

NESC: Robust Neural End-2-End Speech Coding with GANs

Nicola Pia, Kishan Gupta, Srikanth Korse, Markus Multrus, Guillaume Fuchs

Fraunhofer IIS Erlangen

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

Abstract

Neural networks have proven to be a formidable tool to tackle the problem of speech coding at very low bit rates. However, the design of a neural coder that can be operated robustly under real-world conditions remains a major challenge. Therefore, we present Neural End-2-End Speech Codec (NESC) a robust, scalable end-to-end neural speech codec for high-quality wideband speech coding at 3 kbps. The encoder uses a new architecture configuration, which relies on our proposed DualPathConvRNN (DPCRNN) layer, while the decoder architecture is based on our previous work Streamwise-StyleMelGAN. Our subjective listening tests on clean and noisy speech show that NESC is particularly robust to unseen conditions and signal perturbations.

Index Terms: neural speech coding, Generative Adversarial Network, residual quantization.

1. Introduction

Very low bit rate speech coding is extremely challenging for classical coding techniques. The paradigm usually employed is parametric coding, which yields intelligible speech at the cost of poor audio quality and unnatural synthesized speech. Recent advances in neural networks are filling this gap, enabling high-quality speech coding at very low bit rates. We categorize the possible solutions to this problem according to the role played by neural networks.

- level 1 *post-filtering*: A neural network based post-processor is employed at the end of a conventional encoder-decoder chain, in order to improve the quality of the coded speech. This enables the enhancement of existing communication systems with minimal effort.
- level 2 *neural decoder*: A conventional encoder model generates a bitstream, which is decoded using a neural network. This enables backward compatible decoding of existing bitstreams.
- level 3 *end-2-end*: Both encoder and decoder are neural networks, which are trained jointly. The input of the encoder is the speech waveform, and the quantization is jointly learned, hence obtaining directly the optimal bitstream for the signal.

Level 1 approaches such as [1, 2, 3, 4, 5] are minimally invasive, as they can be deployed over existing pipelines. Unfortunately they still suffer typical unpleasant artifacts, which are especially challenging to eliminate.

The first published level 2 speech decoder was based on WaveNet [6], and served as a proof of concept. Several follow-up works [7, 8] improved quality and computational complexity, and [9] presented LPCNet, a low complexity decoder which synthesizes good quality clean speech at 1.6 kbps. We have shown in our previous work [10] that the same bitstream used in LPCNet can be decoded by Streamwise-StyleMelGAN

(SSMGAN), a feed-forward GAN model, which provides significantly better quality.

All of these models produce high-quality clean speech, but are not robust in the presence of noise and reverberation. Lyra [11] was the first model to directly address this problem. Overall, it seems that the generalization capabilities and the quality of level 2 models are partly weakened by the limitations of the classical representation of speech at the encoder side.

Many approaches tackling the problem from the perspective of a level 3 solution were proposed [12, 13, 14, 15], but these models usually do not target very low bit rates.

SoundStream [16] was the first fully end-to-end approach, operating at low bit rates and robust under many different noise conditions. It is built on a U-Net convolutional encoder-decoder, without skip connections, and using a residual quantization in the bottleneck. According to the authors' evaluation SoundStream is stable under a wide range of real-life coding scenarios. Moreover, it permits to synthesize speech at bit rates ranging from 3 kbps to 12 kbps. Finally, SoundStream works at 24 kHz, implements a noise reduction mode, and can also code music. More recently the work [17] presented another level 3 solution using a different set of techniques.

We present NESC, a new model capable of robustly coding wideband speech at 3 kbps. The architecture behind NESC is fundamentally different from SoundStream and is the main aspect of novelty of our approach. The encoder architecture is based on our proposed DPCRNN, which uses a sandwich of convolutional and recurrent layers to efficiently model intra-frame and inter-frame dependencies. The DPCRNN layer is followed by a series of convolutional residual blocks with no downsampling and by a residual quantization. The decoder architecture is composed of a recurrent neural network followed by the decoder of SSMGAN.

