

A STREAMWISE GAN VOCODER FOR WIDEBAND SPEECH CODING AT VERY LOW BIT RATE

Ahmed Mustafa, Jan Bütthe*, Srikanth Korse, Kishan Gupta, Guillaume Fuchs, Nicola Pia

Fraunhofer IIS, Erlangen, Germany

{ahmed.mustafa.ahmed, srikanth.korse, kishan.gupta, guillaume.fuchs, nicola.pia}@iis.fraunhofer.de

ABSTRACT

Recently, GAN vocoders have seen rapid progress in speech synthesis, starting to outperform autoregressive models in perceptual quality with much higher generation speed. However, autoregressive vocoders are still the common choice for neural generation of speech signals coded at very low bit rates. In this paper, we present a GAN vocoder which is able to generate wideband speech waveforms from parameters coded at 1.6 kbit/s. The proposed model is a modified version of the StyleMelGAN vocoder that can run in frame-by-frame manner, making it suitable for streaming applications. The experimental results show that the proposed model significantly outperforms prior autoregressive vocoders like LPCNet for very low bit rate speech coding, with computational complexity of about 5 GMACs, providing a new state of the art in this domain. Moreover, this streamwise adversarial vocoder delivers quality competitive to advanced speech codecs such as EVS at 5.9 kbit/s on clean speech, which motivates further usage of feed-forward fully-convolutional models for low bit rate speech coding.

Index Terms— GAN vocoder, StyleMelGAN, neural speech synthesis, LPCNet, speech coding

1. INTRODUCTION

Despite decades of extensive work classical speech coders offer very low quality at bit rates under 3 kbit/s. New techniques based on the use of neural networks showed breakthrough advancements in this area in recent years, enabling compression factors much higher than conventional approaches, while maintaining acceptable quality. Neural speech coders are based on the classical encoder-decoder scheme: the encoder analyzes the input signal and extracts a set of acoustic features, which are then quantized, coded and transmitted; the decoder reconstructs the input signal using the information contained in the received bit stream. In neural speech coders a generative neural network plays the role of the decoder (i.e., neural vocoder), as illustrated in Figure 1. It was demonstrated [1, 2] that conditioning a neural vocoder with coded acoustic parameters could produce natural wideband speech at bit rates lower than 2 kbit/s.

In recent years neural vocoders [3, 4, 5, 6] have revolutionized fields such as text-to-speech, voice conversion and speech enhancement, generating speech of unprecedented high quality. Most of these solutions however, are not suitable for speech coding purposes. This is mainly due to their high computational complexity or very slow generation speed, with clear quality degradation when using coarsely quantized conditioning features.

*This work was done while the author was at Fraunhofer. Now he works at Fantasma, email: Jan.Buethe@gmx.net



Figure 1: High-level block-diagram of a neural speech coder.

Neural vocoders based on generative adversarial networks (GANs) [7] were recently shown to be competitive and viable alternatives to autoregressive and flow-based models for speech synthesis applications [8, 9, 10]. However, they are by design not suited for streaming or real-time speech communication, since they take the advantage of heavy parallelization for processing large blocks of conditioning information at once. This permits efficient generation of speech waveforms in one shot, but exploits the advantage of having the acoustic features encoding information about future samples, which are not available in a streaming scenario because of the high algorithmic delay they would cause. Moreover, GAN vocoders work particularly well with homogeneous speech representations such as mel-spectrograms, whereas speech coding applications primarily use non-homogeneous (e.g., parametric) speech representations that may not easily condition GAN vocoders for high-quality signal generation.

To solve the above-mentioned issues, our contributions in this work are twofold:

- We propose Streamwise StyleMelGAN (SSMGAN), a modified StyleMelGAN vocoder for frame-by-frame generation of wideband speech at low delay, with reasonable computational complexity.
- We demonstrate that SSMGAN is able to generate high-quality speech even when conditioned with a parametric and highly compressed representation provided by the encoder of LPCNet [2], which delivers a 1.6 kbit/s bitstream to our StyleMelGAN-based vocoder.

