

SELF-GUIDANCE: TRAINING VQ-VAE DECODERS TO BE ROBUST TO QUANTIZATION ARTIFACTS FOR HIGH-FIDELITY NEURAL SPEECH CODEC

Anonymous authors

Paper under double-blind review

ABSTRACT

Neural speech codecs, predominantly based on Vector-Quantized Variational Autoencoders (VQ-VAEs), serve as fundamental audio tokenizers for speech large language models (SLLMs). However, their reconstruction fidelity is limited by quantization errors introduced during latent space discretization. Existing solutions typically increase model complexity through larger codebooks or hierarchical quantization, which subsequently intensify the modeling challenge for downstream SLLMs. Inspired by the key insight that the codec decoder produces superior output from continuous pre-quantize embeddings, we propose a novel self-guided training mechanism that addresses this problem by enhancing decoder robustness rather than modifying the quantization process. Our method introduces an additional training objective that aligns the decoder’s intermediate features when processing both quantized tokens and continuous pre-quantized embeddings through a feature-mapping loss. Extensive experiments on XCodec2 demonstrate that self-guidance consistently improves reconstruction quality across various codebook sizes and quantization techniques (FSQ, SimVQ, [multi-codebook VQ](#)), achieving state-of-the-art performance for low-bitrate speech codecs. The method requires minimal additional training cost and no inference-time modifications, offering an efficient solution for high-fidelity neural audio coding. Remarkably, our approach enables a 4× reduction in codebook size while maintaining comparable fidelity. Downstream text-to-speech experiments confirm that this reduction significantly improves LLM-based synthesis performance by simplifying the token modeling space.

1 INTRODUCTION

Audio codecs serve as essential tools for audio compression, originally designed to encode continuous audio signals like human speech into sequences of reconstructable discrete codes, enabling efficient data transmission and storage (Wu et al., 2024a). Recently, neural speech codecs, pioneered by SoundStream (Zeghidour et al., 2021) and EnCodec (Défossez et al., 2022), leverage the Vector-Quantized Variational AutoEncoder (VQ-VAE) (Van Den Oord et al., 2017; Esser et al., 2021) architectures to achieve high-fidelity reconstruction at compression ratios significantly exceeding traditional codecs. This breakthrough facilitates the integration of large language models (LLMs) in speech processing and generation, where the discretized audio tokens could be directly adopted in the standard next-token-prediction frameworks of LLMs. Benefiting from large-scale speech modeling with LLMs, numerous studies have advanced downstream tasks, including text-to-speech generation (Wang et al., 2023; Yang et al., 2023b) and interactive multimodal large language models (MLLMs) (Défossez et al., 2024; Zhan et al., 2024).

The transformation from continuous audio to discrete tokens in a VQ-VAE is enabled by a latent vector quantizer. This component maps continuous latent vectors from the encoder to entries in a finite codebook via nearest-neighbor search (i.e., vector quantization) (Van Den Oord et al., 2017; Yu et al., 2021; Mentzer et al., 2023). The corresponding codebook embeddings are then passed to the decoder to reconstruct the audio waveform.

054 However, quantization is inherently lossy. As noted in prior work (Liu et al., 2024) and confirmed by
055 our preliminary experiments (Section 3.2), the decoder produces higher-fidelity audio when using the
056 continuous, pre-quantized latents compared to the quantized tokens. This performance gap confirms
057 that quantization error constitutes a major obstacle to high-fidelity reconstruction, as it restricts the
058 information available to the decoder.

059 To suppress quantization error, existing neural codecs typically employ strategies such as hierarchical
060 quantization with multiple residual codebooks (Zeghidour et al., 2021; Yang et al., 2023a) or simply
061 scaling up the codebook size (Parker et al., 2024; Xin et al., 2024; Wu et al., 2024b; Ye et al.,
062 2025a). While effective for compression, these approaches introduce significant challenges for
063 downstream LLM modeling. Hierarchical codebooks require complex mechanisms to fit within
064 auto-regressive transformer frameworks (Wang et al., 2023; Yang et al., 2023b; Défossez et al., 2024).
065 On the other hand, unlike text tokenizers that use hierarchical subword units (e.g., BPE) Dubey et al.
066 (2024), a larger audio codebook expands a flat, unstructured vocabulary, exponentially increasing the
067 complexity of autoregressive sequence modeling (Ye et al., 2025b).

068 Thus, reducing quantization error often involves intricate codec designs that inadvertently transfer
069 complexity to downstream LLMs. In this paper, we shift the focus from modifying the quantizer or
070 latent space to enhancing the decoder itself. Our core idea is to **guide the decoder to narrow the
071 output gap between the pre-quantized latent vectors and the quantized tokens**. By aligning the
072 decoder’s outputs for these two inputs, we directly mitigate the artifacts introduced by quantization,
073 thereby relieving the quantizer of the sole burden of error elimination.

074 To this end, we propose a novel learning scheme for VQ-VAE-based codecs, which we call **self-
075 guidance**. During training, the decoder receives both the quantized token embeddings and the
076 continuous pre-quantized latent vectors. We then apply a feature mapping loss between the decoder’s
077 intermediate features or outputs for these two paths. This additional objective uses the high-fidelity
078 output from the pre-quantized latents as a target, guiding the decoder to produce similar, high-quality
079 features when driven by the quantized tokens. Consequently, the decoder becomes more robust to
080 quantization artifacts, enhancing the final reconstructed audio’s fidelity.

081 We implement our self-guidance approach on the state-of-the-art single-codebook neural speech
082 codec XCodec2 (Ye et al., 2025b), applying the feature mapping loss to the outputs of the decoder’s
083 transformer backbone. Experiments on LibriSpeech show consistent reconstruction improvements
084 across various codebook sizes and quantization techniques (e.g., FSQ, SimVQ). Notably, we achieve
085 comparable reconstruction quality with only a quarter of the original codebook size. The benefits of a
086 smaller codebook are further demonstrated in downstream text-to-speech LLM experiments. Audio
087 samples are available on our demo website.¹

088 Our main contributions are as follows:

- 090 1. We propose a novel self-guidance mechanism for VQ-VAEs that directs the decoder to
091 mitigate the detrimental effects of quantization error on reconstruction fidelity.
- 093 2. We apply self-guidance to the XCodec2 model, achieving state-of-the-art reconstruction
094 performance for low-bitrate speech codecs.
- 096 3. Through extensive experiments, we demonstrate that the improvements generalize across
097 different codebook sizes, vector quantization methods, and **codec model architectures**.
- 099 4. We provide statistical evidence confirming that self-guidance primarily regulates the decoder
100 rather than the encoder.
- 102 5. We show that self-guidance reduces the codec’s dependency on large codebooks, **yielding
103 a higher compression rate, as well as significant benefits for downstream LLM-based
104 applications**.

107 ¹<https://sgvqvae.github.io/sgvqvae-demo>

2 RELATED WORKS

2.1 VECTOR QUANTIZATION

VQ-VAE (Van Den Oord et al., 2017) introduced discrete latent representations for generative models, and VQ-VAE2 (Razavi et al., 2019) enhanced representation richness through hierarchical architectures. VQGAN (Esser et al., 2021) integrated adversarial networks, establishing a fundamental VQ framework for high-quality generative models such as Stable Diffusion (Rombach et al., 2022). Nevertheless, these methods encounter representation collapse when dealing with large codebook sizes, which restricts their scalability.

To tackle this issue, DALL-E (Ramesh et al., 2021) employs Gumbel-Softmax sampling to activate more codes during training, although only a small subset of codes is used for quantization during inference (Zhang et al., 2023). VQGAN-FC (Yu et al., 2021) mitigates collapse by reducing latent dimensionality and applying L2 normalization. Finite scalar quantization (FSQ) (Mentzer et al., 2023) and its variant Look-up free quantization (LFQ) (Yu et al., 2023) project latents to low-dimensional spaces (e.g., binary codes), but this comes at the cost of model capacity, as performance degrades when codebooks are small or collapse is not severe. Recently, VQGAN-LC (Zhu et al., 2024a) and SimVQ (Zhu et al., 2024b) enable stable training with codebook sizes up to 100k by incorporating a linear projector for the codebook.

2.2 NEURAL CODEC

In early neural codec model studies, SoundStream (Zeghidour et al., 2021) utilized residual vector quantizers (RVQs) to distribute the codec model’s total bitrate across multiple codebooks, preventing codebook size explosion. However, this hierarchical design complicates downstream applications due to the multiple tokens within each frame, necessitating additional flattening or joint modeling.

In recent years, single-codebook codecs have emerged as a simpler and more efficient alternative, demonstrating strong performance at low bitrates (Li et al., 2024; Guo et al., 2024; Ji et al., 2024; Xin et al., 2024; Della Libera et al., 2025). For instance, BigCodec (Xin et al., 2024) employs larger model sizes and advanced learning objectives to achieve high-fidelity audio decoding from a single quantizer of frame rate 80Hz. Despite these advancements, the reconstruction fidelity of BigCodec on perspective metrics significantly degrades at lower frame rates. While high frame rate incurs longer audio token sequences, resulting in a quadratic increase in downstream LLM computation cost, and the language modeling complexity (Wang et al., 2024).