We show that data augmentation can significantly improve robustness against a wide range of different types of noises and reverberation. We extensively test our model with many types of signal perturbations and unseen speakers as well as unseen languages. Moreover, we analyze the unsupervised speech signal clustering achieved by the latent representation. Our contributions are the following:

- We introduce NESC, a new end-to-end neural codec for speech.
- We present the DPCRNN layer, which offers an efficient way of exploiting intra and inter-frame dependencies, for learning a latent representation suitable for quantization.
- We analyze some interesting clustering behaviour exhibited by NESC's quantized latent.
- We demonstrate NESC's robustness against many types of noise and reverberation, via objective and subjective evaluations.

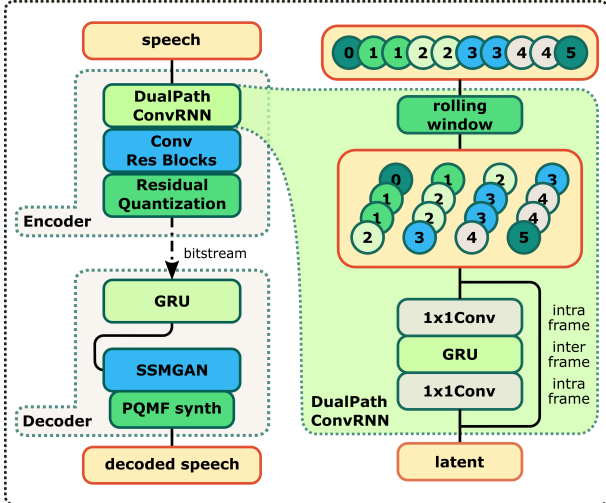


Figure 1: NESc’s high level architecture and the DPCRNN.

2. Proposed Architecture

As illustrated in Fig. 1, the proposed model consists of a learned encoder, a learned quantization layer and a recurrent pre-net followed by a SSMGAN decoder.

The encoder architecture counts 2.09 M parameters, whereas the decoder has 3.93 M parameters. The encoder rarely reuses the same parameters in computation, as we hypothesize that this favors generalization. It runs around 40x faster than real time on a single thread of an Intel(R) Core(TM) i7-6700 CPU at 3.40GHz. The decoder runs around 2x faster than real time on the same architecture, despite only having double as many parameters as the encoder. Our implementations and design are not optimized for inference speed.

2.1. Encoder

The encoder architecture relies on our newly proposed DPCRNN, which was inspired by [18]. This layer consists of a rolling window operation followed by a 1x1-convolution with 512 channels, a GRU with 128 hidden dimensions, and finally another 1x1-convolution with 256 channels, all activated via LeakyReLUs. The rolling window transform reshapes the input signal of shape $[1, t]$ into a signal of shape $[s, f]$, where s is the length of a frame and f is the number of frames. We use frames of 10 ms with 5 ms from the past frame and 5 ms lookahead. For 2 s of audio at 16 kHz this results in $s = 80 + 160 + 80 = 320$ samples and $f = 200$. The 1x1-convolutional layers model the time dependencies within each frames, i.e. intra-frame dependencies, whereas the GRU models the dependencies between different frames, i.e. inter-frame dependencies. This approach allows us to avoid downsampling via strided convolutions or interpolation layers, which in early experiments were shown to strongly affect the final quality of the audio synthesized by SSMGAN.

The rest of the encoder architecture consists of 4 residual blocks, each of which consists of a 1d-convolution with kernel size 3 followed by a 1x1-convolution, both with 256 channels and activated via LeakyReLUs. The use of the DPCRNN provides a compact and efficient way to model the temporal dependencies of the signal, hence making the use of dilation or other tricks for extending the receptive field of the residual blocks un-

necessary.

