# AudioCodecBench: A Comprehensive Benchmark for Audio Codec Evaluation

**Lu Wang[1], Hao Chen[1*], Siyu Wu[1*], Zhiyue Wu[1*†], Hao Zhou[2], Chengfeng Zhang[3], Ting Wang[4], Haodi Zhang[1]**

[1] Shenzhen University [2] Nankai University [3] Zhejiang University [4] East China Normal University
{wanglu, hdzhang}@szu.edu.cn, twang@sei.ecnu.edu.cn

## Abstract

Multimodal Large Language Models (MLLMs) have been widely applied in speech and music. This tendency has led to a focus on audio tokenization for Large Models (LMs). Unlike semantic-only text tokens, audio tokens must both capture global semantic content and preserve fine-grained acoustic details. Moreover, they provide a discrete method for speech and music that can be effectively integrated into MLLMs. However, existing research is unsuitable in the definitions of semantic tokens and acoustic tokens. In addition, the evaluation of different codecs typically concentrates on specific domains or tasks, such as reconstruction or Automatic Speech Recognition (ASR) task, which prevents fair and comprehensive comparisons. To address these problems, this paper provides suitable definitions for semantic and acoustic tokens and introduces a systematic evaluation framework. This framework allows for a comprehensive assessment of codecs' capabilities which evaluate across four dimensions: audio reconstruction metric, codebook index (ID) stability, decoder-only transformer perplexity, and performance on downstream probe tasks. Our results show the correctness of the provided suitable definitions and the correlation among reconstruction metrics, codebook ID stability, downstream probe tasks and perplexity.

**Code**: https://github.com/wuzhiyue111/Codec-Evaluation
**Dataset**: https://huggingface.co/datasets/LeBeGut/Audio CodecBench

## Introduction

Discrete audio tokens have received attention for their potential to bridge the domains of text and audio, playing an important role in the development of Multimodal Large Language Models (MLLMs) (Liu et al. 2023; Team 2025). The process of generating discrete token is compressing the original waveform into a finite set of vectors. However, MLLMs focus more on semantic in the text domain, but need to focus on both semantic and acoustic in the audio domain, resulting in a modality gap between text and audio. Semantic tokens are often obtained through the quantization hidden states from Self-supervised Learning (SSL) models. These tokens fixed patterns in the same semantic informations so
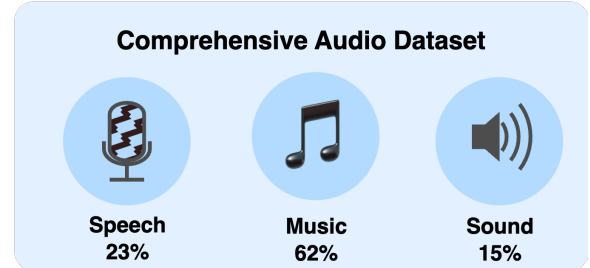
---

Figure 1: AudioCodecBench data distribution overview.

that the fixed patterns are easier to be modeled by downstream tasks (Baevski et al. 2020). Acoustic tokens are often obtained by training the neural audio codec (Codecs) in an end-to-end manner with the goal of high-fidelity reconstruction. These tokens focus more on the absolute distance between audio sampling points. This absolute distance definitely contains semantic, but this part of the semantic is difficult to be modeled in downstream tasks and is more suitable for reconstruction (Borsos et al. 2023; Zeghidour et al. 2021).

The core task of Large Language Models (LLMs) is to predict next token in a sequence. This mechanism requires that its input must be a series of discrete tokens. Therefore, researchers always adopt the discrete quantization methods (Mentzer et al. 2023; Yang et al. 2023). These methods aim to approximate a large, continuous vector space with a finite, discrete set of representative vectors, mapping high-dimensional continuous signals in a finite codebook. Therefore, the signal can be translated effectively into token sequences that LLMs can understand and generate. These discrete methods function as a clustering process to generate codebook indices. Whether these indices represent semantics or acoustics depends on the encoder. However, despite growing research on discrete tokens, there is still no unified framework to evaluate and compare the performance of different token types.

To address these shortcomings, this paper introduces a systematic, multi-dimensional benchmark for codec evaluation. This benchmark comprehensively assesses codec capabilities across four key experiments: **Reconstruction**, to assess audio reconstruction fidelity; **ID Sensitivity**, to evaluate
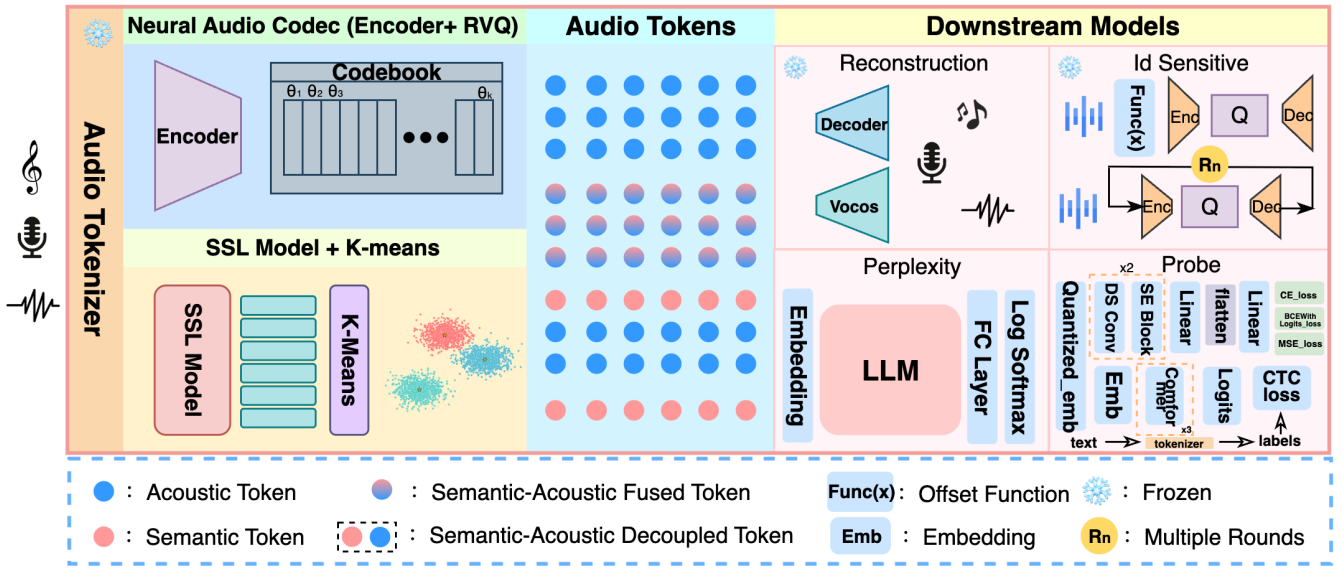
Figure 2: The proposed AudioCodecBench framework. Users provide pre-trained codec and obtain token-level outputs through encoding and quantization. Different types of tokens are input into different evaluation task components for multi-dimensional task evaluation.

codebook ID stability under noisy conditions; **Perplexity**, to measure the impact of different token sequences on Large Models (LMs) modeling and **Probe**, to evaluate downstream task performance. The distribution of datasets is illustrated in Figure 1. We hope that this benchmark will offer a more comprehensive comparison of various audio tokenization methods. Our contributions include the following:

- We provide suitable definitions of semantic and acoustic features. And base on their combination further define fused features and decoupled features.
- We evaluate four features across a variety of tasks in our benchmark. This benchmark considers multiple evaluation metrics, and also covers three audio domains: speech, music and sound.
- We explor the correlation between various task metrics and perplexity.

## Related Work

### Audio Representation

SpeechTokenizer (Zhang et al. 2024) distinguishes between "Semantic token" and "Acoustic token". Semantic token originates from SSL models like BEST-RQ (Chiu et al. 2022), HuBERT (Hsu et al. 2021), Wav2Vec2 (Baevski et al. 2020) and WavLM (Chen et al. 2022). These models typically employ BERT-like structures and MLM loss to capture global contextual information, and it is often assume that semantics can be equated with performance on the Automatic Speech Recognition (ASR) task. However, we think that semantics is not only responded by ASR performance. In contrast, acoustic tokens are generated by codecs like En-Codec (Défossez et al. 2022), SoundStream and DAC (Kumar et al. 2024) employ VQ-VAE driven by reconstruction

loss to achieve high-fidelity reconstruction. This concept of audio representation provides a foundation for systematically analyzing the information types of discrete tokens.

To leverage the strengths of both token types, subsequent research explores different paradigms. SemantiCodec (Liu et al. 2024) and XY-Tokenizer (Gong et al. 2025) employs a dual-encoder architecture to decouple acoustic and semantic tokens by reconstruction loss and k-means clustering. In contrast, models like XCodec (Ye et al. 2024) (Yuan et al. 2025) directly concatenate the two token types at the feature level. Meanwhile, SpeechTokenizer and Mimi (Défossez et al. 2024) introduce a "semantic distillation" approach. It uses an SSL model to guide the encoder of codec so that its discrete tokens carry both acoustic and semantic content in the first codebook. With the development of these different representation methods, establishing a fair and comprehensive evaluation becomes a significant challenge.