To address this challenge, XCodec (Ye et al., 2025a) and FocalCodec (Della Libera et al., 2025) integrate pretrained self-supervised audio encoders to support the reconstruction performance on perspective metrics. Thanks to the stabilized quantizer like FSQ, TS3Codec (Wu et al., 2024b) and XCodec2 (Ye et al., 2025b) extend the codebook size to over 2^{16} to further boost the model performance, achieving state-of-the-art performance on single-codebook codec of frame rate around 50Hz. However, the drastically extended codebook size poses a significant challenge to the language modeling of downstream LLMs. This issue motivates the development of this paper to explore an approach that relieves the existing codec model’s dependency on the large codebook.

2.3 SELF-DISTILLATION

Self-distillation is a specific paradigm within knowledge distillation where the student and teacher are instances of the same model architecture, or even the same model itself Zhang et al. (2019); Furlanello et al. (2018). This is often achieved by using a model’s own outputs from a previous training iteration or a differently initialized copy as the teacher’s knowledge, leveraging consistency regularization to improve the student’s generalization and calibration Mobahi et al. (2020). Our method is conceptually related to self-distillation, which also uses feature-mapping losses for model guidance. However, self-guidance presents a distinct contribution by specifically targeting the decoder’s robustness to quantization error in VQ-VAEs—a previously under-explored bottleneck. Our core innovation lies in using the pre-quantized features as an internal guide to explicitly align the decoder’s manifolds, establishing a new training paradigm for speech codecs. Furthermore, unlike self-distillation, which

requires a pre-trained teacher, our approach is implemented within a single, end-to-end training process, offering greater practicality and efficiency.

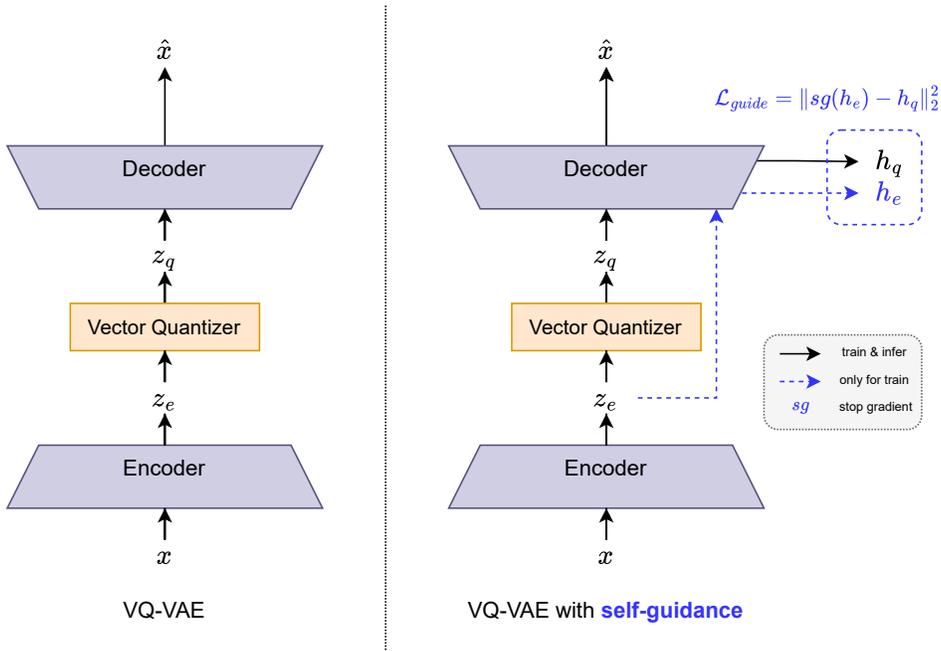


Figure 1: Illustration of the VQ-VAE architecture and the proposed self-guidance (SG) mechanism with the introduced feature mapping loss \mathcal{L}_{guide} .

3 PRELIMINARY: THE EFFECT OF QUANTIZATION IN NEURAL CODEC

3.1 REVISITING THE VQ-VAE FRAMEWORK

The Vector-Quantized Variational Autoencoder (VQ-VAE) framework forms the foundation of modern neural audio codecs. As illustrated in Figure 1 (left), the architecture consists of three main components: an encoder, a vector quantizer, and a decoder. The encoder processes an input audio signal x to produce a sequence of continuous latent embeddings $z_e \in \mathbb{R}^{d_e}$, where d_e is the latent dimension. The vector quantizer then maps each embedding in z_e to the nearest entry in a finite codebook $\mathcal{Q} \subset \mathbb{R}^{d_e}$. This operation produces a sequence of quantized token embeddings z_q . Finally, the decoder reconstructs the audio signal \hat{x} from z_q . During training, gradients are propagated through the non-differentiable quantization operation using straight-through estimation (STE), which copies gradients from z_q directly to z_e .

The quantization process inherently introduces error as it projects continuous latent vectors onto a discrete codebook. The quantization error can be quantified as:

$$e_q = \|z_e - z_q\|_2 \tag{1}$$

This error represents the information loss incurred during discretization.

3.2 OBSERVATION OF QUANTIZATION ARTIFACTS

The quantization error e_q introduces information loss that propagates to the decoder, resulting in reconstruction fidelity degradation, known as the quantization artifacts. This phenomenon is evident even though the decoder is exclusively trained on quantized inputs during standard VQ-VAE training.

As shown in Table 1, according to the findings from Liu et al. (2024) on the EnCodec model (Défossez et al., 2022), when the decoder processes the continuous pre-quantized latents z_e instead

Table 1: Comparing the reconstruction performance of neural speech codec models with different decoder inputs.

Codec model	bitrate	decoder input	STOI \uparrow	WER \downarrow	SIM \uparrow
<i>Ground Truth</i>			1.00	2.4	1.000
Encodec	6kbps	z_q	0.88	2.7	0.861
Encodec	6kbps	z_e	0.95	2.7	0.922
BigCodec	1.04kbps	z_q	0.93	3.6	0.841
BigCodec	1.04kbps	z_e	0.95	2.9	0.872

of the quantized tokens z_q , reconstruction quality improves significantly. This observation aligns with our evaluations of BigCodec (Xin et al., 2024).

These results demonstrate that quantization artifacts substantially limit reconstruction quality, presenting a major obstacle for achieving optimal performance in neural codecs. Thus, it is desirable to enhance the decoder’s robustness to the quantization error, enabling the generation of high-fidelity samples from the post-quantized latents despite the quantization error.

4 METHODOLOGY

To mitigate quantization artifacts in neural speech codecs, we propose a novel learning scheme called **self-guidance** (SG) for VQ-VAE decoders. This section details the self-guidance mechanism and explains our rationale for applying it to the XCodec2 model to construct a high-fidelity neural speech codec.

4.1 SELF-GUIDANCE MECHANISM

The self-guidance mechanism is designed to enhance the decoder’s ability to compensate for the information loss caused by quantization error in the input tokens z_q . Specifically, we aim to enable the decoder to produce similar outputs from both the quantized tokens z_q and the continuous pre-quantized latents z_e .

While the vanilla VQ-VAE reconstruction loss implicitly guides the decoder toward this objective by using the original input x as a target, our preliminary analysis indicates that this alone is insufficient to fully address quantization artifacts. This suggests the need for more explicit guidance during training.

Inspired by our preliminary findings, we propose using the pre-quantized latent z_e itself as an internal guidance signal. As illustrated in Figure 1 (right), during training we introduce an additional forward pass that feeds z_e to the decoder. We then extract intermediate hidden features from both paths: h_e from the z_e branch and h_q from the z_q branch. Specifically, the feature h is obtained as the output of the final Transformer block in XCodec2 decoder, which is then fed to an ISTFT head to reconstruct waveform. We introduce a feature-mapping loss $\mathcal{L}_{\text{guide}}$ to align these features:

$$\mathcal{L}_{\text{guide}} = \|\text{sg}(h_e) - h_q\|_2^2 \quad (2)$$

where $\text{sg}(\cdot)$ denotes the stop-gradient operation. This loss term is added to the original VQ-VAE objectives to form an end-to-end self-supervised training process.

The self-guidance mechanism introduces minimal computational and architectural overhead:

- **Training:** Only an additional forward pass through the decoder with z_e is required, with no gradient computation needed for this branch.
- **Inference:** No modifications are required; the decoder operates exclusively on z_q as in standard VQ-VAE models.

4.2 NEURAL SPEECH CODEC MODEL

To validate the effectiveness of self-guidance, we apply it to XCodec2, a state-of-the-art neural speech codec that has demonstrated strong performance in low-bitrate speech encoding and downstream speech generation tasks (Boson AI, 2025; Ye et al., 2025b).