2.2. Quantization

The encoder architecture produces a latent vector of dimension 256 for each packet of 10 ms. This vector is then quantized using a learned residual vector quantizer based on Vector-Quantized VAE (VQ-VAE) [19] as in [16]. In a nutshell, this quantizer learns multiple codebooks on the vector space of the encoder latent packets. The first codebook approximates the latent output of the encoder $z = E(x)$ by the closest entry of the codebook z_e . The second codebook does the same on the residual of the quantization, i.e. on $z - z_e$, and so on for the following codebooks. This technique is well-known in classical coding, and permits to use the vector space structure of the latent to code many more points with significantly less complexity than by using a single codebook of equivalent bit rate.

In NESc we use a residual quantizer with three codebooks each at 10 bits to code a packet of 10 ms, hence resulting in a total of 3 kbps. It is important to note that during inference, it is possible to drop one or two codebook indices and get an approximation of the final output. NESc delivers then a scalable bitstream from 1 kbps to 3 kbps. Since we did not explicitly train it for the intermediate bit rates, we did not include them in our final subjective evaluation, even though the signal is coherent and intelligible at each bit rate.

2.3. Decoder

The decoder architecture is composed of a recurrent neural network followed by a SSMGAN decoder. We use a single causal GRU layer as a pre-net in order to prepare the bitstream before feeding it to the SSMGAN decoder. We do not apply significant modifications to the SSMGAN decoder, except for the use of a constant prior signal and the use of 256 conditioning channels provided by the output of the GRU. We refer to [10] for more details on this architecture. Briefly, this is a convolutional decoder, which is based on Temporal Adaptive DE-normalization layers (TADE), similar to the FiLM layers used in [16]. It up-samples the bitstream with very low upsampling scales, and provides the conditioning information at each layer of upsampling, while shaping the signal via softmax-gated tanh activations.

SSMGAN outputs four Pseudo Quadrature Mirror Filterbank (PQMF) [20] sub-bands, which are then synthesized using a synthesis filter. This filter has 50 samples of lookahead, effectively introducing one frame of delay in our implementation. The total delay of our system is then 25 ms, 15 ms from the encoder and the framing and 10 ms from the decoder.

3. Evaluation

3.1. Experimental setup

We train NESc on the 260 hours of speech from the LibriTTS Dataset [21] at 16 kHz. The clean speech samples are augmented by adding background noise from the DNS Challenge Dataset [22] at a random SNR between 0 dB and 50 dB, and also adding reverberation by convolving real or generated room impulse responses (RIRs) from the SLR28 Dataset [23, 24].

The training of NESc is very similar to the training of SSMGAN as described in [10]. We first pre-train encoder and decoder together for around 500k iterations having the spectral reconstruction loss of [25] and the MSE loss as objective. We then turn on the adversarial loss and the discriminator feature

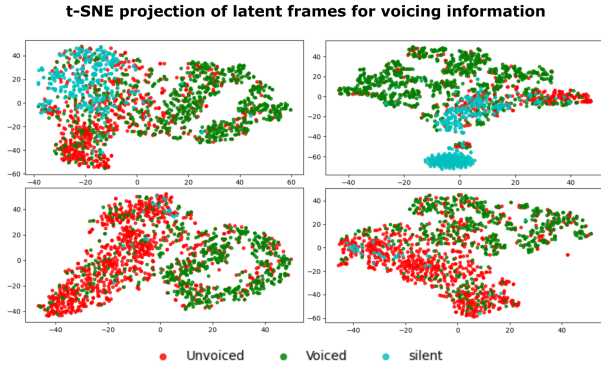


Figure 2: Voiced and unvoiced frames are clustered. Each subplot represents a different speaker selected at random.

