocoding. To allow easier parallel computation, many neural vocoders based on non-autoregressive structures have been proposed [2, 3, 4, 5, 6, 7]. These schemes offer high processing speeds if the device is specialized for parallel computing, such as GPUs. On the other hand, to achieve fast neural vocoding on CPUs, the computational complexity must be reduced drastically. WaveRNN achieves real-time vocoding by replacing huge causal convolutions of WaveNet with a simple GRU, and pruning its weights [8]. LPCNet also introduces signal processing insights into the speech generation process and reduces the number of DNN parameters from that of WaveRNN [9]. Other approaches to reduce the number of DNN inferencing iterations, with multi-sample generation in a sin-

gle forward propagation step [10, 11] and prediction of shortened subband signals instead of waveforms [12, 13].

Quantization and low-rank approximation are promising alternatives to pruning for paring the DNN module’s computational complexity. [13] roughly quadrupled speeds by quantizing neural vocoder weights to 8-bit integers, but it requires quantization error-aware training and an appropriate intrinsic implementation for integers, resulting in high implementation cost and hardware dependency. Although the low-rank approximation shrinks the model size, its speed-up contribution is limited due to the increased number of matrices that must be computed, which requires more matrix-vector product instructions (*e.g.* “*gemv*”) to be called [14]. Pruning [15, 16] used in WaveRNN and LPCNet substitutes the elements of the weight matrix with zeros in the training process. During inference, the calculation of zero weights can be skipped, and hardware dependency is also small because the weights can be treated as floating-point without modification. While excessive pruning leads to quality degradation, regularization is an effective solution. Lasso regularization promotes a sparse model, thus allowing for a smaller gap in the model with and without pruning [17]. However, the order of non-zero weight elements becomes non-contiguous when using Lasso for regularization. To exploit fully the fast single instruction multiple data (SIMD) operations intrinsic to CPUs, the non-zero elements must be contiguous, so Lasso is not optimal for SIMD. The use of group Lasso (gLasso) [18] allows for a sequence of non-zero elements, but the model’s expressiveness is sacrificed due to row- or column-wise weight sparsity.

To ensure both speed and quality, we propose a gLasso-like regularization approach that explicitly sparsifies weights with the group pruning size. The proposed method aligns the weights in both the non-zero and zero regions with enough size to fully occupy the SIMD registers at once. Thus, group pruning across the boundary between regions can be avoided, and computational efficiency can be improved while maintaining model expressiveness. In our experiments, we used multi-sample subband WaveRNN [19] with 70% sparsity. While pruning degraded the naturalness of the synthetic speech of the vocoder without regularization and the one with the use of Lasso and gLasso regularization, regularization with the proposed method and pruning maintained the same naturalness as that achieved before pruning. We also found

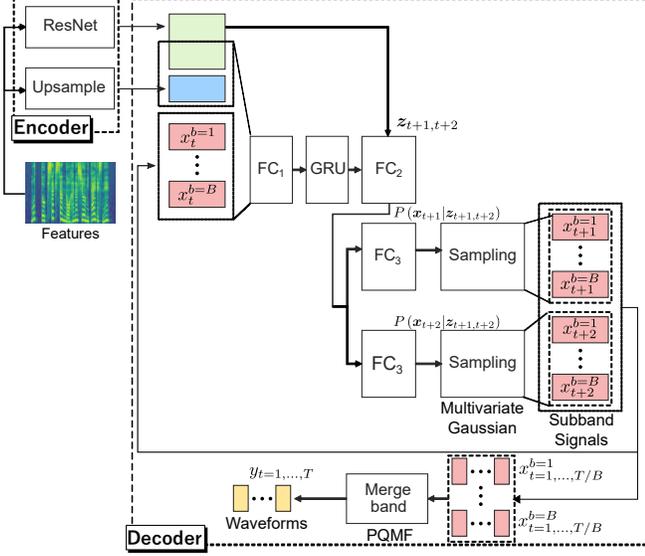


Fig. 1. Overview of multi-sample subband WaveRNN via multivariate Gaussian [19]. This vocoder predicts M subband signals simultaneously in single forward propagation step ($M=2$ is used in this figure).

that our vocoder was faster than the conventional alternatives.