2. RELATED WORKS

The research on neural vocoders is a very active field with new models being presented every few months. For this reason, here we only refer to some of the ones which sparked the most attention. The first family to appear was the one of autoregressive models [3, 5, 6], followed by flow-based models [4], and then GANs [8, 11, 12, 9, 10].

The first work to show the feasibility of low bit rate neural speech coding was [1], using a WaveNet decoder. The decoder network's complexity makes it impossible to deploy it in concrete applications. The complexity issue was partially tackled with a different approach in [13]. Finally the LPCNet model [2] introduced optimizations which made neural speech coding possible on edge

device. Moreover, the coding scheme used in LPCNet has a very low bit rate of 1.6 kbit/s. The coding parameters include acoustic features classically used in parametric speech coding, i.e. the Bark scale cepstrum, the pitch information and the energy. Table 1 describes in detail these parameters and the bit budget allocated to code them.

Coding Parameter	Bits/packet
Pitch lag	6
Pitch modulation	3
Pitch correlation	2
Energy	7
Cepstrum absolute coding	30
Cepstrum delta coding	13
Cepstrum interpolation	3
Total	64

Table 1: LPCNet coding parameters and their bit allocation for a 40 m sec packet

LPCNet’s 1.6 kbit/s decoder is an autoregressive architecture based on WaveRNN generating sample-by-sample wideband speech (16 kHz). It relies on linear prediction to reduce computational complexity, hence generating the signal in the residual linear prediction domain. The decoding step is divided into two parts: a frame-rate network that computes the conditioning for every 10 ms frame using the coded parameters, and a sample-rate network that computes the conditional sampling probabilities. LPCNet predicts the new excitation sample using the previously generated excitation and speech samples, as well as the current linear prediction sample from the 16th-order linear prediction.

More recent work [14] presented a new neural speech decoder (Lyra) compressing speech at 3 kbit/s. The encoder directly codes stacked mel-spectra and the decoder uses noise suppression and variance regularization to improve the quality of out-of-distribution samples. When compared to the proposed solution, Lyra is conditioned on a substantially different bit stream and works under different conditions (e.g. noisy speech).

To the best of our knowledge, there exists no prior GAN vocoder which allows frame-by-frame generation of speech at low delay or which provides high quality speech synthesis conditioned on a coded bit stream.

3. STREAMWISE STYLEMELGAN VOCODER (SSMGAN)

3.1. Baseline StyleMelGAN

StyleMelGAN [10] is a lightweight neural vocoder allowing synthesis of high-fidelity speech with low computational complexity. It employs Temporal Adaptive DE-normalization (TADE) to style a noise vector with the acoustic features of the target speech (e.g., mel-spectrogram) via instance normalization and elementwise modulation. More precisely it learns adaptively the modulation parameters γ and β from the acoustic features, and then applies the transformation $\gamma \odot c + \beta$, where c is the normalized content of the input activation. For efficient training, multiple random-window discriminators adversarially evaluate the speech signal analyzed by a set of Pseudo-Quadrature Mirror Filters (PQMF) [15] filter banks, with the generator regularized by a multi-resolution STFT loss. All

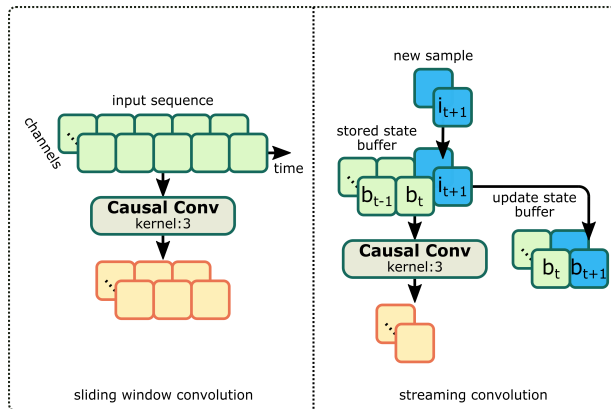


Figure 2: Diagrams for non-streaming convolution (left) and streaming convolution (right)

convolutions in StyleMelGAN are non-causal and run as a moving-average on sliding windows of the input tensors. This results in significant amount of algorithmic delay due to the deep hierarchical structure of the model. In the following, we describe major modifications to this baseline model that enable the generation at very low delay with different acoustic features for conditioning.