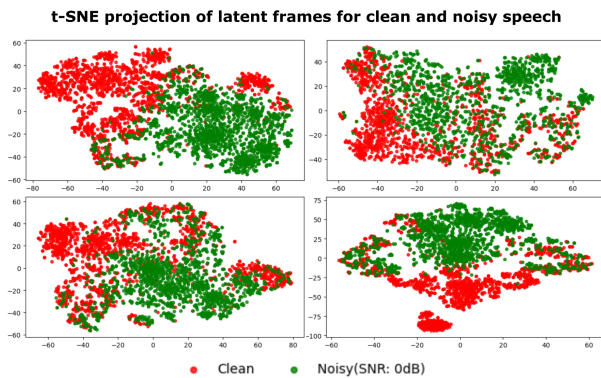


Figure 3: Noisy and clean frames are clustered. Each subplot represents a different speaker selected at random.

losses from [26] and train for another 700k iterations; beyond that we have not seen substantial improvements. The generator is trained on audio segments of 2 s with batch size 64. We use an Adam [27] optimizer with learning rate $1 \cdot 10^{-4}$ for the pre-training of the generator, and decrease the learning rate to $5 \cdot 10^{-5}$ as soon as the adversarial training starts. We use an Adam optimizer with learning rate $2 \cdot 10^{-4}$ for the discriminator.

3.2. Qualitative statistical analysis of the latent

We provide a qualitative analysis of the distribution of the latent in order to give a better understanding of its behaviour in practice. The quantized latent frames are embedded in a space of dimension 256, hence in order to plot their distribution we use their t-SNE projections [28]. For each experiment, we first encode the audio with different recording conditions and we label each frame depending on a priori information regarding its acoustic or linguistic characteristics; finally we look for clusters in the low dimensional projections. Each subplot represents audio coming from different speakers randomly selected from the LibriTTS Dataset. Notice that the model is not trained with any clustering objective, hence any such behaviour shown at inference time is an emergent aspect of the training set up.

In our first experiment we test voicing information using a VAD algorithm to label each frame automatically. Fig. 2 shows

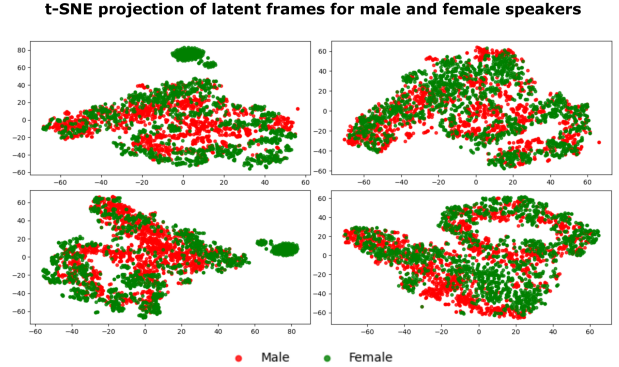


Figure 4: No particular clustering is shown. Each subplot represents a different speaker combination selected at random.

a clear clustering of voiced, unvoiced, and silent frames. We also analyze the clustering ability of the pitch level, but we do not observe a clear trend and therefore we do not add the plots due to lack of space.

In Fig. 3 noise injection during data augmentation is analyzed. We once again notice a clear division between noise frames and clean frames in the latent space, suggesting that the model is using distinct parts of the latent for these distinct modes.

Finally, we test linguistic and speaker dependent characteristics such as gender (Fig. 4), speaker id, and language (figures not shown due to lack of space). In these cases we do not observe any particular clusterings, suggesting that the model is not able to distinguish between these macro-level aspects.

We hypothesize that these clustering behaviours might reflect the compression strategy of the model, which would be in line with well-known heuristics already in use in classical codecs.

3.3. Objective scores

We evaluate NESC using several objective metrics. It is well-known that such metrics are not reliable for assessing the quality of neural codecs [6, 10], as they disproportionately favor waveform-preserving codecs. Nonetheless, we report their values for comparison purposes. We consider ViSQOL v3 [29], POLQA [30] and the speech intelligibility measure STOI [31].

The scores are calculated on two internally curated test sets, the StudioSet and the InformalSet, respectively in Table 1 and 2. The StudioSet is constituted of 108 multi-lingual samples from the NTT Multi-Lingual Speech Database for Telephony, totalling around 14 minutes of studio-quality recordings. The InformalSet is constituted of 140 multi-lingual samples scraped from several datasets including LibriVox, and totalling around 14 minutes of audio recordings. This test set includes samples recorded with integrated microphones, more spontaneous speech, sometimes with low background noise or reverberation from a small room. NESC scores the best among the neural coding solutions across all three metrics.

