# MuCodec: Ultra Low-Bitrate Music Codec

Yaoxun Xu[1,*], Hangting Chen[2,†], Jianwei Yu[†], Wei Tan[2], Rongzhi Gu[2], Shun Lei[1], Zhiwei Lin[1], Zhiyong Wu[1,3,†]

[1] Shenzhen International Graduate School, Tsinghua University, Shenzhen, China
[2] Tencent AI Lab [3] The Chinese University of Hong Kong, Hong Kong SAR, China
xuyx22@mails.tsinghua.edu.cn erichtchen@tencent.com tomasyu@foxmail.com zywu@sz.tsinghua.edu.cn

*Abstract*—Music codecs are a vital aspect of audio codec research, and ultra low-bitrate compression holds significant importance for music transmission and generation. Due to the complexity of music backgrounds and the richness of vocals, solely relying on modeling semantic or acoustic information cannot effectively reconstruct music with both vocals and backgrounds. To address this issue, we propose MuCodec, specifically targeting music compression and reconstruction tasks at ultra low bitrates. MuCodec employs MuEncoder to extract both acoustic and semantic features, discretizes them with RVQ, and obtains Mel-VAE features via flow-matching. The music is then reconstructed using a pre-trained MEL-VAE decoder and HiFi-GAN. MuCodec can reconstruct high-fidelity music at ultra low (0.35kbps) or high bitrates (1.35kbps), achieving the best results to date in both subjective and objective metrics. Code and Demo: https://xuyaoxun.github.io/MuCodec_demo/.

*Index Terms*—Music, Codec, Flow-Matching, Low Bitrate

## I. INTRODUCTION

Music codecs [1]–[3] are a crucial component in the field of audio codec [4]–[6] research. The significance of ultra low-bitrate compression lies in its potential applications, such as music transmission, where the bitrate of MP3 [7] is considerably high, and music generation [8]–[10], where short sequences are highly effective for language model construction. Furthermore, considering the diversity of background, sound events, and vocals in music, achieving high-fidelity reconstruction at ultra low bitrates would signify a substantial advancement in the field of universal audio generation.

Recent music compression techniques based on neural codecs [11]–[15] attempt to compress music directly into discrete tokens. While discrete representations often yield higher compression densities, they inherently suffer from substantial information loss. To reconstruct a more accurate approximation of the original features from discrete tokens, a more robust representation and a stronger decoder are necessary. Common codecs like Encodec [16] and Generative Adversarial Networks(GAN)-based methods [17]–[19] exhibit limitations in achieving particularly low bitrates.

In recent years, some research and works have focused on using semantic modeling to represent musical characteristics and utilizing diffusion [20] for reconstruction, such as SemantiCodec [21] and SEED-TTS [22]. However, these models are

not specifically designed for music-related tasks. Compared to speech tasks, music has a rich background, including instruments like piano and bass, and vocals that should be clearly discernible from the background music. Therefore, it is essential to consider both semantic and acoustic information; focusing solely on one aspect would compromise the overall perceptual quality of the reconstructed audio.

To address these challenges, we propose a flow-matching-based [23] music codec MuCodec. MuCodec uses a specialized feature extractor, MuEncoder, based on the two key aspects of music: vocals and background. The MuEncoder features are then discretized using RVQ and employed as conditions for reconstructing Mel-VAE features via flow-matching. We reconstruct the Mel spectrogram by passing the Mel-VAE features through a pre-trained Mel-VAE decoder [24], and ultimately, the music is reconstructed using HiFi-GAN [25]. Our contributions can be summarized as follows:

- We propose MuCodec, which achieves the lowest bitrate to date while maintaining the highest-quality music reconstruction capabilities.
- MuCodec employs MuEncoder as the feature extractor and Diffusion Transformer (DiT) [26] along with flow-matching-based method for fine-grained music modeling.
- Both subjective and objective experiments demonstrate that MuCodec achieves the best performance to date in music reconstruction tasks at both low and high bitrates.