3.2. Streamwise Convolution

There are two requirements to operate a convolutional model in streaming manner with low algorithmic delay. First, the dependency on future inputs to predict the current output should be as low as possible. We achieve this by enforcing all convolutions in StyleMelGAN to be causal so that the model has zero delay. The second requirement is to generate the output frame by frame, as the new input information is available. This condition is fulfilled in StyleMelGAN by adding an internal memory buffer to the causal convolutions in inference mode, as described in [16] and illustrated in Figure 2. Each causal convolution stores a buffer containing the last input samples used for generating the previous output frame, and then reused once the new input sample is available. By applying the above modifications to StyleMelGAN, we obtain *Streamwise StyleMelGAN (SSMGAN)*, which is able to generate speech signals frame by frame with no delay between the input conditioning features and the output waveform.

3.3. Channel Normalization

It is not feasible to run instance normalization [17] in SSMGAN as the normalization statistics are estimated along the temporal dimension of the input activations. We replace instance normalization with channel normalization [18], that estimates the statistics along the channel dimension instead. Interestingly, we found this normalization maintains the model performance and keeps the training fast. It also avoids the creation of subtle clicking artifacts that sometimes occur when training StyleMelGAN with instance normalization on a multi-speaker dataset.

3.4. Modified TADE Residual Block

The TADE residual blocks are slightly modified from the original model, as shown in Figure 3. The complexity in SSMGAN is re-

duced by using a single TADE conditioning layer and applying the same modulation parameters β and γ twice rather than having two separate TADEs in the residual block. With this modification, the total number of model parameters reduces from 3.86 M to 2.73 M.

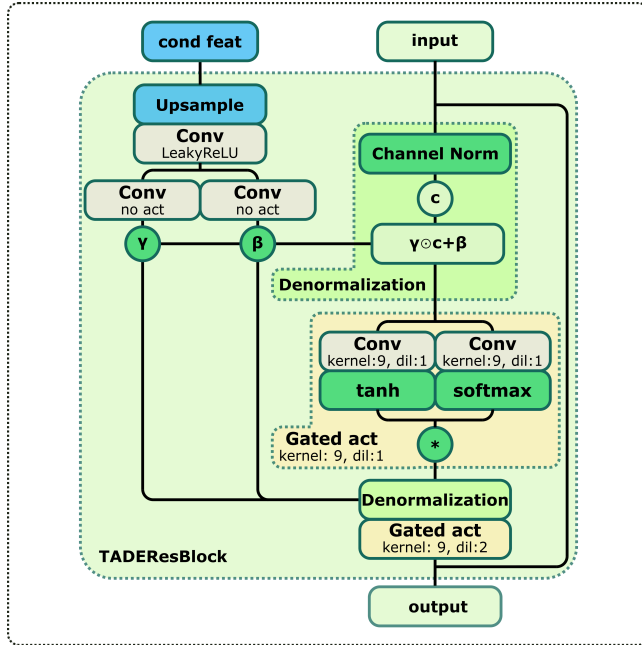


Figure 3: Modified TADE residual block for the SSMGAN.

3.5. Multiband Generation

SSMGAN further reduces the complexity compared to the baseline model by introducing multiband synthesis as in [19, 20]. Rather than synthesizing the whole band of the speech signal in time domain at the output sampling rate f_s , the generator outputs simultaneously different sub-bands sampled at f_s/N Hz, with $N = 4$ and $f_s = 16$ kHz. By design, SSMGAN generates the sub-bands as an N -channels output, which is then fed to a PQMF synthesis filter-bank to obtain a frame of synthesized speech. Since the PQMF uses a filter prototype with 50% of overlap, it incurs a delay of 1 frame.