3.4. Subjective Evaluation

We subjectively test the model on challenging unseen conditions in order to assess its robustness. We select a test set of speech samples from the NTT Dataset comprising unseen

Table 1: Average objective scores for neural decoders on the StudioSet. For all metrics higher scores are better. Confidence intervals are negligible for POLQA and ViSQOL v3, while for STOI they are smaller than 0.02.

Codec	POLQA	STOI	ViSQOL v3
OPUS 6 kbps	1.681	0.480	2.273
EVS 5.9 kbps	3.308	0.553	3.036
SSMGAN 1.6 kbps	2.213	0.536	2.505
NESC 1 kbps	1.534	0.612	2.109
NESC 2 kbps	2.382	0.641	2.615
NESC 3 kbps	2.548	0.643	2.841

Table 2: Average objective scores for neural decoders on the InformalSet. For all metrics higher scores are better. Confidence intervals are negligible for POLQA and ViSQOL v3, while for STOI they are smaller than 0.025.

Codec	POLQA	STOI	ViSQOL v3
OPUS 6 kbps	1.833	0.613	2.357
EVS 5.9 kbps	3.486	0.736	3.071
SSMGAN 1.6 kbps	2.267	0.647	2.476
NESC 1 kbps	1.659	0.745	2.074
NESC 2 kbps	2.492	0.802	2.595
NESC 3 kbps	2.707	0.817	2.822

speakers, languages and recording conditions. In the naming convention "m" stands for male, "f" for female, "ar" for Arabic, "en" for English, "fr" for French, "ge" for German, "ko" for Korean, and "th" for Thai.

We also test the model on noisy speech, for this we select the same speech samples as for the clean speech test, and apply a similar augmentation policy as in Section 3.1. We mix ambient noises (e.g. airport noises, typing noises, ...) at SNRs between 10 dB and 30 dB, and then convolve with room impulse responses (RIR) coming from small, medium and large rectangular rooms. More precisely, "ar/f", "en/f", "fr/m", "ko/m", and "th/f" are convolved with RIRs from small rooms, and hence for these signals the reverberation does not play a big role; whereas the other samples are convolved with RIRs medium and large size rooms. Notice that the SNR interval is different from the one used in training as SNRs lower than 10 dB would have been too challenging to evaluate, and SNRs above 30 dB would have not provided enough perturbation in the signal.

We conduct two MUSHRA [32] listening tests, both involving 13 expert listeners. The test results for clean speech in Fig. 5 show that NESC is on par with SSMGAN and EVS [33] in this case. The test results for noisy speech in Fig. 6 confirm that SSMGAN is not robust to such scenarios, while NESC competes with EVS in this case.

The anchor for the tests is generated using OPUS [34] at 6 kbps, since the quality is expected to be very low at this bit rate. We took EVS [33] at 5.9 kbps nominal bit rate as good quality benchmark for the classical codecs. In order to avoid an influence of CNG frames with different signature on the test, we deactivated the DTX transmission. Finally, we test our solution against our previous neural decoder SSMGAN at 1.6 kbps. SSMGAN at 1.6 kbps is trained on the VCTK Corpus [35], hence the comparison with NESC is not completely fair. Early

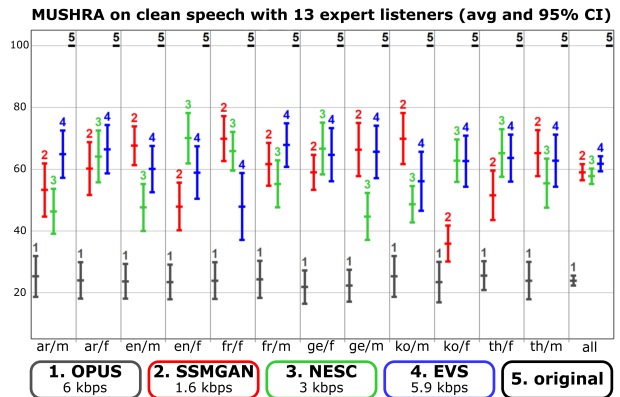


Figure 5: The listening test on clean speech shows that NESC is on par with EVS and SSMGAN.