## II. METHOD

As illustrated in Fig. 1, MuCodec comprises MuEncoder, RVQ, a reconstruction model using flow-matching, Mel-VAE decoder, and HiFi-GAN. MuEncoder is a music extractor, primarily responsible for extracting both acoustic and semantic representations that better capture the characteristics of music. RVQ compresses the representations obtained from MuEncoder. The objective of flow-matching is to reconstruct low-bitrate discrete representations to obtain Mel-VAE features. Subsequently, the pretrained Mel-VAE decoder restores these features into a Mel spectrogram. Finally, the reconstructed music is obtained through a pretrained HiFi-GAN.

### A. MuEncoder

Music reconstruction is more complex than speech or audio events, as it requires modeling both acoustic background and vocals. We design MuEncoder, composed of 13 stacked
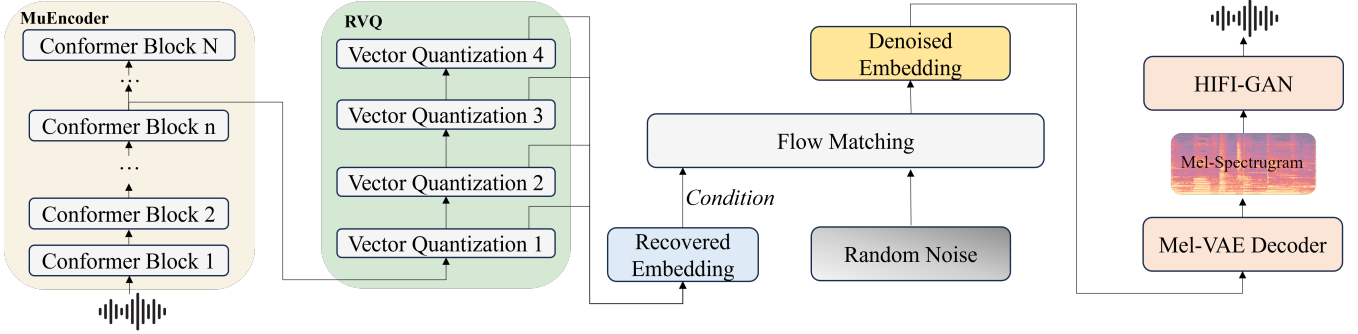
Fig. 1. Framework of the proposed MuCodec.

Conformer blocks, to extract acoustic and semantic features of background music and vocals.

To enable MuEncoder to extract both acoustic and semantic features, we implement a two-stage training process. In the first stage, we use the Mask Language Model constraint [27] to learn to predict masked regions based on unmasked speech signals, allowing MuEncoder to perceive contextual information and enhance representational capabilities. In the second stage, we introduce two constraints: reconstruction and lyrics recognition constraints. Reconstruction constraint aims to make extracted features closer to acoustic features, with two targets: restoring Mel spectrograms and predicting Constant-Q Transform (CQT) [28] features. Lyrics recognition constraint ensures extracted features contain semantic information. These constraints enhance MuEncoder's feature extraction compatibility from both background music and vocal perspectives.

### B. Residual Vector Quantization

In MuCodec, we opt to use Residual Vector Quantization (RVQ) to discretize the MuEncoder features for its ability to compress representations through the residual process and provide more refined approximations using cascaded codebooks.

### C. Flow-Matching

MuCodec employs a flow-matching-based method for reconstruction, as it offers more stable training compared to GAN-based method and requires fewer training steps to achieve better results in ultra low-bitrate reconstruction task. Specifically, we use the discretized MuEncoder representations as a condition and perform finer-grained reconstruction through flow-matching with a Diffusion Transformer.

Instead of choosing the music or its Mel spectrogram as the flow-matching target due to their abundant and complex information, we predict the more manageable and information-rich Mel-VAE features for reconstruction. A pretrained Mel-VAE decoder serves as our Mel spectrogram generator, while a pre-trained HiFi-GAN functions as the music generator.

### D. Discussion

*1) Disentangle:* In music reconstruction tasks, the two most important evaluation aspects are vocals and music background. To better verify the benefits of simultaneously focusing on these two features in music reconstruction tasks, we design

comparative experiments to model these two aspects separately. Specifically, we choose pre-trained HuBERT [30] and MERT [31] models to separately model vocals and music background. HuBERT typically contains richer semantic information, while MERT focuses more on acoustic features.