3.6. Conditioning on Coded LPCNet Features

Finally, we condition SSMGAN with coded parameters in real-time to run as a speech decoder. Instead of providing the mel-spectrogram as an intermediate representation, the coded parameters obtained by the LPCNet encoder at 1.6 kbit/s are introduced to the generator network. The pitch lag was found to be critical for high-quality synthesis, and hence it is processed separately from the rest of the conditioning information. More precisely, the coded cepstral and energy parameters are passed through a simple causal convolutional layer to obtain an 80 channel representation used for conditioning the generation from the prior signal. This prior is not created from latent random noise, but rather from a learned embedding of the pitch lag which is then multiplied elementwise by the pitch correlation. Figure 4 shows the complete architecture of the proposed SSMGAN conditioned on the LPCNet coded parameters. With this setting, SSMGAN can generate wideband speech frames

of 10 m sec length and total delay of 55 m sec, where 45 m sec is introduced by the original extraction of the LPCNet coding packets, while 10 m sec are added by the PQMF synthesis filter-bank.

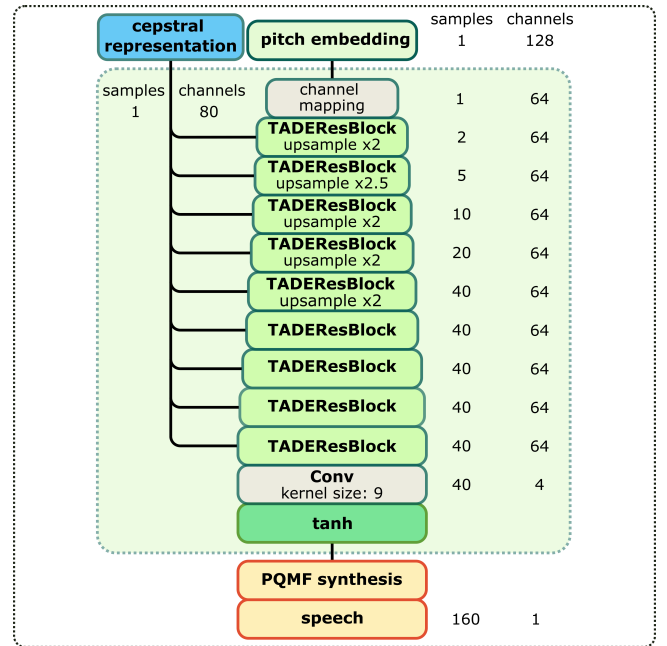


Figure 4: The SSMGAN generator. Dimensions are given for generating 1 frame at 16 kHz sampling rate. The cepstral coefficients pass through a simple convolutional layer to obtain a representation of 80 channels.

4. EXPERIMENTS

4.1. Experimental setup

The training procedure and hyperparameters are very similar to the ones described in [10]. We train SSMGAN using one NVIDIA Tesla V100 GPU on the VCTK corpus [21] at 16 kHz. The conditioning features are calculated as in [6] as described in Section 2. The generator is pretrained for 200k steps using Adam optimizer [22] with learning rate $lr_g = 10^{-4}$, $\beta = \{0.5, 0.9\}$. When starting the adversarial training, we set $lr_g = 5 * 10^{-5}$ and use the multi-scale discriminator described in [8] trained via Adam optimizer with $lr_d = 2 * 10^{-4}$, and same β . The batch size is 32 and for each sample in the batch we extract a segments of length 1 s. The adversarial training lasts for about 1.5 M steps.

4.2. Subjective evaluation

We conducted a subjective listening test following the ITU-R MUSHRA [23] recommendation comparing classical and neural speech coders. The test set is composed of 12 utterances by 10 different speakers in 4 different languages. All speakers and 3 out of 4 languages are unseen during training. Most of the utterances (10 out of the 12) are coming from unseen proprietary databases. The obtained results with 16 expert listeners are shown in Figure 5.