2. MULTI-SAMPLE SUBBAND Wavernn

2.1. Model architecture and its training

Subband WaveRNN [13] reduces sequence length from T to T/B by predicting B -band subband signals instead of a speech waveform. [19] extended it to multi-sample generation for even faster vocoding; Fig. 1 shows an overview. The model consists of an encoder and a decoder, which are responsible for frame rate and sample rate, respectively. Encoders upsample frame-level acoustic features to corresponding samples. The decoder generates predictions of the next time $t + \tau \forall m \in [1, M - 1]$ from the output of the encoder and the previous M subband signals, where m is the index of the number of subband signals to be generated simultaneously. To generate multiple samples in single forward propagation, linear layer FC3s are provided for subband signals of $t + \tau \forall m \in [1, M - 1]$. This module jointly predicts the associations among subband signals because PQMF has band overlaps. Assuming a multivariate Gaussian as FC3's target, this vocoder minimizes the negative log-likelihood given by:

$$\mathcal{L}_{\text{NLL}}(\theta) = - \sum_{t=1}^{T/B} \sum_{m=1}^M \ln \mathcal{N}(\mathbf{x}_{t+m}; \boldsymbol{\mu}(\mathbf{z}_{t+\tau \forall m}, \theta), \boldsymbol{\Sigma}(\mathbf{z}_{t+\tau \forall m}, \theta)), \quad (1)$$

where θ , \mathbf{z}_t , $\mathbf{x}_t \in \mathbb{R}^B$, $\boldsymbol{\mu} \in \mathbb{R}^B$ and $\boldsymbol{\Sigma} \in \mathbb{R}^{B \times B}$ are the DNN model parameters, FC2's output, subband signals, the mean vector and covariance matrix of the multivariate Gaussian, respectively. To guarantee spectral reproducibility,

STFT loss $\mathcal{L}_{\text{STFT}}(\theta)$ [20] is calculated by generating subband signals from the multivariate Gaussian via a reparameterization trick. This is added to Eq. (1) without scaling to optimize the vocoder.

2.2. Weight pruning

For fast vocoding, pruning is performed during training. We performed pruning by gradually increasing sparsity [16] in the same manner as WaveRNN using:

$$d_s = d [1 - \{1 - (s - s_0) / S\}^3], \quad (2)$$

where s , s_0 , and S are the current-, start-, and total- pruning step, respectively. d is the target density, thus the sparsity is defined as $1 - d$. To fully utilize vector algebra with SIMD, we apply group pruning [15]. The group size of FC1, GRU, and FC2 for pruning was set to 16. Taking Intel's AVX2 intrinsic instruction set as SIMD, the dot product can be calculated for 16 elements in two SIMD operations. This is done by putting eight 32-bit float elements on a register and calculating the dot product, which is then applied to and added to the other eight neighboring elements.

2.3. Weight regularization

In order to avoid a degradation in model expressiveness due to pruning, regularization is an efficient approach to sparsify the model in advance. A well-known Lasso regularization term is computed as follows:

$$\mathcal{L}_{\text{Reg}}^{\text{LASSO}}(\theta) = \sum_r \sum_{i=1}^{I_r} \sum_{j=1}^{J_r} |W_r(i, j)|, \quad (3)$$

where r , $\mathbf{W}_r \in \mathbb{R}^{I_r \times J_r}$ are the DNN module's index and the weight matrix to be regularized, respectively. i and j denote the row and column indices.

The regularization term of column-wise group Lasso is given by:

$$\mathcal{L}_{\text{Reg}}^{\text{GLASSO}}(\theta) = \sum_r \sum_{j=1}^{J_r} \|\mathbf{w}_r^j\|, \quad (4)$$

where \mathbf{w}_r^i and $\|\cdot\|$ are the \mathbf{W}_r 's i -th column vector and the L2 norm operator, respectively. The final objective function with regularization is reformulated as:

$$\mathcal{L}(\theta) = \mathcal{L}_{\text{NLL}}(\theta) + \mathcal{L}_{\text{STFT}}(\theta) + \lambda \mathcal{L}_{\text{Reg}}(\theta), \quad (5)$$

where λ is the scale for the regularization term $\mathcal{L}_{\text{Reg}}(\theta)$.

3. PROPOSED SIMD-SIZE AWARE WEIGHT REGULARIZATION

The Lasso and gLasso regularizations described in Section 2.3 could suffer from group pruning of sparse weights across the zero/non-zero boundary, resulting in quality degradation.