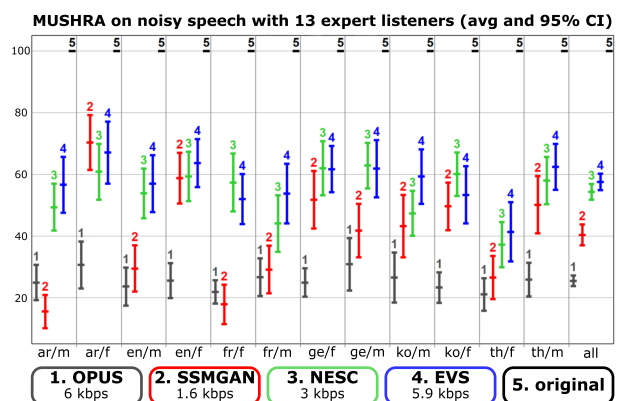


Figure 6: The listening test on noisy speech shows that NESC is robust under very challenging conditions.

experiments showed that training SSMGAN at 1.6 kbps with noisy data is more challenging than expected. We suppose that this issue is due to its reliance on the pitch information, which might be challenging to estimate in noisy environments. For this reason we decided to test NESC against the best neural clean speech decoder that we have access to, namely SSMGAN trained on VCTK, and still add it to the noisy speech test as an additional condition to show its limitations.

Both tests show that NESC is on par with EVS, while having half of its bit rate. Fig. 6 shows the limitations of SSMGAN when working with noisy signals and confirms that the quality of NESC stays high even in these challenging conditions¹.

4. Conclusions

We present NESC, a new GAN model capable of high-quality and robust end-to-end speech coding. It relies on the new DPCRNN as the main building block for efficient and reliable encoding. We test our setup via objective quality measures and subjective listening tests, and show that NESC is robust under various types of noise and reverberation. We show a qualitative analysis of the latent structure giving a glimpse of the internal workings of our codec. Future work will be directed toward further complexity reduction and quality improvements.