The anchor is generated using the Speex speech decoder employed at a bit rate of 4 kbit/s. Two state-of-the-art neural decoders

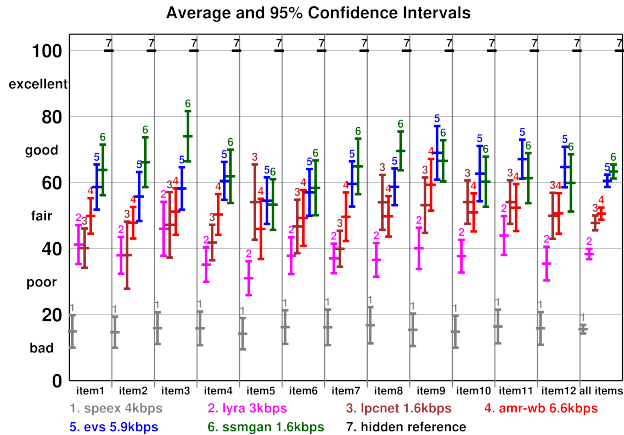


Figure 5: MUSHRA listening test results using t -distribution.

were considered: LPCNet at 1.6 kbit/s and Lyra at 3 kbit/s, as well as two classical but still widely used codecs: AMR-WB [24] at 6.6 kbit/s and the recent 3GPP EVS [25] at 5.9 kbit/s. The condition Lyra at 3 kbit/s was generated using the release v0.0.1 [26] with the default setting. EVS at 5.9 kbit/s works with a variable bit rate (VBR) and that 5.9 kbit/s reflects the average bit rate on active frames. During a long inactive phase, EVS switches to a non-transmission mode (DTX), transmitting only periodically packets at a bit rate as low as 2.4 kbit/s. Since the test items only contain short pauses between sentences, the DTX mode plays a minor role in this test.

LPCNet was trained on the VCTK dataset. One difference from the original work is that we do not apply a domain adaptation by first training on unquantized and then fine-tuning on quantized features, since this was found to make no difference on VCTK. In addition, since VCTK is noisier and much more diverse than the NTT database used in the original work, we removed the data augmentation since it was found to be detrimental to the final quality¹. The publicly available version of the Lyra model was not retained on VCTK, and hence it is not directly comparable with SSMGAN or LPCNet in this case. It was nonetheless taken into consideration as it offers a reproducible benchmark.

4.3. Objective evaluation

Our solution was also compared to the other neural decoders using different objective metrics. Since it is known that objective speech quality models like POLQA [27] are not reliable for non-waveform-preserving coding schemes, and in particular for neural decoders, we also considered the newly introduced objective metric WARP-Q [28], which was designed for this purpose. STOI [29], assessing the speech intelligibility, is also added, and the scores measured on 824 test items of VCTK are reported in Table 2.

SSMGAN at 1.6 kbit/s scores the best among the neural coding solutions across all three metrics, which is in agreement with the subjective listening test. The results of our MUSHRA listening test show moreover that these objective metrics do not fully reflect the perceived quality of the generated speech, disproportionately disfavouring generative models.

¹Check our demo samples at the following url: https://fhgspco.github.io/ssmgan_spco/

Speech decoders	POLQA	STOI	WARP-Q
Speex 4 kbit/s	2.022	0.720	1.074
AMR-WB 6.6 kbit/s	3.202	0.863	0.784
EVS 5.9 kbit/s	3.675	0.890	0.805
LPCNet 1.6 kbit/s	2.628	0.777	0.915
Lyra 3 kbit/s	2.649	0.794	0.958
SSMGAN 1.6 kbit/s	2.719	0.830	0.826

Table 2: Average objective scores for neural decoders. For POLQA-MOS and STOI higher scores are better, while for WARP-Q lower scores are better (confidence intervals are negligible).