*2) Scalability:* Although MuCodec is initially designed for music reconstruction tasks, it can also be easily applied to other types of audio without incorporating any additional training data, such as speech or acoustic events. MuCodec employs two constraints, one to enhance the background modeling of the audio itself and the other to strengthen the semantic modeling of vocals. As a result, MuCodec exhibits good performance in scenarios with pure vocals, pure background, or both vocals and background simultaneously. Our demo webpage exhibits the reconstruction results of different audio types and presents some other experimental outcomes.

### III. EXPERIMENTAL SETUP

To train MuCodec, we utilize a large-scale internal music dataset of Chinese and English songs with a minimum 32kHz sampling rate. We segment the music into fixed 35.84-second lengths during training. For fairness, the test set comprises randomly selected 250 Chinese and 250 English song segments, each 20-30 seconds long with corresponding lyrics.

For the GAN-based method, we use a fully convolutional architecture following Descript Audio Codec(DAC) [29] encoder and decoder, changing the quantizer to RVQ. We match its model size to MuCodec. To further analyze, we also experiment with replacing its encoder directly with MuEncoder.

Considering GANs' weak reconstruction capabilities in low-bitrate scenarios, GAN-based method experiments are trained for 120k steps. In other cases, unless specifically stated, all test models train for 20k steps to demonstrate our approach's effectiveness within reasonable training time. Regarding SemantiCodec, we select two settings with bitrates similar to the high and low bitrates used in our experiments.

To better evaluate the performance of reconstructed music, we adopt both subjective and objective assessments. In subjective evaluations, we randomly select 5 Chinese and 5 English song clips as the test set and invite 10 professional participants to conduct a MUSHRA-inspired [32] listening test. In objective evaluations, we choose two types of metrics corresponding to the two aspects of music: background and

| | Method | CodeBookSize | Token Rate | kbps | VISQOL ↑ | SPK_SIM ↑ | WER (%) ↓ |
|---|---|---|---|---|---|---|---|
| Origin music | — | — | — | — | — | — | 10.92 |
| Low-Bitrate Scenario | DAC [29]+GAN | 1 x 16384 | 25 | 0.35 | 2.94/2.93 | 0.39 | 131.80 |
| | MuEncoder+GAN | 1 x 16384 | 25 | 0.35 | 2.60/2.63 | 0.35 | 97.86 |
| | MuEncoder+Diffusion | 1 x 16384 | 25 | 0.35 | 2.97/2.96 | 0.58 | 89.31 |
| | SemantiCodec [21] | 1 x 32768 | 25 | 0.375 | 1.92/1.92 | 0.52 | 120.17 |
| | **MuEncoder+Flow-Matching (MuCodec)** | 1 x 16384 | 25 | 0.35 | 3.09/3.08 | 0.63 | 68.37 |
| | **MuCodec (200k)** | 1 x 16384 | 25 | 0.35 | **3.19/3.20** | **0.75** | **40.81** |
| High-Bitrate Scenario | DAC [29]+GAN | 4 x 10000 | 25 | 1.33 | 3.00/2.99 | 0.38 | 137.51 |
| | MuEncoder+GAN | 4 x 10000 | 25 | 1.33 | 2.62/2.61 | 0.34 | 62.59 |
| | MuEncoder+Diffusion | 4 x 10000 | 25 | 1.33 | 3.34/3.34 | 0.75 | 43.36 |
| | SemantiCodec [21] | 1 x 16384 | 100 | 1.40 | 1.96/1.96 | 0.68 | 55.17 |
| | **MuEncoder+Flow-Matching (MuCodec)** | 4 x 10000 | 25 | 1.33 | 3.30/3.30 | 0.80 | 34.19 |
| | **MuCodec (200k)** | 4 x 10000 | 25 | 1.33 | **3.43/3.42** | **0.86** | **26.12** |

vocals. We use ViSQOL [33] as an audio quality assessment metric. Since the background can interfere with vocal evaluation, we separate the vocals and background of the generated music using a pre-trained separation model [34]. We then calculate the similarity between the generated vocal part and the original vocal part with a pre-trained speaker similarity model [35] and use Whisper-v2 [36] to compute the Word Error Rate (WER) of the generated vocal part as the vocal clarity evaluation metric.

