

ESC: Efficient Speech Coding with Cross-Scale Residual Vector Quantized Transformers

Yuzhe Gu^{1,2}, Enmao Diao²

¹University of Pennsylvania, Philadelphia, PA

²Duke University, Durham, NC

tracygu@seas.upenn.edu, enmao.diao@duke.edu

Abstract

Existing neural audio codecs usually sacrifice computational complexity for audio quality. They build the feature transformation layers mainly on convolutional blocks, which are not inherently appropriate for capturing local redundancies of audio signals. As compensation, either adversarial losses from a discriminator or a large number of model parameters are required to improve the codec. To that end, we propose Efficient Speech Codec (ESC)¹, a lightweight parameter-efficient codec laid on cross-scale residual vector quantization and transformers. Our model leverages mirrored hierarchical window-attention transformer blocks and performs step-wise decoding from coarse-to-fine feature representations. To enhance codebook utilization, we design a learning paradigm that involves a pre-training stage to assist with codec training. Extensive results show that ESC can achieve high audio quality with much lower complexity, which is a prospective alternative in place of existing codecs.

Index Terms: neural speech coding, vector quantization, efficient discrete learning

1. Introduction

Audio codecs are desired to compress audio signals within minimal computational bits to remove redundancies, while maintaining contents and fidelity in a low-latency manner. Recent advancements in deep learning have demonstrated the superiority of neural codecs over traditional ones, which rely on complex expert design and psycho-acoustic knowledge [1, 2]. Incorporating generative models like WaveNet [3] and SampleRNN [4] into audio codecs has shown promising results. These models act as powerful decoders, generating audio conditioned on intermediate representations produced by traditional codecs [5, 6]. However, their auto-regressive decoding nature typically introduces greater inference latency. Alternatively, some end-to-end neural audio codecs adopt the vector quantization (VQ) framework introduced in [7]. This approach discretizes continuous vectors using a codebook and employs straight-through estimation (STE) [8] to handle the non-differentiable quantization. Among the recent VQ-based audio codecs, notable examples include SoundStream [9], EnCodec [10], and Descript’s audio codec (DAC) [11]. These models typically utilize convolutional encoder and decoder backbones, supplemented by residual vector quantization (RVQ) [12], which recursively quantize previous residuals at the bottleneck representation. They also leverage GANs to discriminate multi-scale waveforms and spectrograms during training [13, 14, 15], producing high-fidelity audio reconstructions. SoundStream serves as the very first universal codec for diverse audio types, while EnCodec enhances

compression by integrating a lightweight language model at the bottleneck. DAC, using similar approaches, alleviates quantization dropout side-effects, a technique enhancing codec scalability, and introduces a periodic inductive bias in its activation functions [16, 17]. They present a remarkable neural audio codec that significantly outperforms EnCodec in both quality and compression ratio. Despite these advancements, we found that DAC relies heavily on its GAN discriminator to produce high-fidelity audios, which introduces inherent optimization challenges in adversarial learning. Furthermore, Descript’s audio codec faces computational limitations due to its large parameter count. To address this, we propose a shift towards a more parameter-efficient codec by adopting a cross-scale residual vector quantization (CS-RVQ) approach initially presented in [18] and replacing conventional convolutional layers with efficient Swin-Transformer Blocks (STBs) [19] as an architectural improvement.

Apart from these challenges, a significant issue in training VQ networks, known as *codebook collapse*, is documented [20]. This term refers to scenarios where a fraction of codebooks are rarely utilized to represent the input vectors. A recent study [21] offers a straightforward explanation for this collapse: an internal codebook covariate shift during the training phase. The frequent adjustments in encoder representations lead to misalignment with the VQ codebooks, resulting in only a subset of codebook parameters being updated. Consequently, VQ layers are prone to divergence, often ending up with a significant number of dead vectors. Various strategies have been proposed in generative modeling to address this, including the use of stochastic VQs [20, 22], reinitializing codewords with K-means centroids [23, 24], and bypassing the ramifications of straight-through estimation (STE) with a finite scalar quantizer [25]. In the context of audio compression, Descript’s codec attempts to mitigate codebook collapse by reducing vector dimensions [26] and normalizing them within a Euclidean ball [23]. However, their findings suggest that the primary factor in enhancing codebook efficiency is the incorporation of adversarial losses. To circumvent the challenges associated with optimizing GANs, we propose a learning paradigm that includes a pre-training stage to benefit codebook learning. Specifically, the main contributions of this paper are as follows:

- We propose a lightweight speech codec, ESC, which achieves efficient compression through the integration of cross-scale residual vector quantization (CS-RVQ) and mirrored non-linear hierarchical Swin Transformer layers.
- ESC attains double the compression ratio of the original TFNet-CSVQ described in [18], while maintaining comparable reconstruction quality to DAC, which is currently recognized as the state-of-the-art in high-fidelity audio codecs.

¹github.com/yzGuu830/efficient-speech-codec

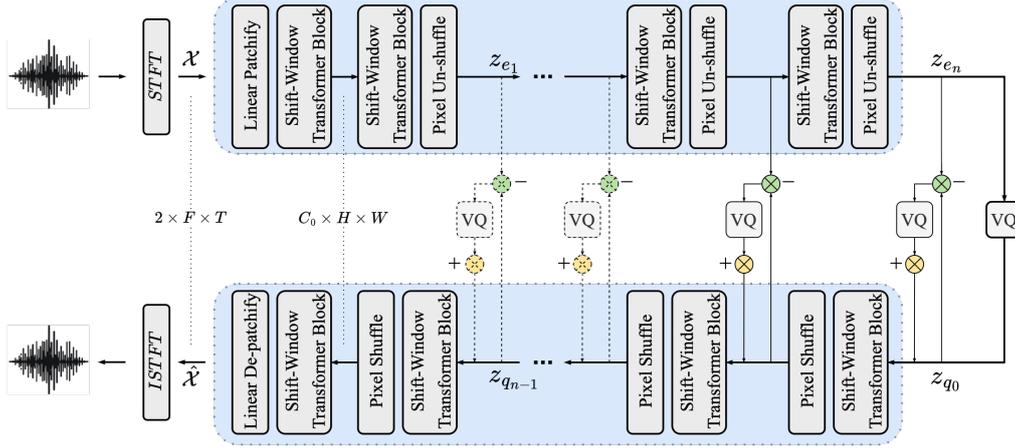


Figure 1: The framework of ESC: Input audio is transformed to a complex STFT \mathcal{X} and linearly embedded into patches. Encoder STBs iteratively halve the frequency resolution and produce hierarchical feature representations. Mirrored decoder STBs iteratively recover the frequency resolution by leveraging coarse-to-fine encoded features. The intermediate residuals between encoder and decoder hidden states are progressively quantized to refine decoding. The entire network is solely composed of efficient transformer blocks. The figure displays a scenario when the deepest 3 of $n + 1$ total bitstreams (solid lines) are transmitted, with others left inactive.

- We introduce a learning paradigm that incorporates a pre-training stage designed to mitigate codebook collapse and enhance codeword utilization. Empirical studies confirm the effectiveness of this approach.

2. Method

2.1. Architecture formulation

Demonstrated in Figure 1, our proposed ESC takes the complex spectrogram $\mathcal{X} \in \mathbb{R}^{2 \times F \times T}$ of audio from Short-Time Fourier Transform (STFT). The real and imaginary parts are treated as separate channels. ESC comprises mirrored encoder and decoder layers $E(\cdot; \phi_1, \dots, \phi_n)$, $D(\cdot; \psi_1, \dots, \psi_n)$, in addition to multi-scale vector quantizers Q_0, Q_1, \dots, Q_n . Each encoder and decoder layer performs downscaling and upscaling to create coarse and fine representations. We denote the hidden states of the t -th frame $\mathbf{x}_t \in \mathbb{R}^{2 \times F}$ after all encoder layers as

$$\mathbf{z}_{e_1}, \dots, \mathbf{z}_{e_n} = E(\mathbf{x}_t; \phi_1, \dots, \phi_n), \quad (1)$$