4.4. Complexity

The main contribution to SSMGAN’s computational complexity stems from the convolutions in the TADEResBlocks and the upsampling layers. If L denotes the channel dimension, K the size of the convolutional kernels, and F the dimension of the input features, then (ignoring activations and lower order terms) the evaluation of a TADEResBlock takes $(F+5L)LK$ multiply accumulate operations (MAC) per output sample. Furthermore, an upsampling layer with kernel size K and channel dimension L takes L^2K MAC. With $L = 64$, $K = 9$, $F = 80$ and TADEResBlock output sampling rates of 100, 200, 500, 1000, 2000, 4000, 4000, and 4000 Hz this accumulates to

$$(80 + 5 \cdot 64) \cdot 64 \cdot 9 \cdot (100 + 200 + 500 + 1000 + 2000 + 4 \cdot 4000) + 64^2 \cdot 9 \cdot (200 + 500 + 1000 + 2000 + 4000) \approx 4.8 \text{ GMACs.}$$

A comparison with other neural vocoders used for neural speech coding is given in Table 3. It should be noted, that the convolutional structure of SSMGAN allows for efficient parallel execution, which gives it a decisive advantage over autoregressive models on GPUs. The current unoptimized PyTorch implementation achieves about real-time frame-by-frame inference using four cores of an Intel(R) Core(TM) i7-6700 3.40GHz CPU. The above complexity calculations show that the next step will be to work on an efficient implementation for mobile devices, which will be the object of a future work.

Model	Complexity
SSMGAN (ours)	4.8 GMACs
LPCNet [2]	1.5 GMACs
Multi-band WaveRNN [19]	2.75 GMACs

Table 3: Complexity of common neural vocoders for speech coding.

5. CONCLUSION

In this paper we introduce SSMGAN, a neural speech decoder generating state-of-the-art quality with low delay, complexity, and working at very low bit rate. We assess the quality against existing neural autoregressive models and modern speech codecs at low bit rate, with both objective scores and subjective listening tests. We show for the first time that GAN-vocoders can perform fast streaming speech synthesis with low algorithmic delay, and that they can achieve high quality synthesis when conditioned on compact parametric speech representations.