In the MuEncoder setting, we employ a 13-layer Conformer [37] model. We assign a weight of 1 to the music reconstruction loss and 0.2 to the lyrics recognition loss, which consists of both CTC Loss and RNN-T Loss.

In the RVQ setting, we design two configurations for high- and low-bitrate scenarios. For low-bitrate scenarios, we employ a single codebook with a size of 16,384 and a bitrate of 0.35kbps. Conversely, for the relatively high-bitrate scenarios, we use four codebooks, each with a size of 10,000, and achieve a bitrate of 1.33kbps.

In the flow-matching setting, we employ a 24-layer Transformer2d model [38] for reconstruction, featuring an attention head dimension of 72, a norm epsilon of 1e-06, and 32 norm groups. We use ada norm single as the norm type and set the number of ada norm embeds to 1000. During the generation process, we utilize sampling via classifier-free guidance [39], specifically setting the guidance scale value to 1.5.

During inference, we choose a denoising step size of 50 for flow-matching to balance reconstruction quality and computational efficiency. We use a pre-trained open-source Mel-VAE decoder and HiFi-GAN for both training and inference. Our experiments run on 8 40G-A100 GPUs with a batch size of 4.

## IV. RESULT

### A. Main Comparation

In this experiment, we offer a thorough comparison of various prevalent reconstruction methods. MuCodec undergoes an in-depth comparative analysis from both objective and subjective assessments, with objective results in TABLE I.

First, it can be observed that DAC+GAN and MuEncoder+GAN method underperform in low-bitrate music reconstruction tasks, despite 120k training steps, which exceed other tasks.

Second, a difference between MuEncoder+Diffusion and MuCodec in low-bitrate music reconstruction tasks can be noticed. While MuEncoder+Diffusion outperforms GAN and MuEncoder+GAN, it falls short compared to MuCodec. This is because MuCodec employs the flow-matching method, which more directly and effectively models the noise-to-target distribution path compared to diffusion methods, achieving better results with fewer reconstruction steps.

Lastly, in the low-bitrate (0.35kbps) scenario, SemantiCodec's performance is subpar. Despite its state-of-the-art performance in acoustic event reconstruction, it lacks a dedicated design for music reconstruction tasks. Hence, its performance significantly decreases when handling more complex music reconstruction tasks compared to MuCodec. Furthermore, SemantiCodec only supports single-channel audio reconstruction at a 16k sampling rate, while MuCodec supports dual-channel audio at a 48k sampling rate, providing a greater advantage in music reconstruction.

At a higher bitrate (1.33kbps), MuCodec's performance continues to surpass other methods, showing the same trend as in the low-bitrate scenario. This demonstrates that MuCodec not only excels in low-bitrate scenarios but also delivers desirable results in high-bitrate music reconstruction tasks.

Moreover, we can observe from the table that when the training steps of MuCodec are increased to 200k, its performance improves further. However, a training step size of 20k already achieves a considerable level, highlighting MuCodec's robust compression and reconstruction capabilities.

Regarding the subjective results in Fig. 2, it is observed that the DAC+GAN method falls short in terms of audio quality at both low and high bitrates, indicating limited fine-grained modeling capability. In contrast, SemantiCodec shows a noticeable improvement over DAC+GAN method and performs better at high bitrates than low bitrates. However, despite its superior performance in acoustic event and speech re-
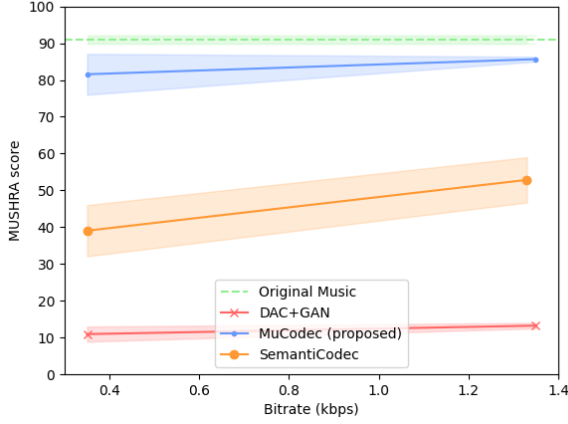
Fig. 2. Listening test results on performance analysis.

construction, the music reconstruction remains unsatisfactory, reflecting the challenges of music reconstruction tasks.