where each \mathbf{z}_e is flattened into a vector prior to quantization. Each VQ is parameterized with a codebook collection of codewords $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_K\}$, allowing for a $\log_2 K$ bit budget. It quantizes a vector $\mathbf{z}_e \in \mathbb{R}^d$ in the euclidean space with

$$Q(\mathbf{z}_e; \mathcal{C}) = \arg \min_{\mathbf{c}_j \in \mathcal{C}} \|\mathbf{z}_e - \mathbf{c}_j\|_2, \quad (2)$$

where the quantized latent of \mathbf{z}_e is the nearest code vector \mathbf{c}_k in \mathcal{C} . The index $k \in \mathbb{Z}$ is the discretized code. The decoder consists of mirrored blocks. Starting from the quantized latents at the bottom layer VQ, it progressively maps the latents back to the original spectrum. This process is further elaborated in Section 2.3. Finally, the reconstructed spectrum is converted back to the one-dimensional audio domain via inverse STFT.

2.2. Window transformer encoder and decoder

As convolutional layers often fall short in effectively capturing redundancies within audio signals, we propose replacing

them with Swin Transformer blocks (STBs) and their associated extended decoding blocks. As illustrated in Figure 1, the complex spectrogram \mathcal{X} is initially split into patches and linearly projected into the space $\mathbb{R}^{C_0 \times H \times W}$, with a patch size of $(\frac{T}{W}, \frac{F}{H})$ across frequency and temporal dimensions. As in [19], each STB consists of a cascade of two interleaved window-attention transformer layers, with the outputs of the prior layer being shifted. This configuration allows STBs to learn both local and global feature dependencies effectively. It also enables efficient computation of attention by partitioning the spectrum into smaller windows.

To maintain temporal resolution, we have adapted the patch merging layer subsequent to the encoder STB by integrating a pixel unshuffle module along only the frequency dimension. For example, the first downsampled resolution of a patch embedding might be $2C_0 \times \frac{H}{2} \times W$. This downsampled embedding is then subject to a linear transformation using a matrix $W_p \in \mathbb{R}^{2C_0 \times C_1}$, which enlarges the hidden size to facilitate the extraction of more abstract information at deeper layers. In the decoder, we employ a mirrored upscaling approach in the STBs. Similar to the encoder layers, each layer is preceded by a pixel shuffle that doubles the frequency resolution, followed by a mirrored linear transformation. At the end of the decoding process, the decoded patch embeddings are reshaped and linearly projected back to produce a recovered spectrum $\hat{\mathcal{X}}$.

2.3. Cross-scale residual vector quantization

For efficient learning of audio features, ESC utilizes multi-scale features through a coarse-to-fine decoding process to generate coded bitstreams at different levels. Laid on the cross-scale vector quantization (CSVQ) [18] paradigm, our codec incorporates a more intuitive cross-scale residual VQ (CS-RVQ) structure without the need for additional networks to fuse encoder and decoder features. As depicted in Algorithm 1 and Figure 1, the decoder is conditioned on encoded features from various scales of resolution, distinguishing it from the commonly used residual VQs (RVQ) that operate only at the lowest scale and tend to overlook high-level information [12]. At the lowest bitstream,

Algorithm 1 Cross-Scale Residual Vector Quantization

- 1: **Input:** multi-scale encoder hidden states $\mathbf{z}_{e_1}, \dots, \mathbf{z}_{e_n}$, vector quantizers Q_0, Q_1, \dots, Q_n , and mirrored decoder D
 - 2: **Initialization:** $\mathbf{z}_{q_0} \leftarrow Q_0(\mathbf{z}_{e_n})$
 - 3: **for** $i = 0, \dots, n - 1$ **do**
 - 4: $\mathbf{z}_{q_i} \leftarrow Q_{i+1}(\mathbf{z}_{e_{n-i}} - \mathbf{z}_{q_i}) + \mathbf{z}_{q_i}$
 - 5: $\mathbf{z}_{q_{i+1}} \leftarrow D_{i+1}(\mathbf{z}_{q_i}; \psi_{i+1})$
 - 6: **end for**
 - 7: **return** $\mathbf{z}_{q_n}, k_0, k_1, \dots, k_n \quad \triangleright k_i$ is discrete code from Q_i
-