XCodec2 comprises several key components: a convolutional encoder, a single-layer finite scalar quantizer (FSQ), and an acoustic decoder. Additionally, it includes a semantic encoder and decoder that form an auxiliary autoencoder operating on Wav2Vec2-BERT features (Barrault et al., 2023), enhancing the semantic content of the encoded latents for improved downstream performance.

A distinctive feature of XCodec2 is its acoustic decoder architecture. Like in TS3Codec (Wu et al., 2024b), rather than using stacked convolutional upsampling blocks, it employs a Transformer backbone followed by an inverse short-time Fourier transform (iSTFT) head (Siuzdak, 2024). This design naturally suggests using the Transformer backbone outputs for computing $\mathcal{L}_{\text{guide}}$ because: (i) the Transformer contains the majority of learnable parameters in the decoder, providing sufficient capacity to benefit from self-guidance; and (ii) the subsequent iSTFT head separates the hidden features from the final waveform generation, preventing potential interference from waveform-level reconstruction losses.

The complete training objective for our enhanced codec is:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{guide}} + \mathcal{L}_{\text{semantic}} + \mathcal{L}_{\text{acoustic}} + \mathcal{L}_{\text{adv}} \quad (3)$$

where:

- $\mathcal{L}_{\text{guide}}$ is the self-guidance feature mapping loss defined in Equation 2, computed using the Transformer backbone outputs;
- $\mathcal{L}_{\text{semantic}}$ is the mean squared error semantic feature reconstruction loss;
- $\mathcal{L}_{\text{acoustic}}$ is the multi-scale mel-spectrogram reconstruction loss;
- \mathcal{L}_{adv} is the adversarial loss from a multi-period discriminator (Kong et al., 2020) and a spectrogram discriminator (Parker et al., 2024).

5 EXPERIMENTS AND ANALYSIS

5.1 EXPERIMENT SETTINGS

Dataset We use the full Librispeech Panayotov et al. (2015) training set for the training of all versions of codec models, which comprises 960 hours of English speech audio at a sampling rate of 16kHz. For evaluation, the *test-clean* subset of LibriSpeech that contains 2620 utterances from 40 speakers is used to assess reconstruction performance.

Implementation details We build our neural codec model based on the official open-source code of XCodec2². The modifications required to implement self-guidance are minimal: (i) adding an additional forward pass in the `forward` function³; (ii) incorporating the computation of $\mathcal{L}_{\text{guide}}$ in the `compute_gen_loss` function⁴. The full modified code script is attached in the supplementary material. Detailed configurations are included in Section A.2. The BigCodec model involved in the preliminary study (Section 3.2) and comparative experiments (Section 5.2) is obtained via training with the official open-source implementation⁵.

Training cost We train all of the codec models on 8 NVIDIA GeForce RTX 4090 GPUs for 600 thousand iterations. The total training time of each codec model is around 237.75 hours. Notably, the self-guidance variant incurs negligible additional training time compared to the baseline

²<https://github.com/zhenye234/X-Codec-2.0>

³https://github.com/zhenye234/X-Codec-2.0/blob/main/lightning_module.py#L146

⁴https://github.com/zhenye234/X-Codec-2.0/blob/main/lightning_module.py#L239

⁵<https://github.com/Aria-K-Alethia/BigCodec>

XCodec2, with differences of only seconds. This efficiency aligns with our design in Section 4.1: the additional forward pass through the acoustic decoder requires no backward propagation, making the computational overhead minimal compared to other components (e.g., discriminators) and gradient synchronization. This demonstrates that the performance gains from self-guidance come at virtually no additional training cost.

Table 2: Comparing reconstruction evaluation results with other existing neural codecs on the LibriSpeech test-clean dataset. (**SG** signifies the proposed self-guidance mechanism; details about each metric are included in Section A.1)

Codecs models	Frame rate	Codebook size(s)	PESQ \uparrow	STOI \uparrow	MCD \downarrow	WER \downarrow	SIM \uparrow	UTMOS \uparrow
<i>Ground Truth</i>			4.64	1.000	0.00	2.5	1.00	4.08
DAC	50Hz	1024 \times 8	2.72	0.940	–	–	0.87	–
DAC	50Hz	1024 \times 2	1.13	0.730	–	–	0.32	–
WavTokenizer	75Hz	4096	2.05	0.886	4.00	6.8	0.59	3.89
BigCodec	80Hz	8192	2.68	0.935	2.93	3.6	0.84	4.11
WavTokenizer	40Hz	4096	1.88	0.868	4.32	8.0	0.57	3.77
BigCodec	40Hz	8192	<u>2.11</u>	<u>0.894</u>	<u>3.72</u>	6.7	0.66	4.05
XCodec2	50Hz	8192	2.03	0.892	3.84	<u>4.1</u>	<u>0.72</u>	4.09
XCodec2+SG	50Hz	8192	2.13	0.898	3.60	3.8	0.73	<u>4.08</u>
TS3Codec	40Hz	65536	2.01	0.893	3.81	4.9	0.61	3.69
TS3Codec	40Hz	131072	2.06	0.897	3.75	4.5	0.63	3.73
TS3Codec	50Hz	65536	2.22	0.909	3.52	3.6	0.68	3.85
TS3Codec	50Hz	131072	2.23	0.910	3.50	3.6	0.68	3.84
XCodec2	50Hz	65536	2.28	0.910	3.57	3.2	0.79	4.06
XCodec2+SG	50Hz	65536	2.39	0.915	3.41	3.2	0.80	4.10

5.2 RECONSTRUCTION PERFORMANCE

We first evaluate the overall reconstruction performance of our proposed model against existing low-bitrate speech codecs. For DAC, WavTokenizer, and TS3Codec, we report results from papers. For BigCodec and XCodec2, we retrain models and include variations with different configurations (XCodec2 default: frame rate = 50 Hz, $|\mathcal{Q}| = 65,536$; BigCodec default: frame rate = 80 Hz, $|\mathcal{Q}| = 8,192$).

As shown in Table 2, our proposed model (XCodec2 with self-guidance) achieves the best performance across most evaluation metrics. For codecs with frame rates of 40–50 Hz, our approach consistently outperforms competitors with similar codebook sizes (8,192 and below, or 65,536 and above), establishing new state-of-the-art performance for low-bitrate speech codecs. **Considering that the BigCodec (159M params) improves PESQ by 0.15 over TFCodec Jiang et al. (2023) (6.37M params) with a 25 \times parameter increase, our method achieves a PESQ improvement of 0.1 over strong baselines with minimal cost (zero inference overhead, no architectural changes).**

Specifically, while the original XCodec2 with codebook size 65,536 shows competitive performance, self-guidance provides further improvements across all metrics. Reducing XCodec2’s codebook size to 8,192 significantly degrades acoustic reconstruction quality (PESQ, STOI, MCD), falling behind BigCodec. However, when augmented with self-guidance, this reduced-size model surpasses BigCodec on all metrics.

5.3 ABLATION STUDIES

We conduct ablation studies to isolate the contribution of self-guidance and evaluate its robustness under different configurations. We compare models trained with and without self-guidance while varying quantizer settings.

Codebook size We experiment with codebook sizes of 8,192, 16,384, and 65,536. As shown in Table 3, self-guidance improves performance across all settings, except for a minor degradation in

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

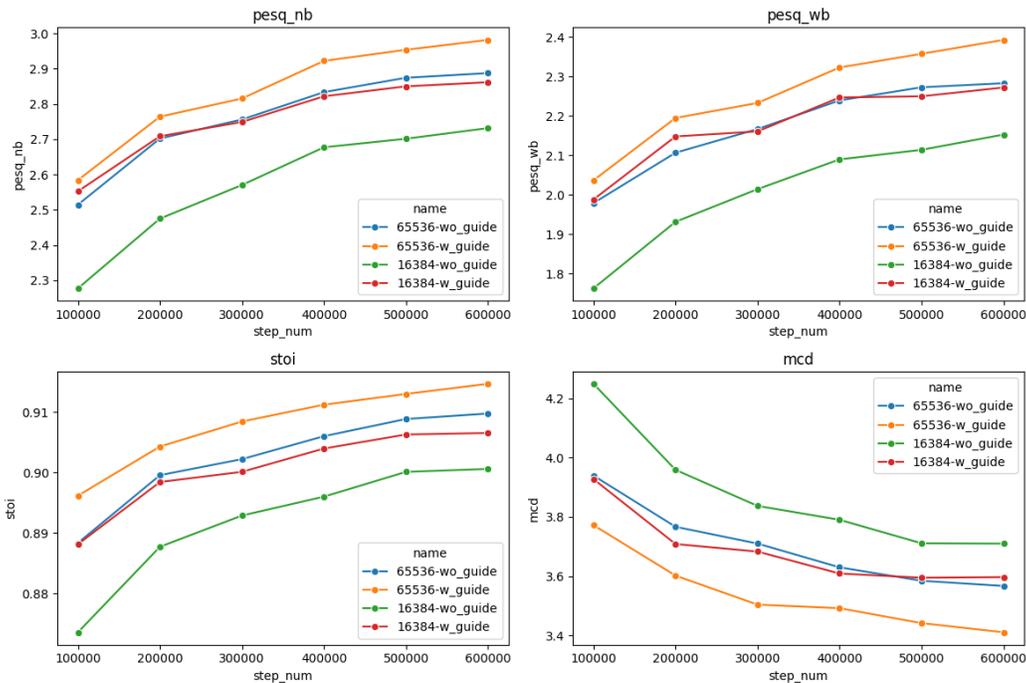


Figure 2: Comparison of the reconstruction performance under various settings along the training process. Horizontal axis is the training iterations. Best viewed in color.