In comparison, our proposed MuCodec achieves excellent reconstruction results at both low and high bitrates, significantly outperforming SemantiCodec and DAC+GAN methods and closely resembling the original music. Moreover, the small difference between low and high bitrate MUSHRA scores suggests that MuCodec already attains a highly desirable reconstruction effect at extremely low bitrates.

### B. Impact of Different MuEncoder Training Losses

In this experiment, we evaluate the impact of MuEncoder on MuCodec under different loss conditions. Experiment #1 uses only the Mask Language Model loss (MLM Loss). Experiment #2 adds reconstruction loss (Recons Loss) to #1, including Mel spectrogram and CQT feature reconstruction loss. Experiment #3 incorporates lyrics recognition loss (ASR loss) based on #2, with specific results shown in TABLE II.

TABLE II
EXPERIMENTAL RESULTS ON DIFFERENT MUENCODER TRAINING LOSSES

| ID | MuEncoder Loss | VISQOL ↑ | SPK_SIM ↑ | WER (%) ↓ |
|---|---|---|---|---|
| #1 | MLM Loss | 2.70/2.71 | 0.587 | 84.41 |
| #2 | #1+Recons Loss | 2.83/2.86 | 0.591 | 87.98 |
| #3 | #2+ASR Loss | **3.09/3.08** | **0.631** | **68.37** |

The results show that compared to #1, ViSQOL and speaker similarity indicators improve in #2 due to the additional reconstruction loss, while WER slightly decreases. This suggests that reconstruction loss enhances audio quality but has limited impact on vocal modeling. Comparing #3 to #2 reveals a significant WER reduction after adding recognition loss, benefiting vocal modeling and providing some support to ViSQOL and speaker similarity. This highlights that introducing reconstruction and recognition losses during training improves MuCodec's performance in music reconstruction tasks.

### C. Influence of MuEncoder Layer Selection

In this experiment, we evaluate MuCodec's performance under different MuEncoder layer conditions, specifically the 3rd, 7th, and 11th layers, with results in TABLE III.

TABLE III
EXPERIMENTAL RESULTS ON THE MUENCODER LAYER SELECTION

| MuEncoder Layer | VISQOL ↑ | SPK_SIM ↑ | WER (%) ↓ |
|---|---|---|---|
| 3 | **3.13/3.14** | **0.656** | 76.65 |
| 7 | 3.09/3.08 | 0.631 | 68.37 |
| 11 | 2.92/2.92 | 0.618 | **63.46** |

The results indicate that music reconstructed with lower MuEncoder layer features has better ViSQOL and speaker similarity indicators. As the number of MuEncoder layers increases, the reconstructed music quality decreases, while vocal clarity improves. This suggests that lower MuEncoder layers have stronger acoustic characteristics aiding in background music reconstruction, and higher layers contain more semantic features supporting vocal reconstruction. Therefore, in practice, the choice of MuEncoder layers needs to balance specific requirements, leading us to select the 7th layer as a balanced option in our experiments.

### D. Validation Experiment for Disentangling Acoustic and Semantic Features

In this experiment, we analyze the comparison between separate modeling and MuEncoder. We select the high-bitrate scenario in disentangle experiments to match MuCodec's high-bitrate case (1.33kbps). Separate HuBERT and MERT experiments use 4 codebooks, each with a size of 10,000, while joint modeling experiments with HuBERT and MERT allocate 2 codebooks for each model, each containing 10,000 elements. Specific results are detailed in TABLE IV.

TABLE IV
EXPERIMENTAL RESULTS ON THE FEATURE EXTRACTOR

| Feature Extractor | Codebook | VISQOL ↑ | SPK_SIM ↑ | WER (%) ↓ |
|---|---|---|---|---|
| HuBERT | 4 | 2.54/2.54 | 0.312 | 80.78 |
| MERT | 4 | 2.82/2.83 | 0.463 | 136.78 |
| HuBERT+MERT | 2+2 | 3.12/3.14 | 0.710 | 58.06 |
| MuEncoder | 4 | **3.30/3.30** | **0.804** | **34.19** |

It can be found that using HuBERT alone results in relatively low ViSQOL and speaker similarity, suggesting its inability to effectively model rich backgrounds. Using MERT alone improves audio quality and background but slightly decreases vocal clarity. Jointly modeling HuBERT and MERT features improves both background and vocal clarity without increasing bitrate. This suggests that jointly modeling vocals and background positively impacts overall music reconstruction but introduces additional computational complexity.