\mathbf{z}_{e_n} is directly quantized by the bottleneck quantizer Q_0 . For higher bitstreams, the residual between each encoded feature and current quantized latent, $\mathbf{z}_{e_{n-i}} - \mathbf{z}_{q_i}$, is quantized by Q_{i+1} . This quantized residual is then added back to \mathbf{z}_{q_i} , and decoded by the following decoder layer function $D_{i+1}(\cdot; \psi_{i+1})$, producing the next quantized latent $\mathbf{z}_{q_{i+1}}$. Subsequently, the residuals at higher resolutions are progressively quantized, forming additional bitstreams conditioned on the preceding ones. This framework enables multi-scale learning and allows the decoder’s to incrementally mitigate quantization errors starting from the bottleneck Q_0 .

During inference, compressing with $s > 1$ bitstreams involves additional forward pass by a subset of the decoder layers, which produces coarse-to-fine residuals quantized by Q_0, Q_1, \dots, Q_{s-1} . When s is set to 1, the model simplifies to a fixed-scale VQ codec operating at the bottleneck. To reconstruct audio, the compressed discrete codes k_1, \dots, k_{s-1} are dequantized using their corresponding codewords. These are then progressively added back to recover the original input frame \mathbf{x}_t .

2.4. Vector quantization module

The proposed Cross-Scale Residual VQ (CS-RVQ) encodes representations at variable scales, which results in large latent dimensions and a propensity for codebook collapse. To effectively optimize the codebook embeddings, we adopt a combination of Product Quantization [27] and vector dimension reduction [26] at each bitstream. Concretely, the vector $\mathbf{z} \in \mathbb{R}^d$ to be quantized is split into a set of l sub-vectors. Each sub-vector is quantized independently by a VQ and subsequently concatenated as follows:

$$\mathbf{z} \equiv \{\mathbf{z}_m \mid \mathbf{z}_m \in \mathbb{R}^{d/l}, m = 1, \dots, l\},$$
$$Q(\mathbf{z}_m; \mathcal{C}_m) = W_{out}^\top \arg \min_{\mathbf{c}_j \in \mathcal{C}_m} \|W_{in}^\top \mathbf{z}_m - \mathbf{c}_j\|_2, \quad (3)$$

$$\text{where } W_{in} \in \mathbb{R}^{\frac{d}{l} \times u}, W_{out} \in \mathbb{R}^{u \times \frac{d}{l}}, u \ll d/l.$$

Following [11], the projected vector $W_{in}\mathbf{z}_j$ and codebook \mathcal{C} are both L_2 normalized before searching for the nearest codewords. This equalizes the scales of encoded vectors and codewords, which improves codebook learning because a larger subset of codewords can receive gradients updates.

2.5. Training

As transformer layers are intrinsically difficult to converge, joint training with VQs often results in sub-optimal parameters. Therefore, we propose a learning paradigm that includes a warm-start to enhance the learning in ESC, as detailed in Algorithm 2. Initially, all VQ layers are deactivated, so no quantization occurs. Only the encoder and decoder are trained using the same CS-RVQ framework, allowing encoded representations to bypass the quantizers and flow directly into the decoder layers,

Algorithm 2 Efficient Optimization of Vector Quantizers

- 1: **repeat**
 - 2: $\hat{\mathbf{x}}_t = D(E(\mathbf{x}_t; \phi_1, \dots, \phi_n); \psi_1, \dots, \psi_n)$
 - 3: $\mathcal{L} = \mathcal{L}_{recon}(\mathbf{x}_t, \hat{\mathbf{x}}_t)$
 - 4: Take gradient descent step on $\nabla_E \mathcal{L}, \nabla_D \mathcal{L}$
 - 5: **until** converged
 - 6: activate all codebooks \mathcal{C} and continue learning as usual
-

in what we term a “pre-training” stage. After encoder and decoder have converged by minimizing reconstruction losses, we proceed to train the entire vector-quantized network as usual. This helps alleviate the distribution shift of encoded vectors as the encoder is optimized in advance. Through a pre-training stage, we aim to stabilize subsequent codebook learning and improve codebook usage. Moreover, pre-training encoders and decoders is simpler, as it avoids the quantization errors typically associated with VQs.