Table 3: Reconstruction evaluation results of the proposed neural speech codec across different codebook sizes. (with SG signifies whether the proposed self-guidance mechanism is applied)

Codebook size	with SG	PESQ-WB↑	PESQ-NB↑	STOI↑	MCD↓	WER↓	SIM↑	UTMOS↑
Ground Truth		4.64	4.54	1.000	0.00	2.49	1.00	4.08
8192	✗	2.03	2.59	0.892	3.84	4.08	0.72	4.09
8192	✓	2.13	2.69	0.898	3.79	3.77	0.73	4.08
16384	✗	2.15	2.73	0.901	3.71	3.47	0.76	3.98
16384	✓	2.27	2.86	0.907	3.70	3.53	0.77	4.08
65536	✗	2.28	2.89	0.910	3.57	3.23	0.79	4.06
65536	✓	2.39	2.98	0.915	3.41	3.15	0.80	4.10

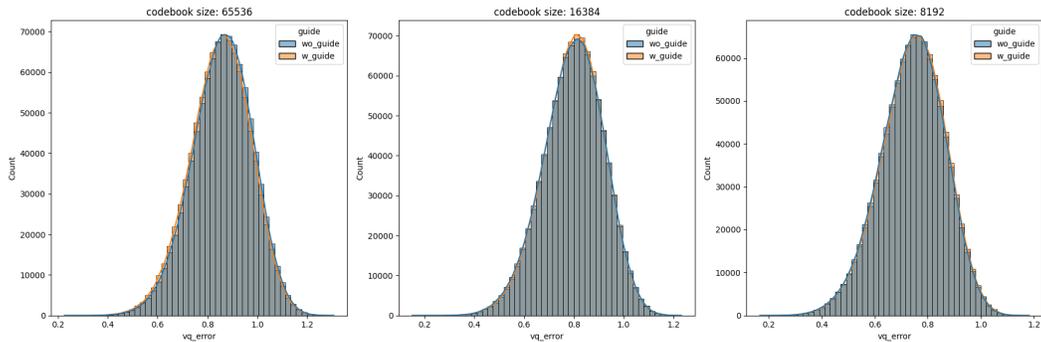
word error rate (WER) at the intermediate size (16,384). Figure 2 shows that with self-guidance, a model with codebook size 16,384 achieves similar performance to the baseline XCodec2 with a 4x larger codebook (65,536) on several metrics.

Type of vector quantizer To assess generalization across quantizer types, we replace the default FSQ quantizer in XCodec2 with SimVQ (VQGAN-FC suffered from codebook collapse and produced unintelligible results). Table 4 shows that self-guidance consistently improves performance with SimVQ, reproducing the minor WER degradation observed with FSQ. Since we only observe slight WER rise of 0.06% at codebook size 16384, there appears to be no systematic trend for sacrificing intelligibility, especially given the consistent gains in the spectral intelligibility metric STOI.

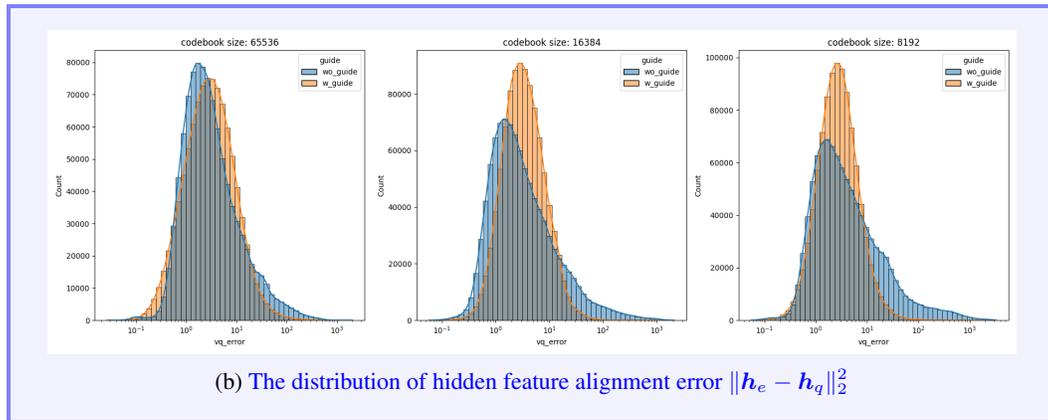
To further demonstrate the generalizability of self-guidance to multi-codebook VQ and different decoder networks, we provide extra experiment results on XCodec2 using ResidualFSQ and the BigCodec codec model in Section A.3 and A.4, which report consistent performance gains.

Table 4: Reconstruction evaluation results of the proposed neural speech codec across different types of vector quantizers (XCodec2 adopts FSQ by default), with codebook size fixed at 16384. (**with SG** signifies whether proposed self-guidance mechanism is applied).

Quantizer	with SG	PESQ-WB \uparrow	PESQ-NB \uparrow	STOI \uparrow	MCD \downarrow	WER \downarrow	SIM \uparrow	UTMOS \uparrow
<i>Ground Truth</i>		4.64	4.54	1.000	0.00	2.49	1.00	4.08
FSQ	\times	2.15	2.73	0.901	3.71	3.47	0.76	3.98
FSQ	\checkmark	2.27	2.86	0.907	3.60	3.53	0.77	4.08
SimVQ	\times	2.10	2.67	0.900	3.63	3.59	0.75	3.85
SimVQ	\checkmark	2.17	2.74	0.904	3.56	3.63	0.76	3.93



(a) The distribution of quantization error $\|z_e - z_q\|_2^2$



(b) The distribution of hidden feature alignment error $\|h_e - h_q\|_2^2$

Figure 3: The histogram of the quantization error e_q and hidden feature alignment MSE on LibriSpeech test-clean dataset with the self-guidance mechanism activated (`w_guide`) or omitted (`wo_guide`) across different codebook sizes (from left to right: 65536, 16384, 8192).

Quantization error Since $\mathcal{L}_{\text{guide}}$ gradients propagate to the encoder via straight-through estimation, we investigate whether performance gains stem from decoder guidance or implicit quantization error reduction. Figure 3a shows quantization error (e_q) distributions on the test-clean dataset for baseline and self-guidance models across different codebook sizes. The closely overlapping distributions demonstrate that **self-guidance enhances reconstruction by improving decoder robustness rather than reducing quantization error.**

Hidden feature alignment We also plot the distributions of the hidden feature alignment error $\|h_e - h_q\|_2^2$ in Figure 3b, which shows obvious divergence between the baseline and proposed

approach. The proposed self-guidance significantly suppresses the error, indicating an effective alignment between the generation results from the pre- and post-quantize latents. This evidence verifies the core motivation of self-guidance to enhance the VQ-VAE decoder’s robustness to quantization error via latent manifold alignment, leading to higher reconstruction fidelity. Detailed statistics of both figures are provided in Section A.5.

5.4 DOWNSTREAM AUTO-REGRESSIVE TTS

Building on our finding that self-guidance enables smaller codebooks to achieve performance comparable to larger ones (Figure 2), we evaluate its impact on downstream autoregressive text-to-speech (TTS) synthesis. We hypothesize that reduced codebook size simplifies the language modeling task, potentially improving final TTS quality.

We conduct rapid TTS experiments using a Qwen2.5-0.5B causal LLM backbone (Qwen et al., 2025) trained on LibriTTS-R (Koizumi et al., 2023). Input text is phonemized before training. Models are supervised fine-tuned for phoneme-to-audio-token generation for 85 epochs. For inference, we use the continual synthesis approach from VALL-E (Wang et al., 2023), providing phoneme sequences and first 3-second audio tokens as prompts for continuation generation. We filter LibriTTS-R test-clean samples longer than 6 seconds and generate continuations using top-k (50) and top-p (0.9) sampling.

As shown in Table 5, the TTS performance strongly correlates with codebook size, where models using a 16,384 codebook significantly outperform those with a 65,536 codebook. This aligns with our hypothesis in Section 1 that a large flat audio token vocabulary is harmful to the downstream LLMs.

Within this context, self-guidance yields the best performance at the 16,384 size. At the 65,536 size, the results of self-guidance are mixed, as the TTS model is primarily hampered by the fundamental language modeling difficulty of the large codebook, which overshadows the fidelity improvement from self-guidance.

Table 5: Downstream text-to-speech continuation performance on the LibriTTS test-clean dataset.