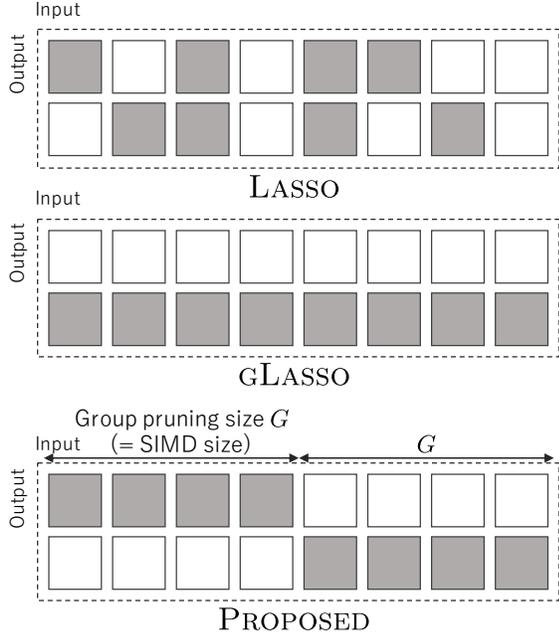


Fig. 2. Comparison between Lasso, gLasso, and our regularization proposal applied to a single weight matrix. The white and gray regions denote non-zero and zero components, respectively.

To overcome this problem, we propose a regularization that assumes group pruning at the SIMD size. The weight matrix \mathbf{W}_r is reshaped in advance into a third-order tensor $\mathcal{W}_r \in \mathbb{R}^{I_r \times J_r / G \times G}$ to make the sparse weight’s block size equal to group pruning size G . Then, the proposed regularization term for this tensor is formulated as follows:

$$\mathcal{L}_{\text{Reg}}^{\text{PROPOSED}}(\theta) = \sum_r \sum_{i=1}^{I_r} \sum_{j=1}^{J_r} \sqrt{\sum_{g=1}^G \mathcal{W}_r(i, j, g)^2}. \quad (6)$$

Figure 2 compares regularized weights when Lasso, gLasso, and our proposed method (PROPOSED) are applied. The white and gray regions of the matrices in this figure are the non-zero and zero components, respectively. This prior matrix-to-tensor conversion and computing Eq. (6) contribute to SIMD size group-wise sparsification. Therefore, group-wise sparsified weights can be preferentially pruned, reducing the probability of pruning portions contain non-zero elements. Since weights are continuous in SIMD size, weights can always be assigned to aligned memory, which is advantageous in terms of computing efficiency.

4. EXPERIMENTS

4.1. Setup

We used speech data uttered by a Japanese professional female speaker. The sampling frequency was 22.05 kHz. 200

utterances were extracted as evaluation data (18.3 minutes), and the remainder were used for training and validation (30.6 hours).

Eighty-dimensional logarithmic mel-spectrograms were used as the conditioning feature of the neural vocoder. The analysis frame shift was 5 ms¹. The ResNet of the encoder has ten residual blocks, each consisting of 1D-convolution with 128 units, batch normalization, and activation. ReLU was used for all activations, and the simultaneous generation sample was set to two likewise [19]. This multi-sample vocoder occasionally failed to predict accurate variance parameters which yielded clicking sounds. To avoid this problem, we 1) eliminated variance outliers and 2) clipped sampled results in a similar way to [10]. The number of training steps was 5000k, and parameters for pruning steps in Eq. (2) were set to $S_0 = 2000k$ and $S = 2500k$. We used $d = 0.3$ and λ in Eq. (5) to 1.0×10^{-4} for all regularization methods. This guarantees that the model sparsity is at least 70%. If the number of zero elements increases due to regularization, the sparsity could be higher than this. The vocoder’s optimization was performed using RAdam [21], with $\alpha = 1.0 \times 10^{-4}$, $\beta = (0.9, 0.999)$, and $\varepsilon = 1.0 \times 10^{-8}$.