In contrast, using only MuEncoder yields better music reconstruction results than separate HuBERT+MERT and requires modeling only one type of feature. This makes it more suitable for modeling, prediction, and music generation tasks.

### V. CONCLUSION

To better address the challenge of ultra low-bitrate music reconstruction, we propose MuCodec, which achieves the lowest bitrate to date while maintaining excellent reconstruction music quality. MuCodec employs the MuEncoder feature

extractor that considers both acoustic and semantic features of music, then the features are discretized using RVQ and finely reconstructed to Mel-VAE features via a flow-matching approach. The music is then reconstructed through a pretrained Mel-VAE decoder and HiFi-GAN. In both subjective and objective experiments, MuCodec significantly surpasses the current best results, realizing high-quality music reconstruction at an ultra low-bitrate scenario.

## REFERENCES

[1] Stuart Cunningham and Iain McGregor, "Subjective evaluation of music compressed with the acer codec compared to aac, mp3, and uncompressed pcm," *International Journal of Digital Multimedia Broadcasting*, vol. 2019, no. 1, pp. 8265301, 2019.

[2] Jean-Marc Valin, Gregory Maxwell, Timothy B Terriberry, and Koen Vos, "High-quality, low-delay music coding in the opus codec," *arXiv preprint arXiv:1602.04845*, 2016.

[3] Martin Dietz, Markus Multrus, Vaclav Eksler, Vladimir Malenovsky, Erik Norvell, Harald Pobloth, Lei Miao, Zhe Wang, Lasse Laaksonen, Adriana 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.

[4] Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi, "Soundstream: An end-to-end neural audio codec," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 495–507, 2021.

[5] Jean-Marc Valin, Timothy B Terriberry, Christopher Montgomery, and Gregory Maxwell, "A high-quality speech and audio codec with less than 10-ms delay," *IEEE transactions on audio, speech, and language processing*, vol. 18, no. 1, pp. 58–67, 2009.

[6] Yi-Chiao Wu, Israel D Gebru, Dejan Marković, and Alexander Richard, "Audiodec: An open-source streaming high-fidelity neural audio codec," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.

[7] Sergio Casas, Abbas Sadat, and Raquel Urtasun, "Mp3: A unified model to map, perceive, predict and plan," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14403–14412.

[8] Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez, "Simple and controllable music generation," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[9] Shuochen Gao, Shun Lei, Fan Zhuo, Hangyu Liu, Feng Liu, Boshi Tang, Qiaochu Huang, Shiyin Kang, and Zhiyong Wu, "An end-to-end approach for chord-conditioned song generation," 2024.

[10] Shun Lei, Yixuan Zhou, Boshi Tang, Max W. Y. Lam, Feng Liu, Hangyu Liu, Jingcheng Wu, Shiyin Kang, Zhiyong Wu, and Helen Meng, "Songcreator: Lyrics-based universal song generation," 2024.

[11] Xiaofei Wang, Manthan Thakker, Zhuo Chen, Naoyuki Kanda, Sefik Emre Eskimez, Sanyuan Chen, Min Tang, Shujie Liu, Jinyu Li, and Takuya Yoshioka, "Speechx: Neural codec language model as a versatile speech transformer," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.

[12] Yi-Chiao Wu, Israel D Gebru, Dejan Marković, and Alexander Richard, "Audiodec: An open-source streaming high-fidelity neural audio codec," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.

[13] Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, et al., "Neural codec language models are zero-shot text to speech synthesizers," *arXiv preprint arXiv:2301.02111*, 2023.

[14] Zhihao Du, Shiliang Zhang, Kai Hu, and Siqi Zheng, "Funcodec: A fundamental, reproducible and integrable open-source toolkit for neural speech codec," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 591–595.

[15] Yang Ai, Xiao-Hang Jiang, Ye-Xin Lu, Hui-Peng Du, and Zhen-Hua Ling, "Apcodec: A neural audio codec with parallel amplitude and phase spectrum encoding and decoding," *arXiv preprint arXiv:2402.10533*, 2024.