### SSL and Codec Benchmark

Evaluation of discrete audio representations presents a diverse challenge. SSL benchmarks like SUPERB (wen Yang et al. 2021) and MARBLE (Yuan et al. 2023) respectively evaluate representation performance on downstream tasks in the domains of speech and Music Information Retrieval. HEAR (Turian et al. 2022) further extends the downstream tasks to multiple domains of speech, environment sounds and music. Similar to HEAR, ARCH (La Quatra et al. 2024) introduces diverse datasets and offers a more extensible cross-domain evaluation framework than HEAR. However, a common limitation of these benchmarks is that they focus on downstream tasks, ignoring other evaluation aspects such as audio reconstruction and LM perplexity. Other methods of evaluation aspects like Code Drift (O'Reilly et al. 2025) evaluates the stability of multi-round reconstruction, while

| Feature Type | Model | Sample Rate | #Codebooks | Codebook Size | #Params | Bitrate (kbps) | Token Rate |
|---|---|---|---|---|---|---|---|
| **Acoustic** | DAC | 24kHz | 8 | 1024 | 74.7M | 6kbps | 75 |
| | EnCodec | 24kHz | 8 | 1024 | 14.9M | 6kbps | 75 |
| | WavTokenizer | 24kHz | 1 | 4096 | 103M | 0.48kbps | 40 |
| **Semantic** | HuBERT | 16kHz | - | - | 94.4M | - | 50 |
| | Qwen2Audio | 16kHz | - | - | 636M | - | 25 |
| **Semantic and Acoustic Fused** | SpeechTokenizer | 16kHz | 8 | 1024 | 80.9M | 5.33kbps | 50 |
| | Mimi | 24kHz | 8 | 2048 | 39.4M | 1.1kbps | 12.5 |
| | XCodec | 16kHz | 8 | 1024 | 123M | 4kbps | 50 |
| | YuE | 16kHz | 8 | 1024 | 123M | 4kbps | 50 |
| **Semantic and Acoustic Decoupled** | SemantiCodec | 16kHz | 2 | 8192 | 507M | 1.3kbps | 100 |

Table 1: The relevant feature types and parameters of the audio codecs and SSL models.

Codec-SUPERB (Wu et al. 2024) evaluates reconstruction fidelity. DASB (Mousavi et al. 2024) systematically probes discrete tokens in speech tasks.

To consolidate these diverse evaluation methods, researchers develop comprehensive toolkit like VERSA (Shi et al. 2025), and compile survey (Mousavi et al. 2025) to integrate existing methods within a unified framework. However, these evaluation methods typically evaluate the performance of discrete tokens from diverse tasks. As a result, they do not define the different types of information of semantic and acoustic. And exploring the different types connect to different tasks. Therefore, there is an urgent need to bridge this gap.

This paper first establishes a suitable definition of "semantic" that **must be strictly described by text**. Based on this, this paper further defines four different information types and compares the performance of discrete tokens of these four information types under different tasks. Through comprehensive experimental analysis, we explore the information types of various discrete tokens, providing insights to support the design of more effective audio representations.

## Evaluation Framework

### Overall Architecture

In the reconstruction task, we process an original audio signal through the encoder, quantizer, and decoder pipeline to reconstruct waveform, and use metrics like Perceptual Evaluation of Speech Quality (PESQ) (Rix et al. 2001), Short-Time Objective Intelligibility (STOI) (Taal et al. 2010) to evaluate the codec's ability to encode acoustic details; while using Word Error Rate (WER) and Character Error Rate(CER) to evaluate semantic preservation in acoustic details. The codec with higher metric scores is considered to have tokenization more focused on accurately reconstructing acoustic details.

The ID sensitivity experiment consists of two subtasks, as shown in the upper right section of the downstream model in Figure 2. The first task is multi-round reconstruction, we use the output of the $(n)th$ round as the input for the $(n+1)th$ round. The second task is the temporal shift sta-

bility experiment. We simulate signal phase shift by introducing millisecond-level time shifts into the original audio, and reconstruct this shifted audio. We define **ID sensitivity** as the stability of discrete tokens under noise interference. For both subtasks, we calculate the unchanged rate of IDs in the same codebook after the process to evaluate the representation's robustness. Higher stability indicates lower ID sensitivity, and conversely, lower stability indicates higher ID sensitivity.

For the perplexity experiment, we extract the sequence of discrete IDs from the codec, then train a small LM using the Cross-Entropy loss to predict next audio-only tokens. As shown in the lower left section of the downstream model in Figure 2. We use the perplexity of this LM as the evaluation metric to evaluate the adaptability of the discrete ID sequence for LM modeling. A lower perplexity indicates that the sequence is more amenable to LM modeling and also implies that it may contain richer semantics.

In the downstream probe model, we design two structures to evaluate the generalization of discrete tokens through various downstream tasks. As shown in the lower right section of the downstream model in Figure 2. The first is a lightweight network composed of SE-Blocks (Hu et al. 2019) (channel attention) and depthwise separable convolutions (Chollet 2017). This network compresses both the temporal and feature dimensions of the embedding after quantization and then makes predictions using task-specific heads. For the ASR task, we design a different approach to measure the alignment between the representation and text. The extracted discrete IDs are fed through an embedding layer into a three-layer Conformer network (Gulati et al. 2020), and the model is trained end-to-end using the Connectionist Temporal Classification (CTC) loss (Graves et al. 2006).

### Audio Feature Classification

We review existing definitions of audio representations (acoustic and semantic), but find these definitions fail to cover the current diverse features. **Therefore, we propose that a semantic feature must be strictly describable by text.** On this basis, we divide audio features into four categories.

| Audio Type | Task | Dataset | Metric |
|---|---|---|---|
| **Music** | Genre Classification(GC) | GTZAN (Tzanetakis and Cook 2002) | Accuracy |
| | Key Detection(KD) | GiantSteps Key (Knees et al. 2015) | Accuracy |
| | Emotion Detection(ED) | Emomusic (Soleymani et al. 2013) | $R^2_{\text{Valence}}$ & $R^2_{\text{Arousal}}$ |
| | | MTG MoodTheme (Bogdanov et al. 2019) | ROC-AUC & PR-AUC/AP |
| | Vocal Technique Detection(VTD) | VocalSet (Wilkins et al. 2018) | Accuracy |
| | Pitch Classification(PC) | NSynth (Engel et al. 2017) | Accuracy |
| | Music Tagging(MT) | MagnaTagATun (Law et al. 2009) | ROC-AUC & PR-AUC/AP |
| | | MTG Top50 (Bogdanov et al. 2019) | ROC-AUC & PR-AUC/AP |
| | Instrument Classification(IC) | NSynth (Engel et al. 2017) | Accuracy |
| | | MTG Instrument (Bogdanov et al. 2019) | ROC-AUC & PR-AUC/AP |
| | Singer Identification(SI) | VocalSet (Wilkins et al. 2018) | Accuracy |
| **Speech** | Automatic Speech Recognition(ASR) | Common Voice (Ardila et al. 2020) | WER,CER |
| | Emotion Detection(ED) | MELD (Poria et al. 2019) | Accuracy |
| **Sound** | Vocal Sound Classification(VSC) | VocalSound (Gong, Yu, and Glass 2022) | Accuracy |
| | Environmental Sound Classification(ESC) | ESC-50 (Trowitzsch et al. 2020) | Accuracy |

Table 2: The task, dataset and evaluation metric for the downstream probe.The following text will use abbreviations to replace the full names of various tasks, datasets, and evaluation materials. Dataset-related GiantSteps Key: GS, Emomusic: EMO, MTG MoodTheme: MTGMT, VocalSet: VST, NSynth: NS, MagnaTagATun: MTT, MTG Top50: MTGT, MTG Instrument: MTGI, Common Voice: CV, VocalSound: VSD. Material-related ROC-AUC & PR-AUC: RA.

**1) Acoustic feature**: The discrete feature **cannot be described by text**. These features originate from codecs optimized for reconstruction, representing the quantized encoding of acoustic details, such as environmental noise, vocal fold vibration and air vibration.

**2) Semantic feature**: The discrete feature extracted from MLM within SSL frameworks **must be strictly defined by text**. They aim to capture high-level and abstract information, such as the transcribed text of speech, the emotion or key of music and the human voice in music.

**3) Semantic-Acoustic fused feature**: The discrete features is **fused with both text-describable semantics and text-indescribable acoustic information**. For instance, features representing a specific speaker's voice simultaneously contain textual content and unique acoustic details.