4.2. Weight heatmap comparison

Figure 3 shows FC1 (described in Fig. 1) weight heatmaps before and after pruning for each method. The horizontal and vertical axes are the input and output dimensions, respectively. Comparing w/o REG (w/o PRUNE) and w/o REG (w/ PRUNE), we can see that the elements have been replaced with zero leaving large components. As described in Fig. 1, the recurrent multiple subband signals are fed to FC1, so these weights are particularly large on the left-side of the heatmap as if to focus on their signals (e.g. “Detail 1” in Fig. 3). Other right-side weight components are responsible for receiving encoder’s outputs (e.g. “Detail 2”). These results revealed that w/o REG (w/ PRUNE) disregarded or neglected the encoder’s conditioned outputs. Although LASSO (w/o PRUNE) was not overly dependent on the most recent sample, its weights had discontinuous order. The GLASSO (w/o PRUNE) yield poor sparsification, because under the constraint of sparsifying entire columns, it was difficult to find compact representations. These regularizations, like w/o REG (w/o PRUNE), are prone to degrade the synthesized speech’s quality because group pruning is applied across the boundaries between zero and non-zero. On the other hand, the proposed method PROPOSED (w/o PRUNE) was sparse and its weights were continuously aligned in group pruning size $G = 16$ in the column direction before pruning. Since there is no drastic difference in the heatmaps between with and without pruning, we found

¹Although we also investigated the commonly used frame shift of 12.5 ms in our preliminary experiments, we chose to set it to 5 ms because it better reproduced the pitch of synthetic speech. If a faster inference speed is preferred, the frame shift can be set to 12.5 ms like other studies for lower encoder computational complexity.

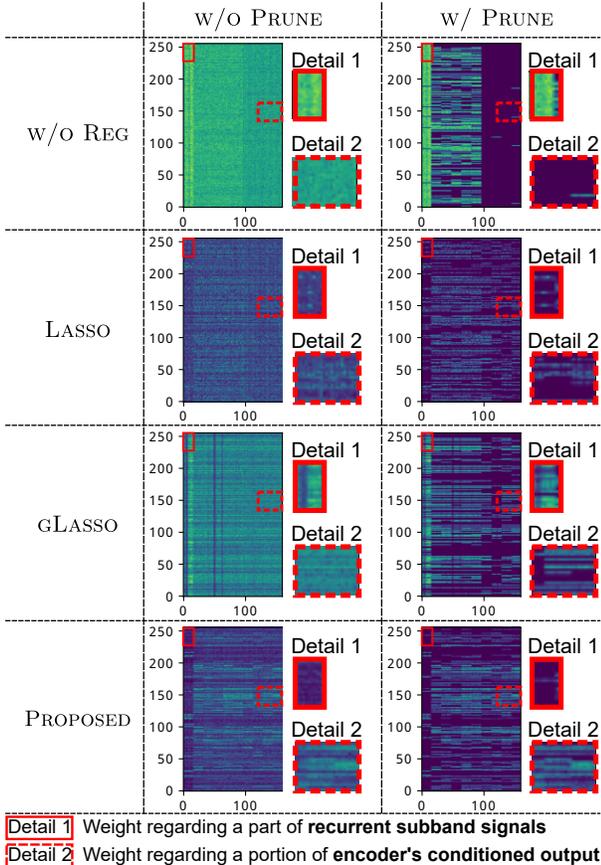


Fig. 3. FC1’s weight heatmap comparison. From top to bottom: w/o REG, LASSO, GLASSO, and PROPOSED. The right and left heatmaps show without and with pruning, respectively. The left-side of the heatmaps receive the recurrent multiple subband signals (e.g. “Detail 1”). Other right-side of ones accept the encoder’s conditioned output (e.g. “Detail 2”). The comparison of “Detail 1” and “2” between w/o and w/ PRUNE, demonstrated that PROPOSED kept the minimum difference among them.

that pruning can be done without sacrificing the model’s expressiveness.

4.3. Subjective evaluations

We subjectively evaluated the naturalness of synthetic speech by using mean opinion score (MOS) on a five-point scale ranging from 5: very natural to 1: very unnatural. Sixty listeners participated in the test via crowdsourcing. They evaluated ten sentences for each method, randomly selected from all 200 evaluation data, for a total of sixty sentences. These participants were different evaluators for analytic resynthesis and TTS.