[16] Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi, "High fidelity neural audio compression," *arXiv preprint arXiv:2210.13438*, 2022.

[17] Srikanth Korse, Nicola Pia, Kishan Gupta, and Guillaume Fuchs, "Postgan: A gan-based post-processor to enhance the quality of coded speech," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 831–835.

[18] Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar, "High-fidelity audio compression with improved rvqgan," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[19] Arijit Biswas and Dai Jia, "Audio codec enhancement with generative adversarial networks," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 356–360.

[20] Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising diffusion probabilistic models," *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.

[21] Haohe Liu, Xuenan Xu, Yi Yuan, Mengyue Wu, Wenwu Wang, and Mark D Plumbley, "Semanticodec: An ultra low bitrate semantic audio codec for general sound," *arXiv preprint arXiv:2405.00233*, 2024.

[22] Philip Anastassiou, Jiawei Chen, Jitong Chen, Yuanzhe Chen, Zhuo Chen, Ziyi Chen, Jian Cong, Lelai Deng, Chuang Ding, Lu Gao, et al., "Seed-tts: A family of high-quality versatile speech generation models," *arXiv preprint arXiv:2406.02430*, 2024.

[23] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le, "Flow matching for generative modeling," *arXiv preprint arXiv:2210.02747*, 2022.

[24] Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D Plumbley, "Audioldm 2: Learning holistic audio generation with self-supervised pretraining," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.

[25] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae, "Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis," *Advances in neural information processing systems*, vol. 33, pp. 17022–17033, 2020.

[26] William Peebles and Saining Xie, "Scalable diffusion models with transformers," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 4195–4205.

[27] Chung-Cheng Chiu, James Qin, Yu Zhang, Jiahui Yu, and Yonghui Wu, "Self-supervised learning with random-projection quantizer for speech recognition," in *International Conference on Machine Learning*. PMLR, 2022, pp. 3915–3924.

[28] Christian Schörkhuber and Anssi Klapuri, "Constant-q transform toolbox for music processing," in *7th sound and music computing conference, Barcelona, Spain*. SMC, 2010, pp. 3–64.

[29] Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar, "High-fidelity audio compression with improved rvqgan," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[30] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM transactions on audio, speech, and language processing*, vol. 29, pp. 3451–3460, 2021.

[31] Yizhi Li, Ruibin Yuan, Ge Zhang, Yinghao Ma, Xingran Chen, Hanzhi Yin, Chenghao Xiao, Chenghua Lin, Anton Ragni, Emmanouil Benetos, et al., "Mert: Acoustic music understanding model with large-scale self-supervised training," *arXiv preprint arXiv:2306.00107*, 2023.

[32] Michael Schoeffler, Fabian-Robert Stöter, Bernd Edler, and Jürgen Herre, "Towards the next generation of web-based experiments: A case study assessing basic audio quality following the itu-r recommendation bs. 1534 (mushra)," in *1st Web Audio Conference*, 2015, pp. 1–6.

[33] Andrew Hines, Jan Skoglund, Anil C Kokaram, and Naomi Harte, "Visqol: an objective speech quality model," *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2015, pp. 1–18, 2015.

[34] Simon Rouard, Francisco Massa, and Alexandre Défossez, "Hybrid transformers for music source separation," in *ICASSP 23*, 2023.

[35] Siqi Zheng, Luyao Cheng, Yafeng Chen, Hui Wang, and Qian Chen, "3d-speaker: A large-scale multi-device, multi-distance, and multi-dialect corpus for speech representation disentanglement," *arXiv preprint arXiv:2306.15354*, 2023.

[36] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever, "Robust speech recognition via large-scale weak supervision," in *International conference on machine learning*. PMLR, 2023, pp. 28492–28518.

[37] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al., "Conformer: Convolution-augmented transformer for speech recognition," *arXiv preprint arXiv:2005.08100*, 2020.

[38] Alexey Dosovitskiy, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[39] Jonathan Ho and Tim Salimans, "Classifier-free diffusion guidance," *arXiv preprint arXiv:2207.12598*, 2022.