**4) Semantic-Acoustic decoupled feature**: The discrete features that **separates text-describable semantics and text-indescribable acoustic information into independent codebooks**. For the same audio clip of 'Hello', it outputs two independent token streams: one representing the text-describable information 'Hello', and the other representing text-indescribable acoustic details such as the speaker's unique acoustic signature.

Based on the definitions of the four feature classes, the codecs and SSL models evaluated in this paper are classified accordingly. Table 1 provides a summary of these models, detailing the model feature types they generate and key technical specifications such as sample rate, bit rate, and token rate.

## Experiments and Analysis

We evaluate the performance of eight codecs and two SSL models. Their relevant attributes are listed in Table 1. We use the first 8 codebooks to evaluate the performance of the multi-codebook codecs.

### Reconstruction

We conduct reconstruction experiment on the LibriTTS test-other (Zen et al. 2019) and GTZAN test datasets. In Table 3, the results are rounded to the required precision for each metric. The left side of each metric result is the speech dataset result, and the right side is the music dataset result. Since Mimi and SpeechTokenizer are not trained on music datasets, they are not evaluated on music dataset experiments.

On the speech dataset, acoustic codecs such as DAC and EnCodec achieve the highest reconstruction fidelity. Codecs that integrate semantics like XCodec and YuE demonstrate the suboptimal performance, while WavTokenizer performs the worst. The result suggests that semantics may affect the reconstruction of acoustic details. Although WavTokenizer's discrete tokens are acoustic, its reconstruction quality is weak. We speculate that to balance compression bitrate and reconstruction quality, **small codebook size and few codebooks limit the variety of combinations for the discrete tokens, which weakens the ability of these tokens to capture acoustic details**.

All reconstruction metrics are lower on the music dataset compared to the speech dataset. This is because music contains more intricate harmonic structures and richer dynamic variations than speech. Therefore, music is more difficult

| Codec | PESQ↑ | Spk-Sim↑ | WER (GT/REC)↓ | CER (GT/REC)↓ | STOI↑ |
|---|---|---|---|---|---|
| DAC | **3.69 / 2.66** | **0.965** / - | 0.155 / 0.202 \| - / - | 0.09 / 0.125 \| - / - | **0.94 / 0.86** |
| EnCodec | 3.21 / 2.27 | 0.919 / - | 0.155 / 0.198 \| - / - | 0.09 / 0.114 \| - / - | 0.93 / 0.85 |
| Mimi | 2.77 / - | 0.928 / - | 0.155 / 0.287 \| - / - | 0.09 / 0.173 \| - / - | 0.88 / - |
| SemantiCodec | 2.64 / 1.32 | 0.907 / - | 0.155 / 0.318 \| - / - | 0.09 / 0.195 \| - / - | 0.86 / 0.60 |
| WavTokenizer | 2.17 / 1.14 | 0.743 / - | 0.155 / 0.494 \| - / - | 0.09 / 0.325 \| - / - | 0.83 / 0.49 |
| SpeechTokenizer | 2.97 / - | 0.924 / - | 0.155 / 0.216 \| - / - | 0.09 / 0.120 \| - / - | 0.89 / - |
| XCodec | 3.23 / 1.85 | 0.942 / - | **0.155 / 0.185** \| - / - | **0.09 / 0.106** \| - / - | 0.91 / 0.76 |
| YuE | 3.17 / 1.84 | 0.938 / - | 0.155 / 0.195 \| - / - | 0.09 / 0.113 \| - / - | 0.90 / 0.75 |

Table 3: Reconstruction results of difference codecs in LibriTTS test-other dataset and GTZAN test dataset.
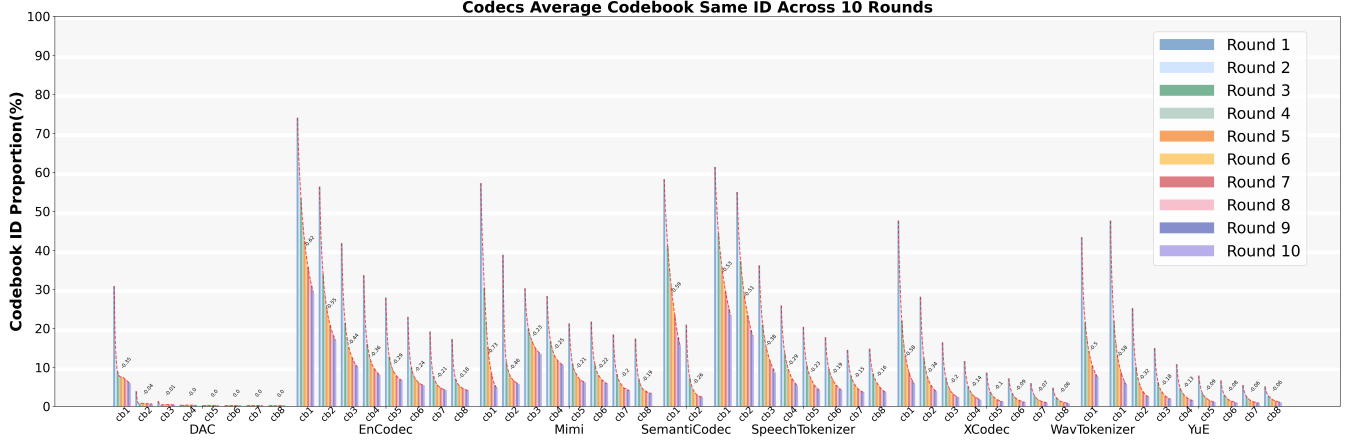


Figure 3: The percentage of the same ID in each codebook of the codecs after multi-round reconstruction, cb stands for codebook.

to model and reconstruct with high-fidelity. Notably, the performance of WavTokenizer and SemantiCodec decreases significantly. This result further highlights the limitations of small codebook size and the single or dual-codebook quantization strategies. **Small codebook size and few codebooks limit the possibility of token combinations to represent the acoustic details of music, thus reducing reconstruction fidelity**. In particular, WavTokenizer exhibits poor modeling capabilities for music, resulting in a decrease in subjective listening quality after reconstruction.

### ID sensitivity

We evaluate ID sensitivity through multi-round ($n = 10$) reconstruction and time shift task. The results are shown in Figure 3. Detailed results for different codecs are shown in Appendix A. After multi-round reconstruction, the codebook IDs of all codecs shift compared to the first round. Codecs focusing on acoustic reconstruction show higher ID stability (**lower slope** of the decrease rate of the same ID). The result indicates that they can accurately reconstruct the signal, including some possible noise. In contrast, codecs that integrate semantics exhibit lower ID stability (**higher slope**). The result shows that these codecs are less sensitive to fitting noise during reconstruction and focus more on ensuring semantic. Although EnCodec generates tokens that are mainly acoustic, its multi-round reconstruction perfor-

mance is similar to the codecs integrating semantics. This may be attributed to EnCodec's integration of LSTM modules during encoding, which capture long-context dependencies, enhancing the stability of multi-round reconstruction.

Inspired by Code Drift (O'Reilly et al. 2025), we select 2ms as the experimental setting for time shift task, the results are shown in Figure 4. Detailed results for different codecs are shown in Appendix B. The result demonstrates that the token sequences of acoustic codecs are sensitive to temporal changes, as they focus on reconstruction and attempt to encode all acoustic details, including slight timing shifts. And codecs that integrate semantics focus more on stable content features, thus demonstrating greater robustness to slight timing shifts. **Codecs that integrate semantics outperform the acoustic codecs on the same ID ratio metric, which indicates that semantic-dominant tokens are more robust to slight timing shifts.**