Codec model	Codebook size	UTMOS \uparrow	WER \downarrow	SIM \uparrow
XCodec2	65536	3.33	33.03	0.58
XCodec2+SG	65536	3.39	35.07	0.58
XCodec2	16384	3.51	28.78	0.56
XCodec2+SG	16384	3.58	28.02	0.58

6 CONCLUSION

We proposed self-guidance, a novel training mechanism for VQ-VAE-based neural speech codecs that enhances decoder robustness to quantization artifacts. By aligning the decoder’s outputs for quantized and continuous latent representations through an additional feature-mapping loss, our method improves reconstruction fidelity without modifying the inference process. Experiments demonstrate that self-guidance consistently enhances performance across various codebook sizes and quantization techniques, enabling comparable quality with 4 \times smaller codebooks. Downstream TTS results confirm that this reduction simplifies language modeling for LLMs, improving synthesis quality. Our approach provides an effective and efficient solution to mitigate quantization errors, advancing high-fidelity neural audio compression.

Future work could explore applying self-guidance to other VQ-VAE domains beyond speech, such as music or general audio processing, and investigating its combination with more advanced quantization techniques.

Ethics statement This work presents research on neural audio codecs, which have significant potential for positive applications in speech compression, communication, and generative modeling. However, we acknowledge several ethical considerations:

1. **Positive Impacts:** Our method enables higher-quality audio compression at lower bitrates, which could improve accessibility and efficiency in telecommunication, hearing assistance

540 devices, and low-bandwidth applications. The reduction in codebook complexity may also
541 decrease computational requirements for downstream applications.

- 542 2. **Potential Misuse:** Like other audio generation technologies, neural codecs could potentially
543 be misused for creating deepfake audio or other deceptive content. However, our work
544 focuses specifically on reconstruction quality rather than generative capabilities. The codec
545 itself does not generate novel content without being integrated into a full generative system.
546

547 **Reproducibility statement** As stated in Section 5.1, the proposed approach is adapted from the
548 official open-source code of XCodec2. The modified code script that implements self-guidance is
549 attached in the supplementary material. We have also described the computational requirements of
550 our experiments in Section 5.1. Details about the model configuration are included in Section A.2 for
551 confirmation.

552 REFERENCES

- 553
554
555 Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Mark Duppenhaler,
556 Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim, et al. Seamless: Multilingual
557 expressive and streaming speech translation. *arXiv preprint arXiv:2312.05187*, 2023.
- 558 Boson AI. Higgs Audio V2: Redefining Expressiveness in Audio Generation. <https://github.com/boson-ai/higgs-audio>, 2025. GitHub repository. Release blog available at <https://www.boson.ai/blog/higgs-audio-v2>.
559
560
- 561 Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki
562 Kanda, Takuya Yoshioka, Xiong Xiao, et al. Wavlm: Large-scale self-supervised pre-training
563 for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):
564 1505–1518, 2022.
- 565
566 Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. High fidelity neural audio
567 compression. *arXiv preprint arXiv:2210.13438*, 2022.
- 568
569 Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou,
570 Edouard Grave, and Neil Zeghidour. Moshi: a speech-text foundation model for real-time dialogue.
571 *arXiv preprint arXiv:2410.00037*, 2024.
- 572 Luca Della Libera, Francesco Paissan, Cem Subakan, and Mirco Ravanelli. Focalcodec: Low-bitrate
573 speech coding via focal modulation networks. *arXiv preprint arXiv:2502.04465*, 2025.
- 574
575 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
576 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
577 *arXiv e-prints*, pp. arXiv-2407, 2024.
- 578 Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image
579 synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
580 pp. 12873–12883, 2021.
- 581
582 Tommaso Furlanello, Zachary Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar.
583 Born again neural networks. In *International conference on machine learning*, pp. 1607–1616.
584 PMLR, 2018.
- 585 Yiwei Guo, Zhihan Li, Chenpeng Du, Hankun Wang, Xie Chen, and Kai Yu. Lscodex: Low-bitrate
586 and speaker-decoupled discrete speech codec. *arXiv preprint arXiv:2410.15764*, 2024.
- 587
588 Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhota, Ruslan Salakhutdinov,
589 and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked
590 prediction of hidden units. *IEEE/ACM transactions on audio, speech, and language processing*,
591 29:3451–3460, 2021.
- 592 Shengpeng Ji, Ziyue Jiang, Wen Wang, Yifu Chen, Minghui Fang, Jialong Zuo, Qian Yang, Xize
593 Cheng, Zehan Wang, Ruiqi Li, et al. Wavtokenizer: an efficient acoustic discrete codec tokenizer
for audio language modeling. *arXiv preprint arXiv:2408.16532*, 2024.

- 594 Xue Jiang, Xiulian Peng, Huaying Xue, Yuan Zhang, and Yan Lu. Latent-domain predictive neural
595 speech coding. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:2111–
596 2123, 2023.
- 597 Yuma Koizumi, Heiga Zen, Shigeki Karita, Yifan Ding, Kohei Yatabe, Nobuyuki Morioka, Michiel
598 Bacchiani, Yu Zhang, Wei Han, and Ankur Bapna. Libritts-r: A restored multi-speaker text-to-
599 speech corpus. *arXiv preprint arXiv:2305.18802*, 2023.
- 600 Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. Hifi-gan: Generative adversarial networks for
601 efficient and high fidelity speech synthesis. *Advances in neural information processing systems*,
602 33:17022–17033, 2020.
- 603 Hanzhao Li, Liumeng Xue, Haohan Guo, Xinfu Zhu, Yuanjun Lv, Lei Xie, Yunlin Chen, Hao Yin,
604 and Zhifei Li. Single-codec: Single-codebook speech codec towards high-performance speech
605 generation. *arXiv preprint arXiv:2406.07422*, 2024.
- 606 Alexander H. Liu, Qirui Wang, Yuan Gong, and James R. Glass. Closer look at neural codec
607 resynthesis: Bridging the gap between codec and waveform generation. In *Audio Imagination:
608 NeurIPS 2024 Workshop AI-Driven Speech, Music, and Sound Generation*, 2024. URL <https://openreview.net/forum?id=eO4DmksTv8>.
- 609 Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschannen. Finite scalar quantization:
610 Vq-vae made simple. *arXiv preprint arXiv:2309.15505*, 2023.
- 611 Hossein Mobahi, Mehrdad Farajtabar, and Peter Bartlett. Self-distillation amplifies regularization in
612 hilbert space. *Advances in Neural Information Processing Systems*, 33:3351–3361, 2020.
- 613 Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: an asr corpus
614 based on public domain audio books. In *2015 IEEE international conference on acoustics, speech
615 and signal processing (ICASSP)*, pp. 5206–5210. IEEE, 2015.
- 616 Julian D Parker, Anton Smirnov, Jordi Pons, CJ Carr, Zack Zukowski, Zach Evans, and Xubo Liu.
617 Scaling transformers for low-bitrate high-quality speech coding. *arXiv preprint arXiv:2411.19842*,
618 2024.
- 619 Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan
620 Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang,
621 Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin
622 Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi
623 Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan,
624 Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL
625 <https://arxiv.org/abs/2412.15115>.
- 626 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen,
627 and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine
628 learning*, pp. 8821–8831. Pmlr, 2021.
- 629 Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with
630 vq-vae-2. *Advances in neural information processing systems*, 32, 2019.
- 631 Antony W Rix, John G Beerends, Michael P Hollier, and Andries P Hekstra. Perceptual evaluation
632 of speech quality (pesq)-a new method for speech quality assessment of telephone networks and
633 codecs. In *2001 IEEE international conference on acoustics, speech, and signal processing.
634 Proceedings (Cat. No. 01CH37221)*, volume 2, pp. 749–752. IEEE, 2001.
- 635 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
636 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-
637 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 638 Takaaki Saeki, Detai Xin, Wataru Nakata, Tomoki Koriyama, Shinnosuke Takamichi, and Hi-
639 roshi Saruwatari. Utmos: Utokyo-sarulab system for voicemos challenge 2022. *arXiv preprint
640 arXiv:2204.02152*, 2022.
- 641
- 642
- 643
- 644
- 645
- 646
- 647