To enable bitrate scalability, we sample the number of transmitted bitstreams $s \sim \text{Uniform}\{1, \dots, n\}$ at a rate p during training. p is a hyperparameter that balances the reconstruction quality at different bitrates, as proposed by [11]. The reconstruction loss \mathcal{L}_{recon} consists of an L_2 distance on \mathcal{X} and a multi-scale mel-spectrogram loss same as that in [11], denoted by \mathcal{L}_{stft} and \mathcal{L}_{mel} . To optimize cross-scale codebooks, we use the standard combination of commitment and codebook loss \mathcal{L}_{ct} and \mathcal{L}_{cb} with straight through estimators in [7], each is averaged across l groups and summed across transmitted s bitstreams to form \mathcal{L}_{vq} . Exponential moving average (EMA) is not adopted as [11] points out that it fails to mitigate codebook collapse. The overall training objective is the summation of \mathcal{L}_{recon} and \mathcal{L}_{vq} , as follows:

$$\mathcal{L}_{recon} = \lambda_1 \mathcal{L}_{mel} + \lambda_2 \mathcal{L}_{stft}, \quad (4)$$

$$\mathcal{L}_{vq} = \lambda_3 \mathcal{L}_{cb} + \lambda_4 \mathcal{L}_{ct}. \quad (5)$$

3. Experiments

3.1. Experimental setup

We extract 150 hours 3-second multilingual speech audio clips from DNS Challenge dataset [28] for training. For evaluation, we use 1158 10-second speech clips from the LibriSpeech [29], Multilingual LibriSpeech [30], and AIShell [31] datasets. We use AdamW optimizer with a learning rate of 0.0001. Final models are trained for 400k steps with a batch size of 36, where the pre-training stage comprises 75k steps. After pre-training, we decay the learning rate by 0.999996 at each step and set p to 0.75. The weighted parameters $\lambda_1, \lambda_2, \lambda_3, \lambda_4$ are configured as 0.25, 1.0, 1.0 and 0.25 respectively.

3.2. Performance evaluation on reconstruction quality

We experimented two versions of our codec:

1. **ESC:** Our final codec has a patch size of (3, 2) and encodes two overlapped STFT frames together. The STFT window and hop lengths are set at 20 ms and 5 ms. Each bitstream is composed of three 10-bit codebooks in product quantization, yielding a bitrate of 1.5kbps. The ESC codec is scalable up to a maximum of 6 bitstreams from 6 STBs in both encoder and decoder layers. It uses GELU activation and LayerNorm in the attention modules with a window size of 4. The hierarchical dimensions of our STBs range from $C_0 = 72$ to $C_5 = 384$, with a total of 8.4M parameters.

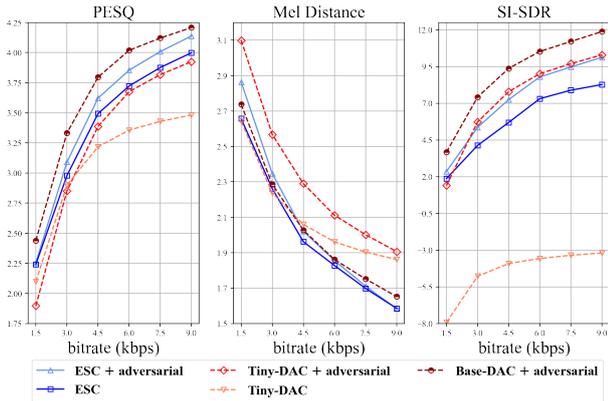


Figure 2: Evaluation results of codecs on our test dataset: dashed lines represent DAC models and solid lines represent ESC models, with x -axis being transmission bit per second and y -axis being PESQ (↑), Mel-Distance (↓) and SI-SDR (↑).