6. REFERENCES

- [1] W. B. Kleijn, F. S. C. Lim, A. Luebs, J. Skoglund, F. Stimberg, Q. Wang, and T. C. Walters, "WaveNet Based Low Rate Speech Coding," in *ICASSP 2018, IEEE International Conference on Acoustics, Speech and Signal Processing*, 2018, pp. 676–680.
- [2] J. Valin and J. Skoglund, "A Real-Time Wideband Neural Vocoder at 1.6 kb/s Using LPCNet," in *INTERSPEECH 2019, 20th Annual Conference of the International Speech Communication Association*, 2019, pp. 3406–3410.
- [3] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, *et al.*, "WaveNet: A Generative Model for Raw Audio," *arXiv:1609.03499*, 2016.
- [4] R. Prenger, R. Valle, and B. Catanzaro, "WaveGlow: A Flow-based Generative Network for Speech Synthesis," in *ICASSP 2019, IEEE International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 3617–3621.
- [5] N. Kalchbrenner, E. Elsen, K. Simonyan, Noury, *et al.*, "Efficient neural audio synthesis," in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, J. Dy and A. Krause, Eds., vol. 80. PMLR, 10–15 Jul 2018, pp. 2410–2419.
- [6] J. Valin and J. Skoglund, "LPCNet: Improving Neural Speech Synthesis through Linear Prediction," in *ICASSP 2019, IEEE International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 5891–5895.
- [7] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, *et al.*, "Generative Adversarial Nets," in *Advances in NeurIPS 27*, 2014, pp. 2672–2680.
- [8] K. Kumar, R. Kumar, de T. Boissiere, L. Gestein, *et al.*, "MelGAN: Generative Adversarial Networks for Conditional Waveform Synthesis," in *Advances in NeurIPS 32*, 2019, pp. 14 910–14 921.
- [9] J. Kong, J. Kim, and J. Bae, "Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33, 2020, pp. 17 022–17 033.
- [10] A. Mustafa, N. Pia, and G. Fuchs, "Stylemelgan: An efficient high-fidelity adversarial vocoder with temporal adaptive normalization," in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 6034–6038.
- [11] R. Yamamoto, E. Song, and J. Kim, "Parallel WaveGAN: A Fast Waveform Generation Model Based on Generative Adversarial Networks with Multi-Resolution Spectrogram," in *ICASSP 2020, IEEE International Conference on Acoustics, Speech and Signal Processing*, 2020, pp. 6199–6203.
- [12] M. Bińkowski, J. Donahue, *et al.*, "High fidelity speech synthesis with adversarial networks," in *International Conference on Learning Representations*, 2020.
- [13] J. Klejsa, P. Hedelin, C. Zhou, R. Fejgin, and L. Villemoes, "High-quality Speech Coding with SampleRNN," in *ICASSP 2019, IEEE International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 7155–7159.
- [14] W. Kleijn, A. Storus, M. Chinen, T. Denton, F. Lim, A. Luebs, J. Skoglund, and H. Yeh, "Generative speech coding with predictive variance regularization," in *ICASSP 2021, IEEE International Conference on Acoustics, Speech and Signal Processing*, 2021.
- [15] T. Q. Nguyen, "Near-perfect-reconstruction pseudo-QMF banks," *IEEE Transactions on Signal Processing*, vol. 42, no. 1, pp. 65–76, 1994.
- [16] O. Rybakov, N. Kononenko, N. Subrahmanya, M. Visontai, and S. Laurenzo, "Streaming Keyword Spotting on Mobile Devices," in *Proc. Interspeech 2020*, 2020, pp. 2277–2281.
- [17] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Instance normalization: The missing ingredient for fast stylization," *arXiv:1607.08022*, 2016.
- [18] F. Mentzer, G. D. Toderici, M. Tschannen, and E. Agustsson, "High-fidelity generative image compression," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [19] C. Yu, H. Lu, N. Hu, M. Yu, *et al.*, "DurIAN: Duration Informed Attention Network for Speech Synthesis," in *Proc. Interspeech 2020*, 2020, pp. 2027–2031. [Online]. Available: <http://dx.doi.org/10.21437/Interspeech.2020-2968>
- [20] G. Yang, S. Yang, K. Liu, *et al.*, "Multi-band melgan: Faster waveform generation for high-quality text-to-speech," in *2021 IEEE Spoken Language Technology Workshop (SLT)*, 2021, pp. 492–498.
- [21] J. Yamagishi, C. Veaux, and K. MacDonald, "CSTR VCTK Corpus: English Multi-speaker Corpus for CSTR Voice Cloning Toolkit," 2019.
- [22] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *ICLR*, 2015.
- [23] R. BS.1534, "Method for the subjective assessment of intermediate quality levels of coding systems," Tech. Rep., 2003.
- [24] 3GPP, "Speech codec speech processing functions; Adaptive Multi-Rate - Wideband (AMR-WB) speech codec; Transcoding functions," 3rd Generation Partnership Project (3GPP), TS 26.190, 12 2009. [Online]. Available: <http://www.3gpp.org/ftp/Specs/html-info/26190.htm>
- [25] —, "TS 26.445, EVS Codec Detailed Algorithmic Description; 3GPP Technical Specification (Release 12)," 3rd Generation Partnership Project (3GPP), TS 26.445, 12 2014. [Online]. Available: <http://www.3gpp.org/ftp/Specs/html-info/26445.htm>
- [26] (2021) Google/lyra. [Online]. Available: <https://github.com/google/lyra/releases/tag/v0.0.1>
- [27] J. Beerends, C. Schmidmer, J. Berger, M. Obermann, R. Ullmann, J. Pomy, and M. Keyhl, "Perceptual Objective Listening Quality Assessment (POLQA), the third generation ITU-T standard for end-to-end speech quality measurement part I — temporal alignment," *journal of the audio engineering society*, vol. 61, no. 6, pp. 366–384, june 2013.
- [28] W. A. Jassim, J. Skoglund, M. Chinen, and A. Hines, "WARP-Q: Quality Prediction For Generative Neural Speech Codecs," in *ICASSP 2021, IEEE International Conference on Acoustics, Speech and Signal Processing*, 2021.
- [29] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "Algorithm for intelligibility prediction of timefrequency weighted noisy speech," *IEEE Trans. Audio Speech Lang. Process.*, vol. pp. 2125–2136, 2011.