- 648 Hubert Siuzdak. Vocos: Closing the gap between time-domain and fourier-based neural vocoders for
649 high-quality audio synthesis. In *The Twelfth International Conference on Learning Representations*,
650 2024. URL <https://openreview.net/forum?id=vY9nzQmQBw>.
651
- 652 Cees H Taal, Richard C Hendriks, Richard Heusdens, and Jesper Jensen. An algorithm for intelli-
653 gibility prediction of time-frequency weighted noisy speech. *IEEE Transactions on audio, speech,*
654 *and language processing*, 19(7):2125–2136, 2011.
- 655 Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in*
656 *neural information processing systems*, 30, 2017.
657
- 658 Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing
659 Liu, Huaming Wang, Jinyu Li, et al. Neural codec language models are zero-shot text to speech
660 synthesizers. *arXiv preprint arXiv:2301.02111*, 2023.
- 661 Hankun Wang, Haoran Wang, Yiwei Guo, Zhihan Li, Chenpeng Du, Xie Chen, and Kai Yu. Why
662 do speech language models fail to generate semantically coherent outputs? a modality evolving
663 perspective, 2024. URL <https://arxiv.org/abs/2412.17048>.
664
- 665 Haibin Wu, Ho-Lam Chung, Yi-Cheng Lin, Yuan-Kuei Wu, Xuanjun Chen, Yu-Chi Pai, Hsiu-Hsuan
666 Wang, Kai-Wei Chang, Alexander H Liu, and Hung-yi Lee. Codec-superb: An in-depth analysis
667 of sound codec models. *arXiv preprint arXiv:2402.13071*, 2024a.
- 668 Haibin Wu, Naoyuki Kanda, Sefik Emre Eskimez, and Jinyu Li. Ts3-codec: Transformer-based
669 simple streaming single codec. *arXiv preprint arXiv:2411.18803*, 2024b.
670
- 671 Detai Xin, Xu Tan, Shinnosuke Takamichi, and Hiroshi Saruwatari. Bigcodec: Pushing the limits of
672 low-bitrate neural speech codec. *arXiv preprint arXiv:2409.05377*, 2024.
673
- 674 Dongchao Yang, Songxiang Liu, Rongjie Huang, Jinchuan Tian, Chao Weng, and Yuexian Zou. Hifi-
675 codec: Group-residual vector quantization for high fidelity audio codec. *CoRR*, abs/2305.02765,
676 2023a. URL <https://doi.org/10.48550/arXiv.2305.02765>.
677
- 678 Dongchao Yang, Jinchuan Tian, Xu Tan, Rongjie Huang, Songxiang Liu, Xuankai Chang, Jiatong
679 Shi, Sheng Zhao, Jiang Bian, Xixin Wu, et al. Uniaudio: An audio foundation model toward
680 universal audio generation. *arXiv preprint arXiv:2310.00704*, 2023b.
- 681 Zhen Ye, Peiwen Sun, Jiahe Lei, Hongzhan Lin, Xu Tan, Zheqi Dai, Qiuqiang Kong, Jianyi Chen,
682 Jiahao Pan, Qifeng Liu, et al. Codec does matter: Exploring the semantic shortcoming of codec
683 for audio language model. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
684 volume 39, pp. 25697–25705, 2025a.
- 685 Zhen Ye, Xinfu Zhu, Chi-Min Chan, Xinsheng Wang, Xu Tan, Jiahe Lei, Yi Peng, Haohe Liu, Yizhu
686 Jin, Zheqi Dai, et al. Llasa: Scaling train-time and inference-time compute for llama-based speech
687 synthesis. *arXiv preprint arXiv:2502.04128*, 2025b.
688
- 689 Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong
690 Xu, Jason Baldridge, and Yonghui Wu. Vector-quantized image modeling with improved vqgan.
691 *arXiv preprint arXiv:2110.04627*, 2021.
692
- 693 Lijun Yu, José Lezama, Nitesh B Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong
694 Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, et al. Language model beats diffusion-
695 tokenizer is key to visual generation. *arXiv preprint arXiv:2310.05737*, 2023.
- 696 Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. Sound-
697 stream: An end-to-end neural audio codec. *IEEE/ACM Transactions on Audio, Speech, and*
698 *Language Processing*, 30:495–507, 2021.
699
- 700 Jun Zhan, Junqi Dai, Jiasheng Ye, Yunhua Zhou, Dong Zhang, Zhigeng Liu, Xin Zhang, Ruibin Yuan,
701 Ge Zhang, Linyang Li, et al. Anygpt: Unified multimodal llm with discrete sequence modeling.
arXiv preprint arXiv:2402.12226, 2024.

Jiahui Zhang, Fangneng Zhan, Christian Theobalt, and Shijian Lu. Regularized vector quantization for tokenized image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18467–18476, 2023.

Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 3713–3722, 2019.

Lei Zhu, Fangyun Wei, Yanye Lu, and Dong Chen. Scaling the codebook size of vqgan to 100,000 with a utilization rate of 99%. *arXiv preprint arXiv:2406.11837*, 2024a.

Yongxin Zhu, Bocheng Li, Yifei Xin, and Linli Xu. Addressing representation collapse in vector quantized models with one linear layer. *arXiv preprint arXiv:2411.02038*, 2024b.

A APPENDIX

A.1 EVALUATION METRICS OF RECONSTRUCTION

We evaluate acoustic fidelity, intelligibility, and naturalness of the speech audio reconstructed by neural codecs using the following metrics:

Perceptual Evaluation of Speech Quality (PESQ) PESQ Rix et al. (2001) compares degraded and reference speech to predict human-perceived quality. We use a Python implementation⁶ to compute wide-band (PESQ-WB) and narrow-band (PESQ-NB) scores, where higher values indicate better quality.

Mel Cepstral Distortion (MCD) MCD measures the difference between mel-frequency cepstral coefficients (MFCCs), a standard metric for speech synthesis quality.

Short-Time Objective Intelligibility (STOI) STOI Taal et al. (2011) evaluates speech intelligibility by comparing temporal envelopes of clean and degraded signals, with scores ranging from 0 (unintelligible) to 1 (perfect intelligibility).

Word Error Rate (WER) WER is calculated using a HuBERT Hsu et al. (2021) speech recognition model finetuned on Librispeech⁷, reporting percentage errors in transcribed words.

Speaker Similarity (SIM) Speaker characteristics are evaluated via cosine similarity between original and reconstructed utterances, using a WavLM-large Chen et al. (2022)-based speaker verification model⁸.

UTMOS UTMOS Saeki et al. (2022) predicts Mean Opinion Score (MOS) for speech naturalness, with scores from 1 (poor) to 5 (excellent). We use a pretrained UTMOS strong model⁹.

A.2 MODEL CONFIGURATION

The detailed model configuration and loss weights are listed in Table 6. Most of the configurations follows the default configuration of XCodec2. Specifically, the weight of the proposed self-guidance feature mapping loss weight λ_{guide} is selected from the best of [1, 5, 10, 15], according to the overall reconstruction performance in test trials.

⁶<https://github.com/ludlows/PESQ>

⁷<https://huggingface.co/facebook/hubert-large-ls960-ft>

⁸https://github.com/microsoft/UniSpeech/tree/main/downstreams/speaker_verification

⁹<https://github.com/tarepan/SpeechMOS>

Table 6: Model configurations

Configuration entry	Value
Acoustic encoder hidden dim	1024
Acoustic encoder convoution blocks	5
Acoustic encoder up ratio	[2, 2, 4, 4, 5]
Acoustic decoder hidden dim	1024
Acoustic decoder Transformer layers	12
Semantic encoder hidden dim	1024
Semantic decoder hidden dim	1024
FSQ scales (codebook size = 65536)	[4, 4, 4, 4, 4, 4, 4, 4]
FSQ scales (codebook size = 16384)	[4, 4, 4, 4, 4, 4, 4]
FSQ scales (codebook size = 8192)	[4, 4, 4, 4, 4, 4, 2]
loss weight $\lambda_{semantic}$	5.0
loss weight $\lambda_{acoustic}$	15.0
loss weight λ_{adv}	1.0
loss weight λ_{guide} (codebook size = 65536)	5.0
loss weight λ_{guide} (codebook size = 16384)	10.0
loss weight λ_{guide} (codebook size = 8192)	10.0
batch size	16
optimizer	AdamW
optimizer betas	[0.8, 0.9]
learning rate warmup steps	1000
learning rate decay steps	500000
learning rate min value	2e-5
learning rate max value	1e-4

Sensitivity analysis on λ_{guide} A sensitivity analysis was conducted on the guidance loss weight (λ_{guide}) to assess its impact on model performance. The results that we draw from XCodec2 with codebook size 16384 at the 200k training step, as detailed in Table 7, indicate a clear trend:

1. When the weight is too small ($\lambda_{guide} = 1$), the guidance effect is negligible, yielding results similar to the baseline.
2. An optimal range is observed between $\lambda_{guide} = 5$ and 10, where the method achieves significant and robust improvements across nearly all metrics, indicating relative insensitivity to small changes within this window.
3. Conversely, values of $\lambda_{guide} \geq 15$ cause the auxiliary loss to dominate the training objective, leading to a performance drop.

This analysis confirms a stable optimal range and validates our selection of $\lambda_{guide} = 10$, providing clear practical guidance for future implementations.

A.3 GENERALIZATION EXPERIMENT ON RESIDUAL FSQ

To evaluate the general applicability of the proposed self-guidance (SG) loss beyond single-codebook models, we integrated it into a Residual FSQ architecture. The model was configured with two FSQ layers, each employing a scale of [4,4,4,4,4] (resulting in a codebook size of 1024 per layer). When integrated into the XCodec2 framework at a 50 Hz frame rate, this configuration yields a total bitrate of 1000 bps. The model was trained for 100k steps. As shown in Table 8, applying the SG loss led to consistent improvements across all objective metrics. This demonstrates that the self-guidance principle is effective not only for standard single-codebook VQ but also for multi-stage residual quantization paradigms.