2. **ESC + adversarial**: This variant integrates adversarial losses, using the same multi-scale waveform and spectrogram discriminator as described in [11]. In this setup, \mathcal{L}_{stft} is replaced with the L_1 feature matching loss [14], following the approach in [11]. The discriminator has ~ 45 M parameters.

For comparison, we selected Descript’s audio codec as a baseline and reproduced it in three variants on our training set:

1. **Base-DAC + adversarial**: This is Descript’s original codec. We expanded the number of codebooks in its RVQs from 12 to 18 to match the bitrate levels of ESC. It is trained with the released configuration and has 74M parameters.
2. **Tiny-DAC + adversarial**: To ensure a fair comparison with ESC, we produced a smaller version of the DAC by reducing its decoder dimension from 1536 to 288. This version has ~ 10 M parameters, comparable in scale to ESC. All other configurations are maintained as in the original one.
3. **Tiny-DAC**: We also reproduced a smaller and non-adversarial version of Descript’s codec to assess the impact of GAN discriminators on improving audio fidelity.

We use the PESQ [32], the L_1 distance between log mel spectrograms of reference and decoded waveforms (mel distance) [11], and scale-invariant source-to-distortion ratio (SI-SDR) [33] as objective metrics to comprehensively assess the audio reconstruction quality. As shown in Figure 2, although ESC does not outperform Base-DAC trained with the same discriminator among all metrics, it consistently achieves superior reconstructions across all bandwidths compared to Tiny-DAC. This superiority is evident in experiments conducted both with and without GANs (i.e., ESC & Tiny-DAC versus ESC+adversarial & Tiny-DAC+adversarial). Furthermore, the DAC models exhibit a significant dependency on the discriminator, as evidenced by a substantial quality drop when Tiny-DAC is trained without adversarial losses. These results demonstrate that transformers combined with cross-scale residual VQs are more effective than CNNs paired with single-scale residual VQs in current neural audio codecs.

3.3. Performance evaluation on complexity

As noted earlier, Descript’s top-performing codec is hindered by computational bottlenecks. We detail the complexity results in Table 1. Latency experiments were conducted using an Intel

Table 1: A comparison of codec complexity in terms of CPU/GPU inference speed and model sizes.

Model	Params.(M)	Memory (MB)	CPU (s)		GPU (s)	
			Enc.	Dec.	Enc.	Dec.
ESC	8.4	30.5	0.78	0.59	0.10	0.06
Tiny-DAC	10	38.2	0.63	0.69	0.08	0.03
Base-DAC	74	282	1.42	6.02	0.08	0.03

Xeon E5-2660 CPU @ 2.60GHz and an NVIDIA RTX TITAN GPU. The encoding and decoding times were averaged across fifty 10-second, 16kHz speech signals at a bitrate of 9kbps. Notably, the ESC codec is smaller and faster on CPUs compared to Base-DAC, with $\times 9$ smaller model size, $\times 2$ encoding speed, and $\times 11$ decoding speed. While the GPU performance of ESC does not outperform DAC, the difference is minimal, with only an additional 0.03 seconds required for a 10-second input audio. Additionally, the encoding speed of ESC can be improved at lower bitrates—a capability not present in DAC.

3.4. Ablation study on a pre-training stage

We evaluate the efficacy of the proposed pre-training stage by monitoring both the bitrate utilization ratio and PESQ scores on a held-out validation set throughout the training process. The utilization rate is calculated as the sum of entropy (in bits) divided by the maximum number of bits from all transmitted bitstreams. We conducted an ablation experiment **ESC-scratch**, which involves training from scratch without freezing the codebooks. We compare its learning trajectory to that of **ESC-pretrain** in a non-adversarial setting. The results, depicted in Figure 3, show a significant incremental gap between the two trajectories. This gap confirms that a pre-training stage can enhance codebook learning (as evidenced by utilization rates approaching 1.0) and simultaneously improve audio quality.