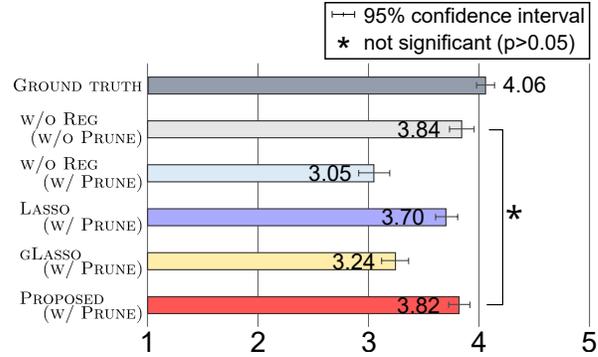


Fig. 4. Mean opinion scores of naturalness. Acoustic features for vocoding were extracted from natural speech.

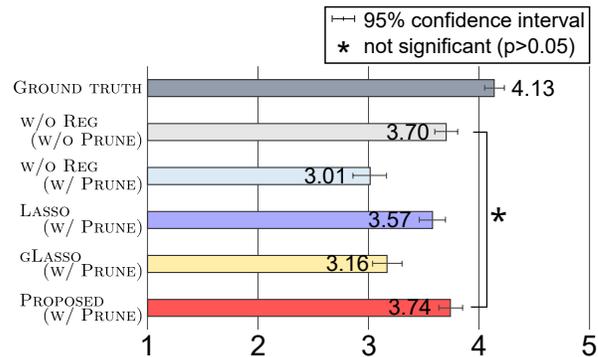


Fig. 5. Mean opinion scores in terms of naturalness. Acoustic features for vocoding were predicted by the TTS model.

4.3.1. Vocoding for extracted acoustic features from natural speech

Figure 4 shows the subjective evaluation results of vocoding with acoustic features extracted from natural speech. w/o REG (w/o PRUNE) yielded lower performance than Ground truth, but still obtained high naturalness. w/o REG (w/ PRUNE) degraded significantly more than w/o REG (w/o PRUNE). This is due to 1) neglecting the conditioning spectral information in pruning, and 2) excessive reliance on recurrent samples led to quality degradation due to mismatches with the training time; as described in Section 4.2. The performance of GLASSO (w/ PRUNE) also falls for the same reason, just not as much as w/o REG (w/ PRUNE). On the other hand, LASSO (w/ PRUNE) showed no more significant degradation from w/o REG, even with pruning. This is due to the fact that model does not rely excessively on recent samples, but focuses more on the conditioning spectral information. PROPOSED (w/ PRUNE) outperforms these pruned models and achieves naturalness comparable to that of w/o REG (w/o PRUNE).

4.3.2. Vocoding for acoustic features predicted by TTS

To investigate robustness against degraded acoustic features, FastSpeech2 [22] as the TTS model was also trained with the

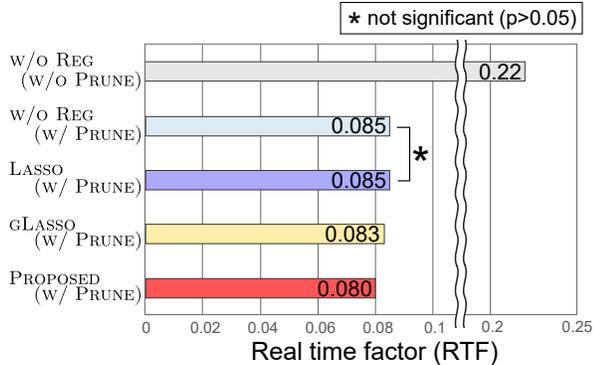


Fig. 6. Average RTFs obtained from all evaluation data.

same data as the neural vocoders. We fed 380 kinds of symbols including phoneme and prosodic information to FastSpeech2. It was optimized via a minimum mean absolute error criterion with 2000k steps by Adam [23] with $\beta = (0.9, 0.98)$, and $\varepsilon = 1.0 \times 10^{-9}$. We followed the same learning rate schedule in [24].

Figure 5 shows the subjective evaluation results. The overall difference in scores between ground truth and synthetic speech was greater when using features extracted from natural speech since acoustic features predicted by TTS were degraded from the original one. The PROPOSED (w/ PRUNE) matched the naturalness of w/o REG (w/o PRUNE). This result confirms that the regularization proposal is robust to acoustic features degraded by TTS.