Table 7: Reconstruction evaluation results of different λ_{guide} values at 200k training steps, with XCodec2 codebook size fixed at 16384.

λ_{guide}	PESQ-WB \uparrow	PESQ-NB \uparrow	STOI \uparrow	MCD \downarrow	WER \downarrow	SIM \uparrow	UTMOS \uparrow
0. (baseline)	1.9309	2.4747	0.8877	3.9591	4.08	0.7088	3.6752
1	1.9219	2.4533	0.8881	3.9527	<u>3.87</u>	0.7148	3.7266
5	<u>2.1166</u>	<u>2.6705</u>	<u>0.8977</u>	<u>3.6796</u>	3.56	0.7488	<u>3.8352</u>
10	2.1474	2.7082	0.8984	3.7086	3.87	<u>0.7428</u>	3.8395
15	2.0409	2.6074	0.8936	3.7796	3.95	<u>0.7374</u>	3.7627
50	1.9462	2.4904	0.8883	3.9035	4.18	0.7073	3.7878
100	1.8779	2.4312	0.8822	3.9648	4.40	0.6845	3.7039

Table 8: Reconstruction evaluation results on XCodec2 with Residual FSQ at 100k step, with codebook size fixed at 1024x2. (**with SG** signifies whether the proposed self-guidance mechanism is applied).

with SG	PESQ-WB \uparrow	PESQ-NB \uparrow	STOI \uparrow	MCD \downarrow	WER \downarrow	SIM \uparrow	UTMOS \uparrow
\times	1.7539	2.2503	0.8768	4.2158	4.30	0.6466	3.3923
\checkmark	1.8594	2.4154	0.8802	4.0819	4.18	0.6747	3.4105

A.4 GENERALIZATION EXPERIMENT ON BIGCODEC

To further validate the generality of our method across different model architectures, we applied the self-guidance loss to the BigCodec framework. We adapted the open-source BigCodec model to operate at a 40 Hz frame rate—comparable to our other experiments—by introducing additional down- and up-sampling blocks in the encoder and decoder. This follows the setting of 40Hz BigCodec presented in Table 2. The model used the default codebook size of 8192, with the VQ mechanism upgraded from vanilla VQ-FC to FSQ for better training stability. Specifically, unlike the Transformer blocks in the XCodec2 decoder, which consistently operate on the same feature frame rate (50Hz), each of the convolutional upsampling blocks in the BigCodec decoder operates on a different feature frame rate (gradually upsampled from 40Hz to 16kHz). Thus, in addition to the self-guidance feature mapping loss on the final block, we also insert self-guidance loss at the end of each previous upsampling block in the BigCodec decoder. For the rest of the training configuration, we follow the official implementation and conduct the evaluation when the training process reaches 100k iterations.

The evaluation results confirm that self-guidance provides a clear and consistent performance gain on this architecture, which features an RNN/CNN-based decoder fundamentally different from the Transformer-based decoder of XCodec2. This result robustly confirms that our method is a general-purpose technique for enhancing decoder robustness, independent of the specific backbone architecture.

Table 9: Reconstruction evaluation results BigCodec (framerate 40Hz) at 100k step, with codebook size fixed at 8192. (**with SG** signifies whether proposed self-guidance mechanism is applied).

with SG	PESQ-WB \uparrow	PESQ-NB \uparrow	STOI \uparrow	MCD \downarrow	WER \downarrow	SIM \uparrow	UTMOS \uparrow
\times	1.6740	2.1795	0.8601	4.3179	11.86	0.4634	3.5694
\checkmark	1.7650	2.3037	0.8655	4.2161	10.98	0.5072	3.8040

A.5 DETAILED STATISTICS OF THE ERROR DISTRIBUTIONS IN FIGURE 3

Here we provide the detailed statistics of the distributions presented in Figure 3. Table 10 contains statistics of quantization error distribution across different codebook sizes (Figure 3a), while Table 11 contains statistics of the hidden feature alignment MSEs (Figure 3b), respectively. These results align with the illustrative demonstrative, where the proposed self-guidance makes little shift to the quantization error distribution, but significantly reduces both the error level and dispersion in hidden feature alignment MSE. Specifically, we also observe that the latter effect becomes more obvious as the codebook size further shrinks.

Table 10: Quantization Error distributions across different codebook sizes.

Codebook size	With self-guidance	Error mean	Error std
65536	No	0.858	0.120
	Yes	0.851	0.121
16384	No	0.798	0.121
	Yes	0.799	0.120
8192	No	0.741	0.120
	Yes	0.744	0.121

Table 11: Hidden MSE distributions across different codebook sizes.

Codebook size	With self-guidance	Error mean	Error std
65536	No	9.439	29.203
	Yes	5.854	13.712
16384	No	13.551	59.137
	Yes	4.958	5.863
8192	No	23.605	109.458
	Yes	4.197	6.865

A.6 RECONSTRUCTION SAMPLES

Here we provide the reconstruction samples on the LibriSpeech *test-clean* dataset to reveal the typical quantization artifacts in the generated audio caused by quantization error, including smeared harmonics (Figure 4), pitch spikes (Figure 5), and oversmoothed harmonics (Figure 6). The audio of the corresponding sampling could be listened to on our demo website.¹⁰

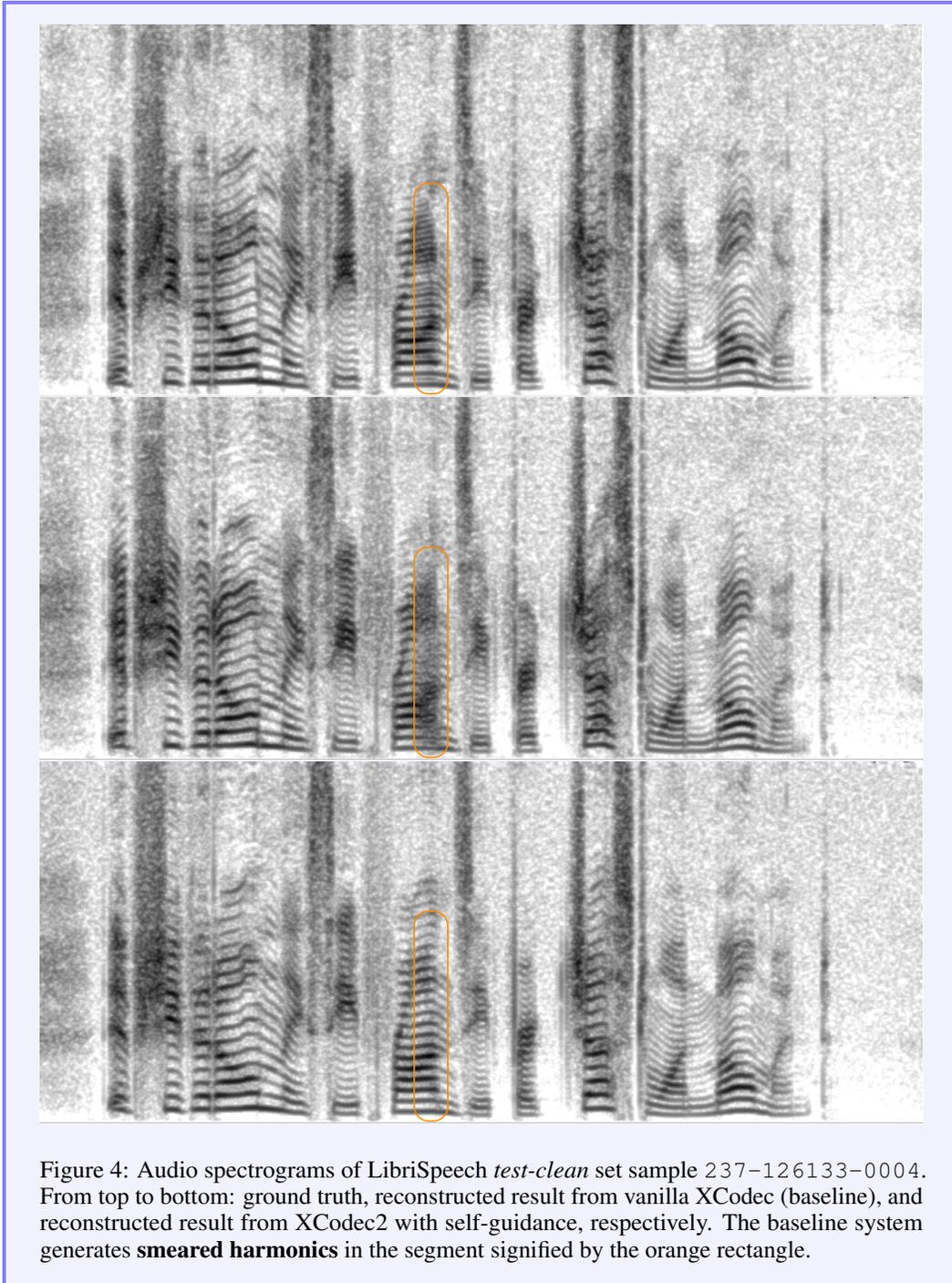
B LLM USAGE

LLM is involved in the production of this paper in the following ways:

1. Polish the human-written manuscripts, correcting grammar and spelling errors, enhancing readability and clarity of the paper.