### Perplexity

We train a 100M LM using Qwen2 architecture (Chu et al. 2024) from scratch to evaluate the modeling efficiency of codecs via validation set perplexity (PPL). For multi-codebook codecs, we apply a parallel evaluation (Yang et al. 2025) to compute PPL for each codebook. To ensure a fair comparison, the PPL values are normalized, and the final PPL is calculated using a mean loss. Be-
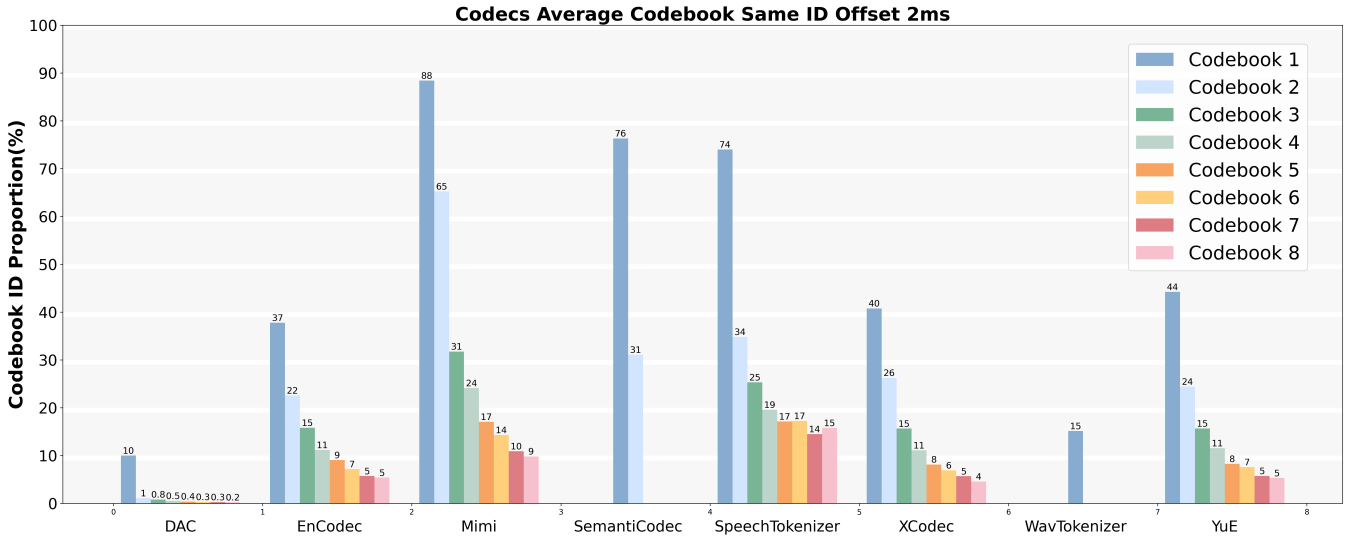
Figure 4: The proportion of identical IDs in each codebook of the codecs after time shift processing and reconstruction.

| Codec | ppl↓ | cb1_ppl | cb2_ppl | cb3_ppl | cb4_ppl | cb5_ppl | cb6_ppl | cb7_ppl | cb8_ppl |
|---|---|---|---|---|---|---|---|---|---|
| DAC | 247 / 194 | 21 / 29 | 147 / 123 | 218 / 152 | 315 / 213 | 396 / 2701 | 483 / 353 | 570 / 413 | 628 / 474 |
| EnCodec | 76 / 141 | 15 / 18 | 33 / 63 | 59 / 111 | 89 / 170 | 111 / 226 | 138 / 287 | 159 / 337 | 173 / 376 |
| WavTokenizer | 105 / 38 | 105 / 38 | - / - | - / - | - / - | - / - | - / - | - / - | - / - |
| X-Codec | 30 / 48 | 10 / 20 | 13 / 20 | 20 / 32 | 31 / 51 | 42 / 65 | 51 / 75 | 62 / 87 | 71 / 100 |
| YuE | 29 / 46 | 9 / 18 | 16 / 29 | 20 / 30 | 29 / 48 | 39 / 60 | 51 / 75 | 55 / 83 | 54 / 76 |
| SpeechTokenizer | 14 / - | 2 / - | 6 / - | 12 / - | 18 / - | 22 / - | 25 / - | 29 / - | 31 / - |
| Mimi | 127 / - | 9 / - | 58 / - | 148 / - | 185 / - | 229 / - | 257 / - | 279 / - | 299 / - |
| SemantiCodec | **8** / **16** | 1 / 1 | 82 / 272 | - / - | - / - | - / - | - / - | - / - | - / - |

Table 4: PPL results of different codecs in Emilia-EN dataset and MTG-Jamendo dataset, cb stands for codebook.

cause PPL scores are directly influenced by the codebook size; larger codebooks typically result in higher PPL. Therefore, we normalize all values to a reference codebook size of 1024. The calculation is as follows:

$$\text{PPL} = \frac{\exp(\mathcal{L}_{CE})}{S_{\text{cb}}/1024} \quad (1)$$

where $\mathcal{L}_{CE}$ is the average cross-entropy loss calculated over the entire token sequence, and $S_{\text{cb}}$ denotes the codec codebook size. The training runs for 100k steps on 8 NVIDIA A6000 GPUs using the Emilia-EN (He et al. 2024) and MTG-Jamendo datasets. Table 4 presents the results, rounded to the nearest integer. The left side of PPL metric result is the speech dataset result, and the right side is the music dataset result. Since Mimi and SpeechTokenizer are not trained on music datasets, they are not evaluated on music dataset experiments.

On the speech dataset, codecs that integrate semantics achieve better results than acoustic codecs. This result demonstrates that **semantic tokens are easier for LMs to model**. Analysis of the multi-codebook codecs' results shows that earlier codebooks have lower PPL, which provides strong support for the conclusion that semantics is beneficial for LM modeling. Although EnCodec mainly gen-

erates acoustic tokens, it achieves unexpectedly low PPL. Mimi uses a semantic teacher to guide its first quantizer, but it fails to achieve the performance of other codecs that integrate semantics. The exact reasons behind these unusual results are still unknown and need further exploration.

The PPL is higher on the music dataset compared to the speech dataset, the finding that is consistent with human intuition. This is because music involves multiple instruments and complex temporal structures. These factors create a larger variety of possible token combinations, making their distribution much sparser than in speech. However, the PPL values for DAC and WavTokenizer on the music dataset are unexpectedly lower than on the speech dataset. We speculate that this is because DAC and WavTokenizer were trained on the MTG-Jamendo dataset but not on the Emilia-EN dataset, so their ppl results are different from other codecs.

**Probe**

In the downstream probe tasks, to ensure fair and reliable results, all experiments are conducted under the same computational budget. For the ASR task, we select Speech2Text (Ott et al. 2019; Wang et al. 2020) as the tokenizer for word segmentation. The related tasks, datasets,

| Task | GC | ED | | | | MT | | | | IC | | | KD | VTD | PC | SI |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Dataset | GTZAN | EMO | | MTGMT | | MTT | | MTGT | | NS | MTGI | | GS | VST | NS | VST |
| Metrics | Acc↑ | $R^2_A$↑ | $R^2_V$↑ | AP↑ | RA↑ | AP↑ | RA↑ | AP↑ | RA↑ | Acc↑ | AP↑ | RA↑ | Acc↑ | Acc↑ | Acc↑ | Acc↑ |
| DAC | 0.58 | 0.47 | 0.06 | 0.08 | 0.65 | 0.20 | 0.79 | 0.14 | 0.69 | 0.60 | 0.11 | 0.64 | 0.09 | 0.38 | 0.47 | 0.42 |
| EnCodec | 0.57 | 0.47 | 0.07 | 0.06 | 0.64 | 0.18 | 0.76 | 0.14 | 0.70 | 0.54 | 0.10 | 0.62 | 0.10 | 0.30 | 0.55 | 0.30 |
| WavTokenizer | 0.42 | 0.46 | 0.07 | 0.06 | 0.63 | 0.17 | 0.74 | 0.14 | 0.70 | 0.54 | 0.11 | 0.64 | 0.09 | 0.29 | 0.44 | 0.13 |
| SemantiCodec | **0.70** | 0.51 | **0.32** | 0.10 | **0.72** | 0.32 | **0.88** | **0.23** | **0.80** | **0.66** | 0.15 | **0.72** | 0.34 | 0.45 | 0.76 | 0.34 |
| XCodec | 0.66 | 0.55 | 0.14 | 0.10 | 0.71 | **0.32** | 0.87 | 0.22 | 0.78 | 0.64 | **0.16** | 0.71 | **0.46** | 0.57 | **0.91** | **0.54** |
| YuE | 0.67 | **0.57** | 0.16 | **0.10** | 0.71 | 0.32 | 0.87 | 0.19 | 0.76 | 0.62 | 0.13 | 0.70 | 0.45 | **0.59** | 0.90 | 0.52 |

Table 5: The results of various detection tasks performed by the codecs across different music datasets.

| Task | ASR | | VSC | ESC | ED |
|---|---|---|---|---|---|
| Dataset | CV | | VSD | ESC-50 | MELD |
| Metrics | WER↓ | CER↓ | Acc↑ | Acc↑ | Acc↑ |
| DAC | 0.53 | 0.23 | 0.54 | 0.33 | 0.48 |
| EnCodec | 0.50 | 0.21 | 0.57 | 0.28 | 0.48 |
| WavTokenizer | 0.58 | 0.29 | 0.52 | 0.14 | 0.48 |
| SemantiCodec | 0.49 | 0.20 | 0.72 | 0.62 | 0.48 |
| Mimi | **0.44** | **0.17** | 0.83 | 0.34 | 0.48 |
| SpeechTokenizer | 0.47 | 0.19 | 0.78 | 0.67 | 0.50 |
| XCodec | 0.47 | 0.19 | 0.73 | 0.64 | 0.49 |
| YuE | 0.47 | 0.19 | 0.78 | 0.64 | 0.52 |
| HuBERT | - | - | 0.88 | 0.53 | 0.50 |
| Qwen2Audio | - | - | **0.95** | **0.98** | **0.59** |

Table 6: The results of various detection tasks performed by the codecs and SSL models across different speech datasets.

| Task | Metric | r |
|---|---|---|
| Reconstruction | $WER_{REC}$ | 0.06 |
| | $CER_{REC}$ | 0.1 |
| | PESQ | -0.35 |
| | Spk_Smi | -0.05 |
| | STOI | -0.35 |
| ID sensitivity | MRC | 0.52 |
| | OS | 0.44 |
| Probe | $WER_{CTC}$ | 0.37 |
| | $CER_{CTC}$ | 0.36 |
| | $VSD_{ACC}$ | 0.55 |

Table 7: Pearson correlation coefficient between PPL and metrics from various speech evaluation tasks.

and evaluation metrics are shown in Table 2. Detailed introductions are shown in Appendix D.