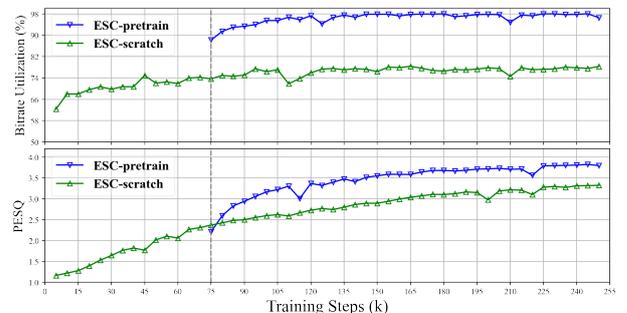


Figure 3: Learning trajectories of ESC model with and without a pre-training stage in the first 250k steps: vertical dashed gray lines denote the cutoff point when the pre-training stage ends.

4. Conclusions

In this paper, we propose an Efficient Speech Codec (ESC) surpassing existing baselines in coding efficiency for multi-lingual speech compression. Extensive evaluations demonstrate that our vector-quantized codec, which utilizes cross-scale learned transformers, outperforms traditional single-scale convolutional models in both reconstruction quality and complexity. As a future work, we anticipate that our codec can be effectively scaled up to accommodate larger, universal audio datasets.