¹⁰<https://sgvqvae.github.io/sgvqvae-demo>

918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971



2. Assisting in implementing the code for dataset preprocessing, as well as the collection and visualization of the evaluation results.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

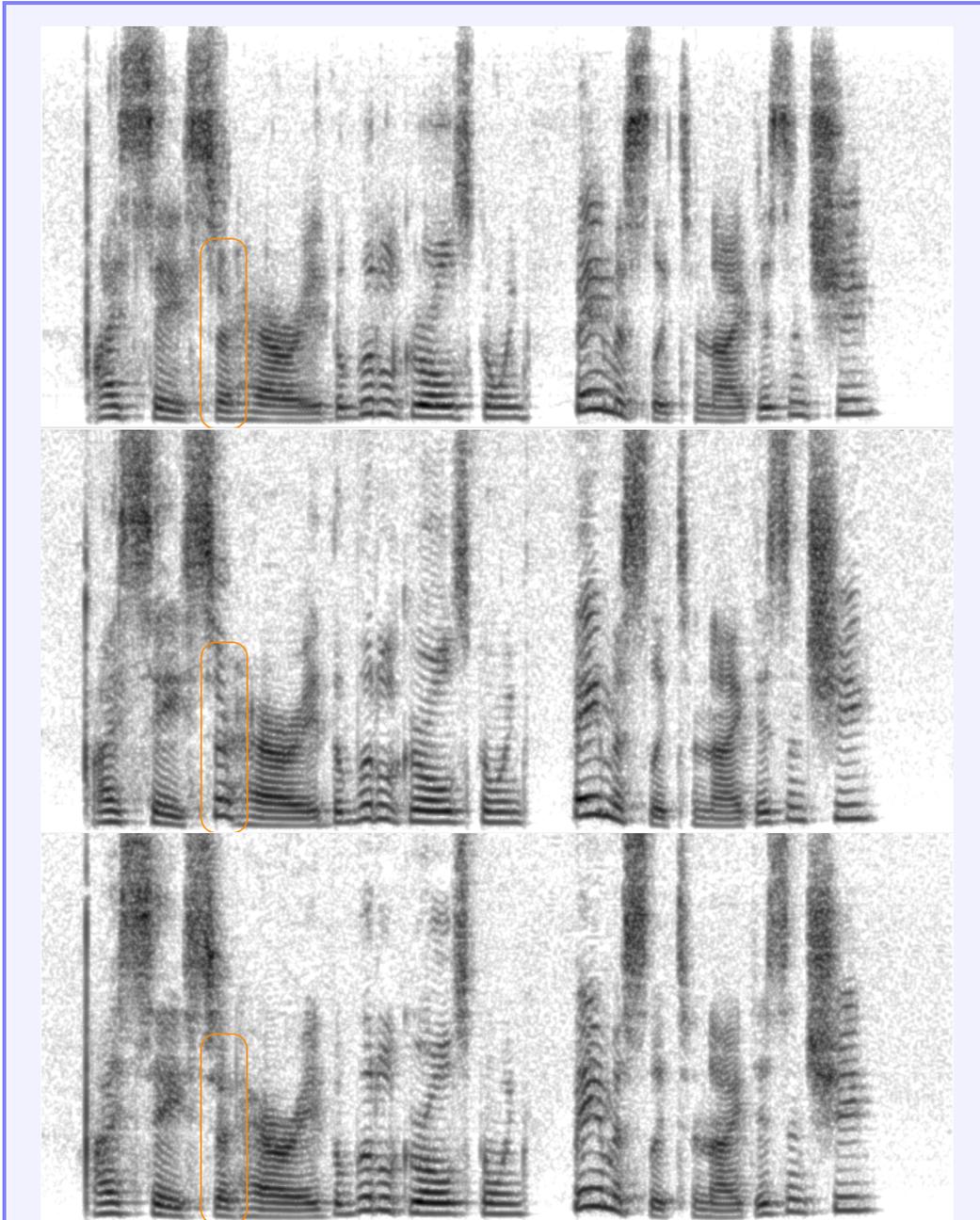


Figure 5: Audio spectrograms of LibriSpeech *test-clean* set sample 4446-2271-0012. From top to bottom: ground truth, reconstructed result from vanilla XCodec (baseline), and reconstructed result from XCodec2 with self-guidance, respectively. The baseline system generates **pitch spike** in the segment signified by the orange rectangle.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

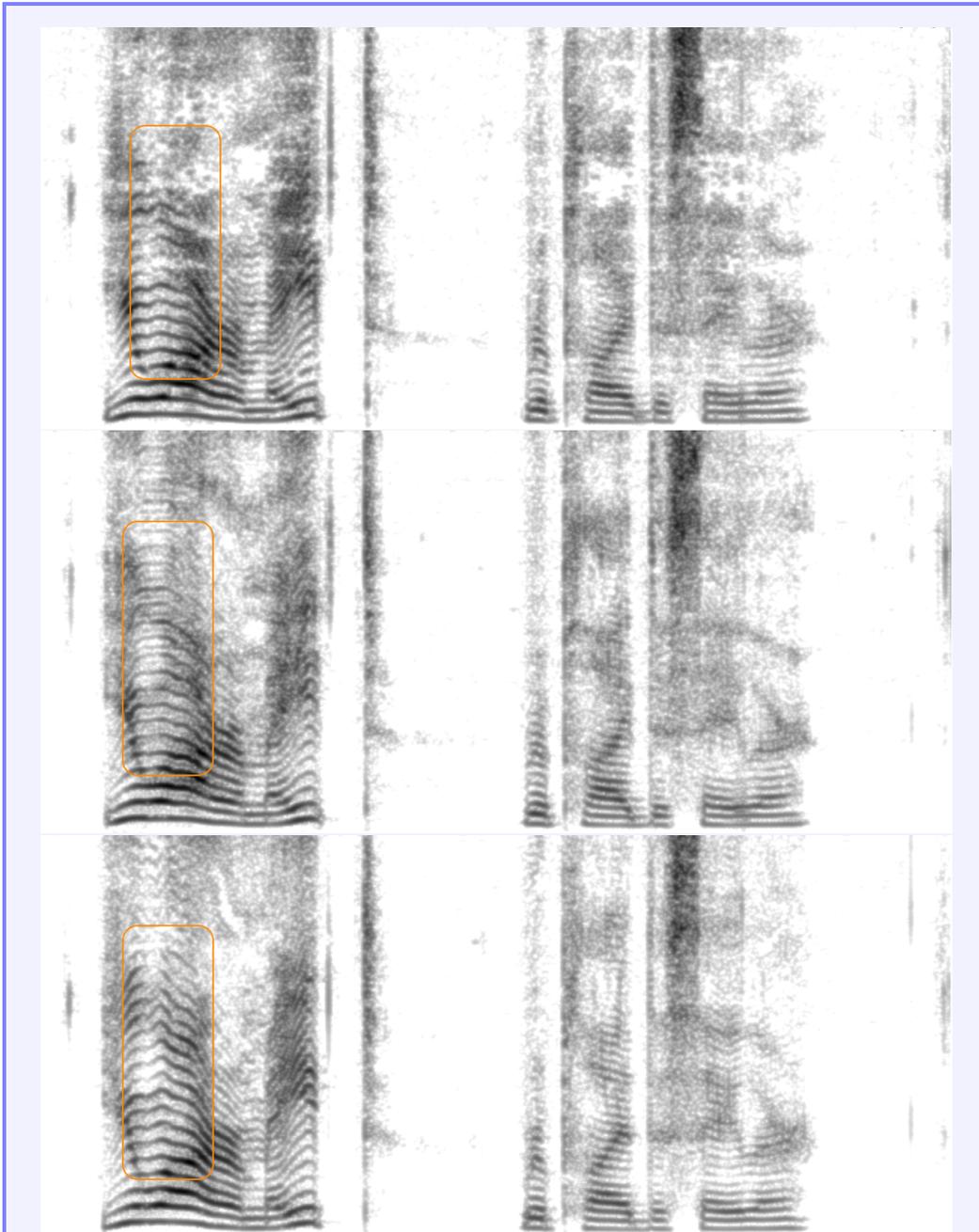


Figure 6: Audio spectrograms of LibriSpeech *test-clean* set sample 8555-284449-0009. From top to bottom: ground truth, reconstructed result from vanilla XCodec (baseline), and reconstructed result from XCodec2 with self-guidance, respectively. The baseline system generates **oversmoothed harmonics** in the segment signified by the orange rectangle.