The results of the music probe task are shown in Table 5. The visualized result is shown in Figure 22 in Appendix C. In the ED task, SemantiCodec's performance on Valence prediction is the best. This result supports the conclusion that Valence is closely associated with semantics (Asgari et al. 2014). Tasks such as MT, GC and KD involve high-level musical structures, SemantiCodec shows advantages in these tasks. Meanwhile, XCodec and SemantiCodec also achieved better performance in IC and PC tasks, which closely related to symbolic music information. Although tasks like SI and VTD rely on acoustic properties like timbre, codecs that combine both semantics and acoustics perform even better than acoustic codecs. This suggests that representation containing both semantics and acoustics makes information such as timbre more easily utilized by downstream models. Among these results, codecs that integrate semantics show better performance compared to acoustic codecs in these tasks. And these results also shows that semantics are crucial for the modeling of downstream tasks.

The results of the speech and sound probe tasks are shown in Table 6. The visualized result is shown in Figure 21 in Appendix C. The SSL models achieve the best performance. Codecs that integrate semantics demonstrate the suboptimal performance. Acoustic codecs perform the worst. WavTokenizer achieves the lowest performance. In the ASR task, codecs that explicitly introduce semantics generally achieve better WER/CER scores than acoustic codecs. In the VSD task, codecs that combine both semantics and acoustics show outstanding performance. It further suggests that timbre information may be effectively retained and utilized in representations that contain both semantics and acoustics. In the ED task, the performance of different codecs is relatively balanced. This suggests that the emotion-related features required for this specific task can be fully fitted by codecs. Notably, in the ESC task, Mimi performs worse than other semantic codecs. This may be attributed to Mimi's use of a larger hop length setting. This setting reduces its temporal resolution and weakens its ability to capture the transient features for ESC task.

In order to explore the impact of various metrics on the PPL value of LMs, we calculate the Pearson correlation coefficients between various task metrics and PPL. We aim to reveal which codec properties or audio features are more beneficial for LM modeling. The correlation coefficient is performed on the speech dataset and the results are shown in Table 7. PPL is positively correlated with CTC probe task metrics, which demonstrates that tokens rich in semantic content are easier for LMs to model. However, it shows a negative correlation with objective acoustic reconstruction

metrics, indicating that overfitting acoustic details may increase the difficulty of LM modeling. Moreover, the reconstructed WER/CER metrics, which should be negatively correlated with PPL, show the opposite of our expected conclusion. We speculate that XCodec and YuE use different methods of integrating semantics. Their reconstruction metrics just below acoustic codecs like DAC and EnCodec. The codebook ID sensitivity metrics show a positive correlation with ppl, which indicates that the LM's modeling capability is more dependent on the semantic content encoded in the token sequence.

## Conlusion

This paper presents a comprehensive, fair and highly reusable evaluation framework for codecs. We first redefine "acoustic" and "semantic" features: **semantic features must be strictly described by text**. Based on this classification, our benchmark systematically evaluates the performance of different discrete tokens across multiple tasks, and breaking the limitation of measuring semantics through ASR performance. Experimental results not only show the potential applications of various representations in MLLMs but also point to a new research direction: training better audio-semantic models by aligning text modality. We are committed to providing an open and fair benchmark and hope to attract more researchers to participate, jointly advancing the field of audio representation learning.

## References

Ardila, R.; Branson, M.; Davis, K.; Henretty, M.; Kohler, M.; Meyer, J.; Morais, R.; Saunders, L.; Tyers, F. M.; and Weber, G. 2020. Common Voice: A Massively-Multilingual Speech Corpus. arXiv:1912.06670.

Asgari, M.; Kiss, G.; Van Santen, J.; Shafran, I.; and Song, X. 2014. Automatic measurement of affective valence and arousal in speech. In *2014 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, 965–969. IEEE.

Baevski, A.; Zhou, Y.; Mohamed, A.; and Auli, M. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. In *Advances in Neural Information Processing Systems*, volume 33, 12449–12460.

Bogdanov, D.; Won, M.; Tovstogan, P.; Porter, A.; and Serra, X. 2019. The mtg-jamendo dataset for automatic music tagging. In *Proceedings of the Machine Learning for Music Discovery Workshop, International Conference on Machine Learning (ICML 2019)*, 1–3.

Borsos, Z.; Marinier, R.; Vincent, D.; Kharitonov, E.; Pietquin, O.; Sharifi, M.; Teboul, O.; Grangier, D.; Tagliasacchi, M.; and Zeghidour, N. 2023. Audiolm: a language modeling approach to audio generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31: 2523–2533.

Chen, S.; Wang, C.; Chen, Z.; Wu, Y.; Liu, S.; Chen, Z.; Li, J.; Kanda, N.; Yoshioka, T.; Xiao, X.; et al. 2022. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6): 1505–1518.

Chiu, C.-C.; Qin, J.; Zhang, Y.; Yu, J.; and Wu, Y. 2022. Self-supervised learning with random-projection quantizer for speech recognition. In Chaudhuri, K.; Jegelka, S.; Song, L.; Szepesvari, C.; Niu, G.; and Sabato, S., eds., *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, 3915–3924. PMLR.

Chollet, F. 2017. Xception: Deep Learning with Depthwise Separable Convolutions. arXiv:1610.02357.

Chu, Y.; Xu, J.; Yang, Q.; Wei, H.; Wei, X.; Guo, Z.; Leng, Y.; Lv, Y.; He, J.; Lin, J.; Zhou, C.; and Zhou, J. 2024. Qwen2-Audio Technical Report. arXiv:2407.10759.

Défossez, A.; Copet, J.; Synnaeve, G.; and Adi, Y. 2022. High Fidelity Neural Audio Compression. arXiv:2210.13438.

Défossez, A.; Mazaré, L.; Orsini, M.; Royer, A.; Pérez, P.; Jégou, H.; Grave, E.; and Zeghidour, N. 2024. Moshi: a speech-text foundation model for real-time dialogue. arXiv:2410.00037.

Engel, J.; Resnick, C.; Roberts, A.; Dieleman, S.; Norouzi, M.; Eck, D.; and Simonyan, K. 2017. Neural audio synthesis of musical notes with wavenet autoencoders. In *Proceedings of the International Conference on Machine Learning (ICML)*, 1068–1077.

Gong, Y.; Jin, L.; Deng, R.; Zhang, D.; Zhang, X.; Cheng, Q.; Fei, Z.; Li, S.; and Qiu, X. 2025. XY-Tokenizer: Mitigating the Semantic-Acoustic Conflict in Low-Bitrate Speech Codecs. arXiv:2506.23325.

Gong, Y.; Yu, J.; and Glass, J. 2022. Vocalsound: A dataset for improving human vocal sounds recognition. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 151–155.

Graves, A.; Fernández, S.; Gomez, F.; and Schmidhuber, J. 2006. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd International Conference on Machine Learning*, ICML '06, 369–376. New York, NY, USA: Association for Computing Machinery. ISBN 1595933832.

Gulati, A.; Qin, J.; Chiu, C.-C.; Parmar, N.; Zhang, Y.; Yu, J.; Han, W.; Wang, S.; Zhang, Z.; Wu, Y.; and Pang, R. 2020. Conformer: Convolution-augmented Transformer for Speech Recognition. arXiv:2005.08100.

He, H.; Shang, Z.; Wang, C.; Li, X.; Gu, Y.; Hua, H.; Liu, L.; Yang, C.; Li, J.; Shi, P.; Wang, Y.; Chen, K.; Zhang, P.; and Wu, Z. 2024. Emilia: An Extensive, Multilingual, and Diverse Speech Dataset for Large-Scale Speech Generation. arXiv:2407.05361.

Hsu, W.-N.; Bolte, B.; Tsai, Y.-H. H.; Lakhotia, K.; Salakhutdinov, R.; and Mohamed, A. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29: 3451–3460.

Hu, J.; Shen, L.; Albanie, S.; Sun, G.; and Wu, E. 2019. Squeeze-and-Excitation Networks. arXiv:1709.01507.

Knees, P.; Faraldo, A.; Herrera, P.; Vogl, R.; Böck, S.; Hörschläger, F.; and Le Goff, M. 2015. Two Data Sets for

Tempo Estimation and Key Detection in Electronic Dance Music Annotated from User Corrections. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 364–370.