4.4. Vocoding speed comparison

The real-time factors (RTFs) were calculated to measure the inference speeds for all methods. RTF is defined by:

$$\text{RTF} := T_{\text{inference}}/T_{\text{data}}, \quad (7)$$

where T_{data} and $T_{\text{inference}}$ are speech length and single-thread inference time measured on an Intel Core i7-8750H CPU 2.20 GHz, respectively.

Figure 6 shows averaged RTFs from all evaluation data. w/o REG (w/ PRUNE) yielded a significant speed improvement over w/o REG (w/o PRUNE) owing to pruning. LASSO (w/ PRUNE) showed no speed improvement from w/o REG (w/ PRUNE). As discussed in Section 4.2, this was attributed to Lasso regularization promoting non-contiguous sparse matrices, which fails to yield a SIMD-friendly contiguous sparse matrix. On the other hand, GLASSO (w/ PRUNE) slightly improved the RTF. The reason is the 70% sparsity by pruning, further increased column-wise contiguous zeroed regions via regularization. The proposed method achieved better RTFs than GLASSO (w/ PRUNE), since it achieved a group-wise sparser matrix than GLASSO, as mentioned in Section 4.2. Our RTF=0.080 is nearly comparable to the one of the recent non-AR HiFi-GAN (v3) [7], which works fast on CPUs

(RTF=0.075). Since they used a higher clock CPU (Intel Core i7 CPU 2.6 GHz) than ours, our proposed method might outperform their speed if on the same CPU. Furthermore, the combination of the other faster approach for WaveRNN (e.g. [25]) and the proposed regularization would be able to provide a significant speed-up compared to HiFi-GAN.

5. CONCLUSION

In this work, we proposed SIMD-size aware group-wise regularization to avoid the quality degradation associated with neural vocoder pruning. We incorporated the regularization proposal into a subband WaveRNN-based vocoder and showed that the regularized weights have group-wise continuous orders suitable for SIMD computation. No major differences in our vocoder’s weight layout were observed via heatmaps of before and after pruning. Subjective evaluations regarding naturalness demonstrated that the proposed pruned vocoder outperformed that with no regularization, Lasso, and group Lasso. In particular, our vocoder achieved comparable naturalness to that achieved without pruning. A speed evaluation also revealed that our vocoder performed significantly faster than the existing alternatives.

Our method can also increase contiguous zeroed region more efficiently than LASSO and GLASSO even if regularized to fit processors with smaller SIMD sizes, e.g. Intel SSE4 and Arm NEON. So we expect to run significantly faster than them without any loss of quality. Our first future work will apply the proposed method for Arm NEON and confirm its efficiency on embedded processors. Our second future work will compare and combine our regularization with AlignReg [26], which has similar concepts proposed for natural language processing. Since the proposed regularization is not limited to neural vocoders, we will also plan to apply it to other computationally expensive models (e.g. RNN-T [27] and BERT [28]).