¹Check our demo samples at: <https://fhgspco.github.io/nesc/>

5. References

- [1] Z. Zhao, H. Liu, and T. Fingscheidt, "Convolutional Neural Networks to Enhance Coded Speech," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 4, pp. 663–678, April 2019.
- [2] J. Skoglund and J. Valin, "Improving Opus Low Bit Rate Quality with Neural Speech Synthesis," in *INTERSPEECH*, 2020.
- [3] S. Korse, K. Gupta, and G. Fuchs, "Enhancement of Coded Speech Using a Mask-Based Post-Filter," in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6764–6768.
- [4] A. Biswas and D. Jia, "Audio Codec Enhancement with Generative Adversarial Networks," in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 356–360.
- [5] S. Korse, N. Pia, K. Gupta, and G. Fuchs, "PostGAN: A GAN-Based Post-Processor to Enhance the Quality of Coded Speech," *arXiv preprint arXiv:2201.13093*, 2022.
- [6] 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.
- [7] C. Gărbacea, A. van den Oord, Y. Li, F. S. Lim, A. Luebs, O. Vinyals, and T. C. Walters, "Low bit-rate speech coding with VQ-VAE and a WaveNet decoder," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 735–739.
- [8] 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.
- [9] 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.
- [10] A. Mustafa, J. Bütche, S. Korse, K. Gupta, G. Fuchs, and N. Pia, "A Streamwise Gan Vocoder for Wideband Speech Coding at Very Low Bit Rate," in *2021 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 2021, pp. 66–70.
- [11] 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.
- [12] S. Morishima, H. Harashima, and Y. Katayama, "Speech coding based on a multi-layer neural network," in *IEEE International Conference on Communications, Including Supercomm Technical Sessions*. IEEE, 1990, pp. 429–433.
- [13] S. Kankanahalli, "End-to-end optimized speech coding with deep neural networks," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 2521–2525.
- [14] K. Zhen, J. Sung, M. Lee, S. Beack, and M. Kim, "Cascaded Cross-Module Residual Learning Towards Lightweight End-to-End Speech Coding," *Proc. Interspeech 2019*, 2019.
- [15] K. Zhen, M. Lee, J. Sung, S. Beack, and M. Kim, "Efficient and scalable neural residual waveform coding with collaborative quantization," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 361–365.
- [16] N. Zeghidour, A. Luebs, A. Omran, J. Skoglund, and M. Tagliasacchi, "SoundStream: An End-to-End Neural Audio Codec," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 1–1, 2021.
- [17] X. Jiang, X. Peng, C. Zheng, H. Xue, Y. Zhang, and Y. Lu, "End-to-End Neural Audio Coding for Real-Time Communications," *arXiv preprint arXiv:2201.09429*, 2022.
- [18] Y. Luo, Z. Chen, and T. Yoshioka, "Dual-Path RNN: Efficient Long Sequence Modeling for Time-Domain Single-Channel Speech Separation," in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 46–50.
- [19] A. van den Oord, O. Vinyals, and K. Kavukcuoglu, "Neural discrete representation learning," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 6309–6318.
- [20] T. Q. Nguyen, "Near-perfect-reconstruction pseudo-QMF banks," *IEEE Transactions on Signal Processing*, vol. 42, no. 1, pp. 65–76, 1994.
- [21] H. Zen, V. Dang, R. Clark, Y. Zhang, R. Weiss, Y. Jia, Z. Chen, and Y. Wu, "LibriTTS: A corpus derived from LibriSpeech for text-to-speech," *arXiv preprint arXiv:1904.02882*, 2019.
- [22] C. Reddy, H. Dubey, V. Gopal, R. Cutler, S. Braun, H. Gampfer, R. Aichner, and S. Srinivasan, "ICASSP 2021 Deep Noise Suppression Challenge," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 6623–6627.
- [23] "OpenSLR 28 Dataset," <https://www.openslr.org/28/>.
- [24] T. Ko, V. Peddinti, D. Povey, M. Seltzer, and S. Khudanpur, "A study on data augmentation of reverberant speech for robust speech recognition," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 5220–5224.
- [25] 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.
- [26] K. Kumar, R. Kumar, de T. Boissiere, L. Gestin *et al.*, "MelGAN: Generative Adversarial Networks for Conditional Waveform Synthesis," in *Advances in NeurIPS 32*, 2019, pp. 14 910–14 921.
- [27] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *ICLR*, 2015.
- [28] L. V. der Maaten and G. Hinton, "Visualizing data using t-SNE," *Journal of machine learning research*, vol. 9, no. 11, 2008.
- [29] M. Chinen, F. S. Lim, J. Skoglund, N. Gureev, F. O’Gorman, and A. Hines, "ViSQOL v3: An open source production ready objective speech and audio metric," in *2020 twelfth international conference on quality of multimedia experience (QoMEX)*. IEEE, 2020, pp. 1–6.
- [30] 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.
- [31] 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.
- [32] R. BS.1534, "Method for the subjective assessment of intermediate quality levels of coding systems," *Tech. Rep.*, 2003.
- [33] 3GPP, "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>
- [34] K. Vos, V. K. K. Sørensen, S. Jensen, and J. Valin, "Voice coding with opus," in *Audio Engineering Society Convention 135*. Audio Engineering Society, 2013.
- [35] J. Yamagishi, C. Veaux, and K. MacDonald, "CSTR VCTK Corpus: English Multi-speaker Corpus for CSTR Voice Cloning Toolkit," 2019.