Kumar, R.; Seetharaman, P.; Luebs, A.; Kumar, I.; and Kumar, K. 2024. High-fidelity audio compression with improved rvqgan. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36.

La Quatra, M.; Koudounas, A.; Vaiani, L.; Baralis, E.; Cagliero, L.; Garza, P.; and Siniscalchi, S. M. 2024. Benchmarking Representations for Speech, Music, and Acoustic Events. In *2024 IEEE International Conference on Acoustics, Speech, and Signal Processing Workshops (ICASSPW)*, 505–509. IEEE.

Law, E.; West, K.; Mandel, M. I.; Bay, M.; and Downie, J. S. 2009. Evaluation of algorithms using games: The case of music tagging. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 387–392.

Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2023. Visual Instruction Tuning. arXiv:2304.08485.

Liu, H.; Xu, X.; Yuan, Y.; Wu, M.; Wang, W.; and Plumbley, M. D. 2024. SemantiCodec: An Ultra Low Bitrate Semantic Audio Codec for General Sound. *IEEE Journal of Selected Topics in Signal Processing*, 18(8): 1448–1461.

Mentzer, F.; Minnen, D.; Agustsson, E.; and Tschannen, M. 2023. Finite Scalar Quantization: VQ-VAE Made Simple. arXiv:2309.15505.

Mousavi, P.; Libera, L. D.; Duret, J.; Ploujnikov, A.; Subakan, C.; and Ravanelli, M. 2024. DASB - Discrete Audio and Speech Benchmark. arXiv:2406.14294.

Mousavi, P.; Maimon, G.; Moumen, A.; Petermann, D.; Shi, J.; Wu, H.; Yang, H.; Kuznetsova, A.; Ploujnikov, A.; Marxer, R.; Ramabhadran, B.; Elizalde, B.; Lugosch, L.; Li, J.; Subakan, C.; Woodland, P.; Kim, M.; yi Lee, H.; Watanabe, S.; Adi, Y.; and Ravanelli, M. 2025. Discrete Audio Tokens: More Than a Survey! arXiv:2506.10274.

O'Reilly, P.; Seetharaman, P.; Su, J.; Jin, Z.; and Pardo, B. 2025. Code Drift: Towards Idempotent Neural Audio Codecs. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1–5.

Ott, M.; Edunov, S.; Baevski, A.; Fan, A.; Gross, S.; Ng, N.; Grangier, D.; and Auli, M. 2019. fairseq: A Fast, Extensible Toolkit for Sequence Modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*.

Poria, S.; Hazarika, D.; Majumder, N.; Naik, G.; Cambria, E.; and Mihalcea, R. 2019. MELD: A Multimodal Multi-Party Dataset for Emotion Recognition in Conversations. arXiv:1810.02508.

Rix, A.; Beerends, J.; Hollier, M.; and Hekstra, A. 2001. Perceptual evaluation of speech quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs. In *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.01CH37221)*, volume 2, 749–752 vol.2.

Shi, J.; jin Shim, H.; Tian, J.; Arora, S.; Wu, H.; Petermann, D.; Yip, J. Q.; Zhang, Y.; Tang, Y.; Zhang, W.; Alharthi, D. S.; Huang, Y.; Saito, K.; Han, J.; Zhao, Y.; Donahue, C.; and Watanabe, S. 2025. VERSA: A Versatile Evaluation Toolkit for Speech, Audio, and Music. arXiv:2412.17667.

Soleymani, M.; Caro, M. N.; Schmidt, E. M.; Sha, C. Y.; and Yang, Y. H. 2013. 1000 songs for emotional analysis of music. In *Proceedings of the 2nd ACM International Workshop on Crowdsourcing for Multimedia*, 1–6.

Taal, C. H.; Hendriks, R. C.; Heusdens, R.; and Jensen, J. 2010. A short-time objective intelligibility measure for time-frequency weighted noisy speech. In *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*, 4214–4217.

Team, C. 2025. Chameleon: Mixed-Modal Early-Fusion Foundation Models. arXiv:2405.09818.

Trowitzsch, I.; Taghia, J.; Kashef, Y.; and Obermayer, K. 2020. The NIGENS General Sound Events Database. arXiv:1902.08314.

Turian, J.; Shier, J.; Khan, H. R.; Raj, B.; Schuller, B. W.; Steinmetz, C. J.; Malloy, C.; Tzanetakis, G.; Velarde, G.; McNally, K.; Henry, M.; Pinto, N.; Noufi, C.; Clough, C.; Herremans, D.; Fonseca, E.; Engel, J.; Salamon, J.; Esling, P.; Manocha, P.; Watanabe, S.; Jin, Z.; and Bisk, Y. 2022. Hear: Holistic evaluation of audio representations. In *Proceedings of the NeurIPS 2021 Competitions and Demonstrations Track*, volume 176, 125–145.

Tzanetakis, G.; and Cook, P. 2002. Musical genre classification of audio signals. *IEEE Transactions on Speech and Audio Processing*, 10(5): 293–302.

Wang, C.; Tang, Y.; Ma, X.; Wu, A.; Okhonko, D.; and Pino, J. 2020. fairseq S2T: Fast Speech-to-Text Modeling with fairseq. In *Proceedings of the 2020 Conference of the Asian Chapter of the Association for Computational Linguistics (AACL): System Demonstrations*.

wen Yang, S.; Chi, P.-H.; Chuang, Y.-S.; Lai, C.-I. J.; Lakhotia, K.; Lin, Y. Y.; Liu, A. T.; Shi, J.; Chang, X.; Lin, G.-T.; Huang, T.-H.; Tseng, W.-C.; tik Lee, K.; Liu, D.-R.; Huang, Z.; Dong, S.; Li, S.-W.; Watanabe, S.; Mohamed, A.; and yi Lee, H. 2021. SUPERB: Speech processing Universal PERformance Benchmark. arXiv:2105.01051.

Wilkins, J.; Seetharaman, P.; Wahl, A.; and Pardo, B. 2018. Vocalset: A singing voice dataset. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 468–474.

Wu, H.; Chung, H.-L.; Lin, Y.-C.; Wu, Y.-K.; Chen, X.; Pai, Y.-C.; Wang, H.-H.; Chang, K.-W.; Liu, A. H.; and yi Lee, H. 2024. Codec-SUPERB: An In-Depth Analysis of Sound Codec Models. arXiv:2402.13071.

Yang, D.; Liu, S.; Guo, H.; Zhao, J.; Wang, Y.; Wang, H.; Ju, Z.; Liu, X.; Chen, X.; Tan, X.; Wu, X.; and Meng, H. 2025. ALMTokenizer: A Low-bitrate and Semantic-rich Audio Codec Tokenizer for Audio Language Modeling. arXiv:2504.10344.

Yang, D.; Liu, S.; Huang, R.; Tian, J.; Weng, C.; and Zou, Y. 2023. HiFi-Codec: Group-residual Vector quantization for High Fidelity Audio Codec. arXiv:2305.02765.

Ye, Z.; Sun, P.; Lei, J.; Lin, H.; Tan, X.; Dai, Z.; Kong, Q.; Chen, J.; Pan, J.; Liu, Q.; Guo, Y.; and Xue, W. 2024. Codec Does Matter: Exploring the Semantic Shortcoming of Codec for Audio Language Model. *arXiv preprint arXiv:2408.17175*.

Yuan, R.; Lin, H.; Guo, S.; Zhang, G.; Pan, J.; Zang, Y.; Liu, H.; Liang, Y.; Ma, W.; Du, X.; Du, X.; Ye, Z.; Zheng, T.; Jiang, Z.; Ma, Y.; Liu, M.; Tian, Z.; Zhou, Z.; Xue, L.; Qu, X.; Li, Y.; Wu, S.; Shen, T.; Ma, Z.; Zhan, J.; Wang, C.; Wang, Y.; Chi, X.; Zhang, X.; Yang, Z.; Wang, X.; Liu, S.; Mei, L.; Li, P.; Wang, J.; Yu, J.; Pang, G.; Li, X.; Wang, Z.; Zhou, X.; Yu, L.; Benetos, E.; Chen, Y.; Lin, C.; Chen, X.; Xia, G.; Zhang, Z.; Zhang, C.; Chen, W.; Zhou, X.; Qiu, X.; Dannenberg, R.; Liu, J.; Yang, J.; Huang, W.; Xue, W.; Tan, X.; and Guo, Y. 2025. YuE: Scaling Open Foundation Models for Long-Form Music Generation. arXiv:2503.08638.

Yuan, R.; Ma, Y.; Li, Y.; Zhang, G.; Chen, X.; Yin, H.; others; and Fu, J. 2023. Marble: Music audio representation benchmark for universal evaluation. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, 39626–39647.