5. References

- [1] J.-M. Valin, K. Vos, and T. Terriberry, "Definition of the opus audio codec," Tech. Rep., 2012.
- [2] M. Dietz, M. Multus, V. Eksler, V. Malenovsky, E. Norvell, H. Pobloth, L. Miao, Z. Wang, L. Laaksonen, A. Vasilache *et al.*, "Overview of the evs codec architecture," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 5698–5702.
- [3] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," *arXiv preprint arXiv:1609.03499*, 2016.
- [4] S. Mehri, K. Kumar, I. Gulrajani, R. Kumar, S. Jain, J. Sotelo, A. Courville, and Y. Bengio, "SampleRNN: An unconditional end-to-end neural audio generation model," in *International Conference on Learning Representations*, 2017. [Online]. Available: <https://openreview.net/forum?id=SkxKPDv5xl>
- [5] W. B. Kleijn, F. S. Lim, A. Luebs, J. Skoglund, F. Stimberg, Q. Wang, and T. C. Walters, "Wavenet based low rate speech coding," in *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2018, pp. 676–680.
- [6] J. Klejsa, P. Hedelin, C. Zhou, R. Fejgin, and L. Villemoes, "High-quality speech coding with sample rnn," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 7155–7159.
- [7] A. Van Den Oord, O. Vinyals *et al.*, "Neural discrete representation learning," *Advances in neural information processing systems*, vol. 30, 2017.
- [8] Y. Bengio, N. Léonard, and A. Courville, "Estimating or propagating gradients through stochastic neurons for conditional computation," *arXiv preprint arXiv:1308.3432*, 2013.
- [9] 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*, vol. 30, pp. 495–507, 2021.
- [10] A. Défossez, J. Copet, G. Synnaeve, and Y. Adi, "High fidelity neural audio compression," *Transactions on Machine Learning Research*, 2023, featured Certification, Reproducibility Certification. [Online]. Available: <https://openreview.net/forum?id=ivCd8z8zR2>
- [11] R. Kumar, P. Seetharaman, A. Luebs, I. Kumar, and K. Kumar, "High-fidelity audio compression with improved RVQGAN," in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. [Online]. Available: <https://openreview.net/forum?id=qjnl1QUfA>
- [12] A. Vasuki and P. Vanathi, "A review of vector quantization techniques," *IEEE Potentials*, vol. 25, no. 4, pp. 39–47, 2006.
- [13] P. Esser, R. Rombach, and B. Ommer, "Taming transformers for high-resolution image synthesis," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 12 873–12 883.
- [14] K. Kumar, R. Kumar, T. De Boissiere, L. Gestin, W. Z. Teoh, J. Sotelo, A. De Brebisson, Y. Bengio, and A. C. Courville, "Melgan: Generative adversarial networks for conditional waveform synthesis," *Advances in neural information processing systems*, vol. 32, 2019.
- [15] J. Kong, J. Kim, and J. Bae, "Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis," *Advances in neural information processing systems*, vol. 33, pp. 17 022–17 033, 2020.
- [16] S. gil Lee, W. Ping, B. Ginsburg, B. Catanzaro, and S. Yoon, "BigVGAN: A universal neural vocoder with large-scale training," in *The Eleventh International Conference on Learning Representations*, 2023. [Online]. Available: <https://openreview.net/forum?id=iTtGCMDEzS>
- [17] L. Ziyin, T. Hartwig, and M. Ueda, "Neural networks fail to learn periodic functions and how to fix it," *Advances in Neural Information Processing Systems*, vol. 33, pp. 1583–1594, 2020.
- [18] X. Jiang, X. Peng, H. Xue, Y. Zhang, and Y. Lu, "Cross-Scale Vector Quantization for Scalable Neural Speech Coding," in *Proc. Interspeech 2022*, 2022, pp. 4222–4226.
- [19] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [20] Y. Takida, T. Shibuya, W. Liao, C.-H. Lai, J. Ohmura, T. Uesaka, N. Murata, S. Takahashi, T. Kumakura, and Y. Mitsufuji, "SQ-VAE: Variational bayes on discrete representation with self-annealed stochastic quantization," in *International Conference on Machine Learning*, 2022.
- [21] M. Huh, B. Cheung, P. Agrawal, and P. Isola, "Straightening out the straight-through estimator: Overcoming optimization challenges in vector quantized networks," in *International Conference on Machine Learning*. PMLR, 2023, pp. 14 096–14 113.
- [22] J. Zhang, F. Zhan, C. Theobalt, and S. Lu, "Regularized vector quantization for tokenized image synthesis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 18 467–18 476.
- [23] A. Łańcucki, J. Chorowski, G. Sanchez, R. Marxer, N. Chen, H. J. Doling, S. Khurana, T. Alummäe, and A. Laurent, "Robust training of vector quantized bottleneck models," in *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2020, pp. 1–7.
- [24] P. Dhariwal, H. Jun, C. Payne, J. W. Kim, A. Radford, and I. Sutskever, "Jukebox: A generative model for music," *arXiv preprint arXiv:2005.00341*, 2020.
- [25] F. Mentzer, D. Minnen, E. Agustsson, and M. Tschannen, "Finite scalar quantization: VQ-VAE made simple," in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=8ishA3LxN8>
- [26] J. Yu, X. Li, J. Y. Koh, H. Zhang, R. Pang, J. Qin, A. Ku, Y. Xu, J. Baldridge, and Y. Wu, "Vector-quantized image modeling with improved VQGAN," in *International Conference on Learning Representations*, 2022. [Online]. Available: <https://openreview.net/forum?id=pfNyExj7z2>
- [27] A. Baevski, S. Schneider, and M. Auli, "vq-wav2vec: Self-supervised learning of discrete speech representations," *arXiv preprint arXiv:1910.05453*, 2019.
- [28] C. K. Reddy, H. Dubey, V. Gopal, R. Cutler, S. Braun, H. Gamper, 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.
- [29] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.
- [30] V. Pratap, Q. Xu, A. Sriram, G. Synnaeve, and R. Collobert, "Mls: A large-scale multilingual dataset for speech research," *arXiv preprint arXiv:2012.03411*, 2020.
- [31] Y. Shi, H. Bu, X. Xu, S. Zhang, and M. Li, "Aishell-3: A multi-speaker mandarin tts corpus and the baselines," *arXiv preprint arXiv:2010.11567*, 2020.
- [32] I. Union, "Wideband extension to recommendation p. 862 for the assessment of wideband telephone networks and speech codecs," *International Telecommunication Union, Recommendation P*, vol. 862, 2007.
- [33] J. Le Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, "Sdr-half-baked or well done?" in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 626–630.