6. REFERENCES

- [1] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu, “WaveNet: A generative model for raw audio,” *CoRR abs/1609.03499*, 2016.
- [2] Aäron van den Oord, Yazhe Li, Igor Babuschkin, Karen Simonyan, Oriol Vinyals, Koray Kavukcuoglu, George Van Den Driessche, Edward Lockhart, Luis C. 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,” *Proc. ICML*, vol. 9, 2018.
- [3] Wei Ping, Kainan Peng, and Jitong Chen, “Clarinet: Parallel wave generation in end-to-end text-to-speech,” *Proc. ICLR*, 2019.
- [4] Bryan Catanzaro Ryan Prenger, Rafael Valle, “WaveGlow: A flow-based generative network for speech synthesis,” *Proc. ICASSP*, pp. 3617–3621, 2019.
- [5] Ryuichi Yamamoto, Eunwoo Song, and Jae-Min Kim, “Parallel WaveGAN: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram,” *Proc. ICASSP*, pp. 6199–6203, 2020.
- [6] Kundan Kumar, Rithesh Kumar, Thibault de Boissiere, Lucas Gestein, Wei Zhen Teoh, Jose Sotelo, Alexandre de Brébisson, Yoshua Bengio, and Aaron C. Courville, “MelGAN: Generative adversarial networks for conditional waveform synthesis,” *Proc. NeurIPS*, 2019.
- [7] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae, “HiFiGAN: Generative adversarial networks for efficient and high fidelity speech synthesis,” *Proc. NeurIPS*, 2020.
- [8] Nal Kalchbrenner, Erich Elsen, Karen Simonyan, Seb Noury, Norman Casagrande, Edward Lockhart, Florian Stimberg, Aäron Van Den Oord, Sander Dieleman, and Koray Kavukcuoglu, “Efficient Neural Audio Synthesis,” *Proc. ICML*, vol. 80, pp. 2410–2419, 2018.
- [9] Jean-Marc Valin and Jan Skoglund, “LPCNet: Improving neural speech synthesis through linear prediction,” *Proc. ICASSP*, pp. 5891–5895, 2019.
- [10] Vadim Popov, Mikhail Kudinov, and Tasnima Sadekova, “Gaussian LPCNet for multisample speech synthesis,” *Proc. ICASSP*, pp. 6204–6208, 2020.
- [11] Patrick Lumban Tobing, Yi-Chiao Wum, Tomoki Hayashi, Kazuhiro Kobayashi, and Tomoki Toda, “Efficient shallow wavenet vocoder using multiple samples output based on laplacian distribution and linear prediction,” *Proc. ICASSP*, pp. 7204–7208, 2020.
- [12] Takuma Okamoto, Kentaro Tachibana, Tomoki Toda, Yoshinori Shiga, and Hisashi Kawai, “Subband WaveNet with overlapped single-sideband filterbanks,” *Proc. ASRU*, pp. 698–704, 2017.
- [13] Chengzhu Yu, Heng Lu, Na Hu, Meng Yu, Chao Weng, Kun Xu, Peng Liu, Deyi Tuo, Shiyin Kang, Guangzhi Lei, Dan Su, and Dong Yu, “DurIAN: Duration informed attention network for speech synthesis,” *Proc. Interspeech*, pp. 2027–2031, 2020.
- [14] Hiroki Kanagawa and Yusuke Ijima, “Lightweight LPCNet-based neural vocoder with tensor decomposition,” *Proc. Interspeech*, pp. 205–209, 2020.
- [15] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf, “Pruning filters for efficient convnets,” *Proc. ICLR*, 2017.
- [16] Michael Zhu and Suyog Gupta, “To prune, or not to prune: exploring the efficacy of pruning for model compression,” *Proc. ICLR*, 2018.
- [17] Markus Thom and Günther Palm, “Sparse activity and sparse connectivity in supervised learning,” *Journal of Machine Learning Research*, vol. 14, no. 1, pp. 1091–1143, 2013.
- [18] Simone Scardapane, Danilo Comminiello, Amir Husain, and Aurelio Uncini, “Group sparse regularization for deep neural networks,” *Neurocomputing*, vol. 241, pp. 81–89, 2017.
- [19] Hiroki Kanagawa and Yusuke Ijima, “Multi-sample subband WaveRNN via multivariate Gaussian,” *Proc. ICASSP*, pp. 8427–8431, 2022.
- [20] Shinji Takaki, Toru Nakashika, Xin Wang, and Junichi Yamagishi, “STFT Spectral loss for training a neural speech waveform model,” *Proc. ICASSP*, pp. 7065–7069, 2019.
- [21] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han, “On the variance of the adaptive learning rate and beyond,” *Proc. ICLR*, 2020.
- [22] Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu, “Fastspeech 2: Fast and high-quality end-to-end text to speech,” *Proc. ICLR*, 2021.

- [23] Diederik P. Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *Proc. ICLR*, 2015.
- [24] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” *Proc. NIPS*, 2017.
- [25] Hiroki Kanagawa, Yusuke Ijima, and Hiroyuki Toda, “Joint modeling of multi-sample and subband signals for fast neural vocoding on CPU,” *Proc. Interspeech*, pp. 1626–1630, 2022.
- [26] James O’ Neill, Sourav Dutta, and Haytham Assem, “Aligned weight regularizers for pruning pretrained neural networks,” *Proc. ACL*, pp. 3391–3401, 2022.
- [27] Alex Graves, “Sequence transduction with recurrent neural networks,” *Proc. ICML*, 2012.
- [28] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina N. Toutanova, “BERT: Pre-training of deep bidirectional transformers for language understanding,” *Proc. NAACL-HLT*, pp. 4171–4186, 2019.