Zeghidour, N.; Luebs, A.; Omran, A.; Skoglund, J.; and Tagliasacchi, M. 2021. SoundStream: An End-to-End Neural Audio Codec. arXiv:2107.03312.

Zen, H.; Dang, V.; Clark, R.; Zhang, Y.; Weiss, R. J.; Jia, Y.; Chen, Z.; and Wu, Y. 2019. LibriTTS: A Corpus Derived from LibriSpeech for Text-to-Speech. arXiv:1904.02882.

Zhang, X.; Zhang, D.; Li, S.; Zhou, Y.; and Qiu, X. 2024. SpeechTokenizer: Unified Speech Tokenizer for Speech Large Language Models. arXiv:2308.16692.

# Appendix A: Multi-round Reconstruction results of different codecs



Figure 5: Multi-round Reconstruction results of DAC.



Figure 6: Multi-round Reconstruction results of EnCodec.



Figure 7: Multi-round Reconstruction results of Mimi.



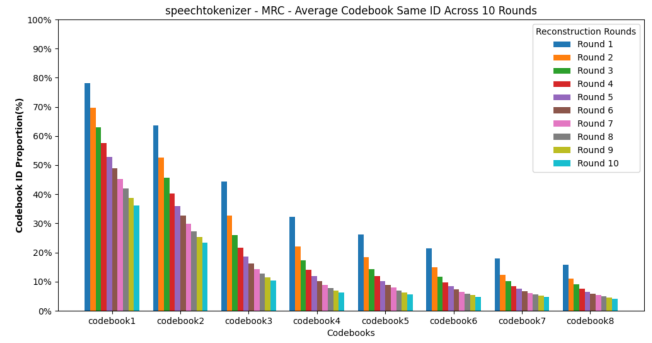Figure 8: Multi-round Reconstruction results of Semanti-Codec.



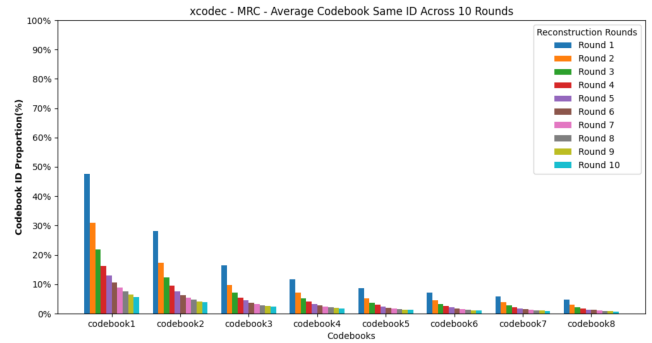Figure 9: Multi-round Reconstruction results of SpeechTok-enizer.



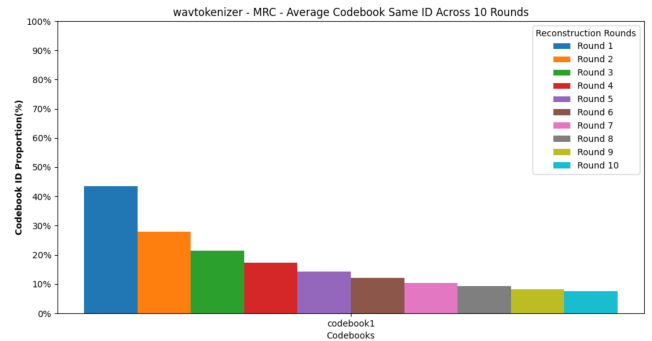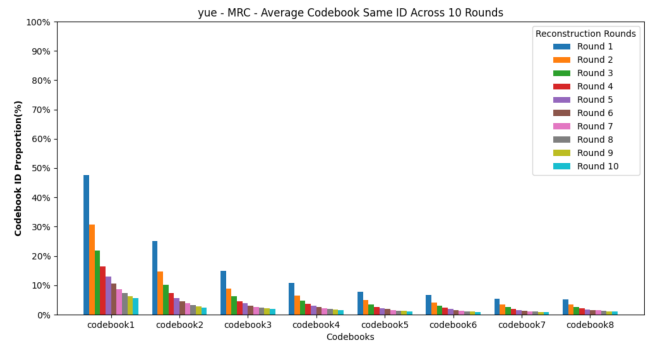Figure 10: Multi-round Reconstruction results of XCodec.



Figure 11: Multi-round Reconstruction results of WavTok-enizer.



Figure 12: Multi-round Reconstruction results of YuE.

# Appendix B: Audio Time Shift results of different codecs



Figure 13: Audio Time Shift results of DAC.



Figure 14: Audio Time Shift results of EnCodec.



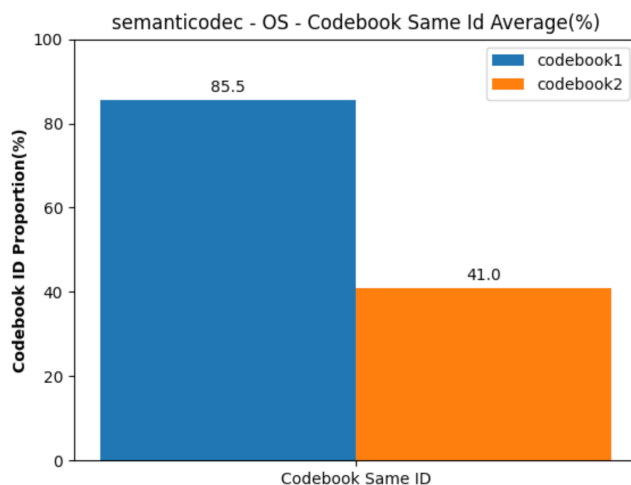Figure 15: Audio Time Shift results of Mimi.



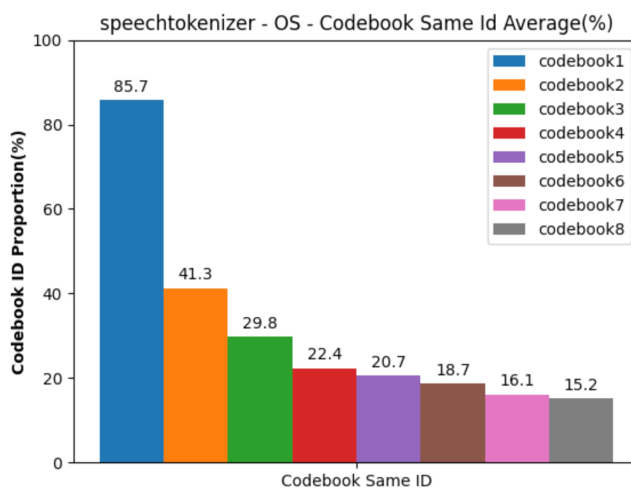Figure 16: Audio Time Shift results of SemantiCodec.



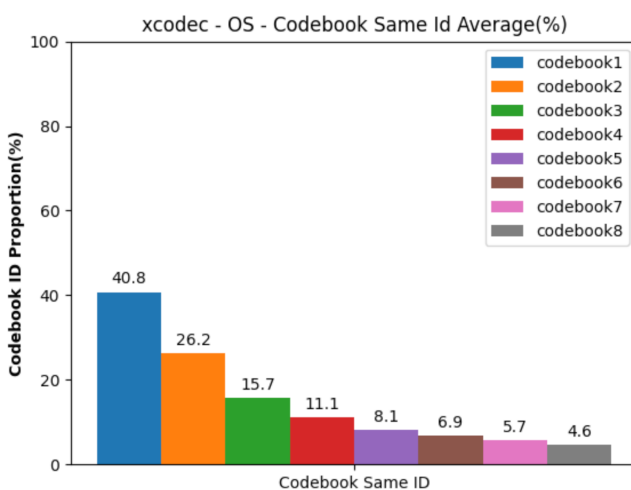Figure 17: Audio Time Shift results of SpeechTokenizer.

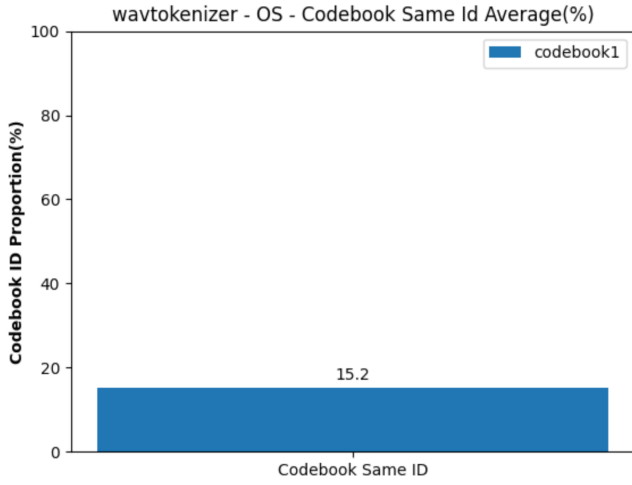

Figure 18: Audio Time Shift results of XCodec.

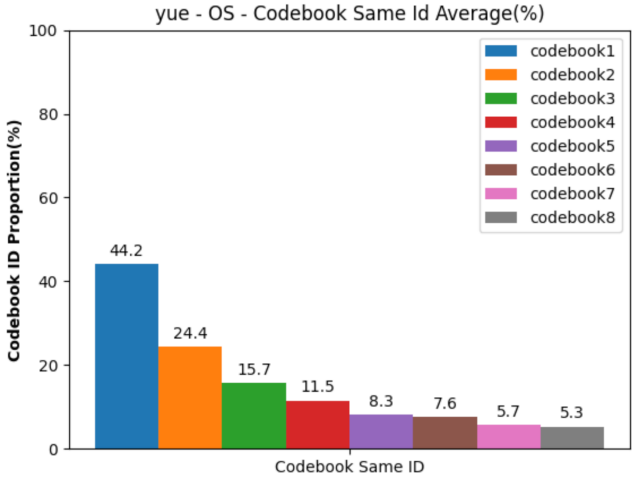Figure 19: Audio Time Shift results of WavTokenizer.



Figure 20: Audio Time Shift results of YuE.

## Appendix C: Visualization of music, speech and sound probe task results
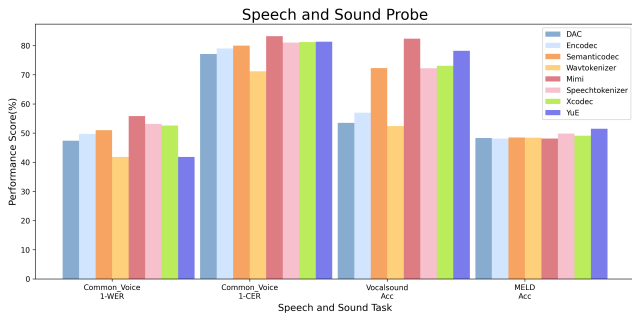


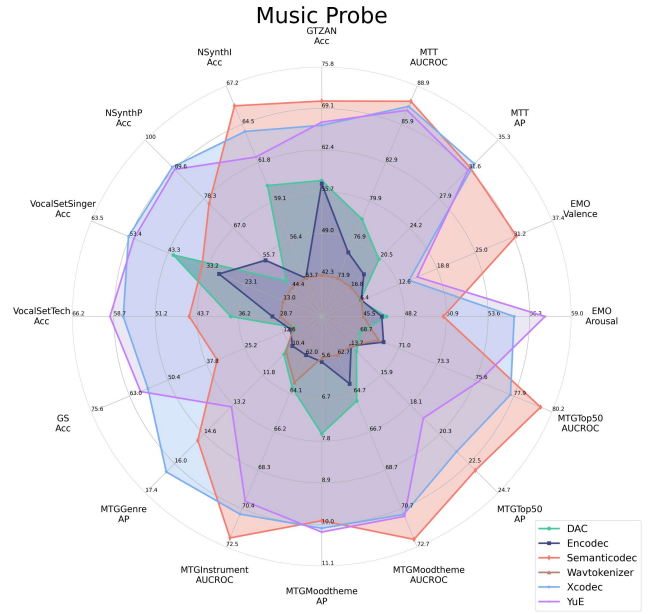Figure 21: Visualization for the speech and sound probe tasks.



Figure 22: Visualization for the music probe tasks.

## Appendix D: Introduction to downstream probe tasks and related datasets

We integrate a comprehensive dataset consisting of 17 sub-datasets from 12 audio collections (mostly derived from the MARBLE benchmark), covering major audio categories of speech, environmental sound, and music. Based on this dataset, we conduct 11 different types of probe tasks to examine the performance of different codec representations across different audio information dimensions, such as emotion, linguistic content, acoustic scene, and speaker identity.

**Genre Classification (GC)**: This task aims to classify music audio into predefined genres (e.g., rock, pop, classical). We use the GTZAN dataset and adopt Accuracy (Acc) as the performance metric. Additionally, we utilize MTG-Genre, a subset of MTG-Jamendo. Considering its longer track durations, we take the first 150 seconds of each track, segment them into 10-second clips, and stack them to serve as the input for the codec. This approach balances computational resources with the evaluation requirements. We use the Area Under the ROC Curve (ROC-AUC) and Average Precision (AP) to evaluate the representation's ability to encode genre information.

**Key Detection (KD)**: The goal of key detection is to predict the musical key of a piece of music, which is defined by its pitch center and mode (e.g., C major, a minor). We use the GiantSteps Key dataset, a collection of electronic dance music containing 24 major and minor keys. We consider the musical key as a global feature of the audio, processing it by stacking 10-second segments as the codec's input. We then use Acc as the evaluation metric to assess the model's ability to capture information about the musical structure.

**Emotion Detection (ED)**: This task focuses on identifying the emotional state or dimension conveyed by the audio (e.g., happiness, sadness, anger). We integrate several

datasets for this purpose: for the Valence and Arousal labels provided by the EmoMusic dataset, we employ a regression model for prediction and use the $R^2$ metric for evaluation. This helps assess the semantic information (high Valence) and acoustic information (high Arousal) embedded in the codec features. For the MTG MoodTheme dataset, which is a multi-label classification task with 59 emotion categories, we use ROC-AUC to evaluate the representation's ability to encode complex musical emotion information. Finally, using the MELD conversational speech dataset, we test the codec's capability to distinguish among seven basic emotions in a realistic context, which is evaluated with Acc.

**Vocal Technique Detection (VTD)**: This task aims to identify specific vocal techniques used by singers in musical compositions. It is a relatively uncommon, fine-grained identification task that focuses on the performance-level details. The main publicly available dataset is VocalSet, which contains recordings of 17 different vocal techniques performed by 20 professional singers, with each audio segment representing one technique category. We use Acc as the metric to evaluate the codec's ability to distinguish these subtle acoustic features.

**Pitch Classification (PC)**: This task aims to classify the main pitch content of a musical audio clip, with the range corresponding to MIDI note numbers 0 to 127 on the chromatic scale. We use the NSynth dataset, which consists of a large number of 4-second monophonic recordings. Due to its monophonic nature, this task can be viewed as a 128-class fine-grained pitch classification problem. It is designed to evaluate the accuracy of the codec's representation of fundamental frequency information, with performance assessed using Acc.

**Music Tagging (MT)**: This is a comprehensive evaluation task in the music domain that requires the model to assign multiple descriptive tags to music clips. These tags may cover various types, such as genre, instrument, and mood. We use the MagnaTagATune and MTG Top50 datasets. Following the MARBLE processing principles, we focus on our evaluation the model's ability to predict the Top 50 most frequent tags within these datasets. Given its multi-label nature, the final performance is measured by the ROC-AUC and the PR-AUC/AP. These metrics are used to evaluate the overall capability of the features in representing musical information.

**Instrument Classification (IC)**: This task aims to identify one or more musical instruments present in an audio recording. In the MARBLE classification system, this is considered an Acoustic-Level task, and its results evaluate the codec's ability to represent fundamental acoustic properties like timbre. For evaluation, we use the NSynth dataset, which contains 11 single-instrument categories and is evaluated using Acc. We also use the MTG Instrument dataset, a multi-label collection with 41 labels, which is evaluated using ROC-AUC and PR-AUC/AP.

**Automatic Speech Recognition (ASR)**: This task focuses on transforming speech signals from audio recordings into textual content. We use the Common Voice dataset, which contains approximately 26119 hours of recordings, including a variety of demographic metadata such as age,

gender, and accent. Among these, about 17127 hours of validated data cover 104 languages, with each language providing the necessary training, development, and test sets required to build a speech recognition model. Word Error Rate (WER) and Character Error Rate (CER) are used as the evaluation metrics.

**Singer Identification (SI)**: This task aims to identify the singer's identity from a short music recording. For this task, we use the VocalSet dataset, a collection containing audio from 20 different singers. We follow the MARBLE-recommended dataset partition (training:validation:test = 12:8:5), and ensure that all singer categories are evenly distributed. Finally, Acc is used to evaluate the codec's ability to distinguish individual vocal features.

**Vocals Sound Classification (VSC)**: This task aims to classify various non-linguistic sounds made by humans. We use the VocalSound dataset, which contains six common non-speech human sounds: laughter, sighs, coughs, throat clearing, sneezes, and sniffs. Since the audio clips in the dataset have non-uniform lengths, we pad all audio to a uniform length before inputting them into the codec. The evaluation for this task is conducted using Acc.

**Environmental Sound Classification (ESC)**: This task focuses on identifying sounds from the environment. We use the ESC-50 dataset, which is a labeled collection of 2000 environmental audio recordings consisting of 5-second-long recordings divided into 50 semantic categories. Since the original dataset does not provide an official standard split, we use a 9:1 ratio to self-partition it into a training set and a test set, with Acc as the metric for the evaluation of this dataset.