

Motivation & Innovations

Long-form, multi-speaker conversational audio (podcasts, audiobooks) demands speaker consistency, natural turn-taking, and expressive cues (breaths, pacing). Prior systems limited to **2 speakers, <10 min**; concatenating single-speaker utterances loses conversational naturalness.

- 1 **7.5 Hz hybrid tokenizers** — decoupled acoustic (σ -VAE) + semantic (ASR), 6–80× more compact than prior codecs
- 2 **LLM + Diffusion Head** — Qwen2.5 for dialogue flow, diffusion head for high-fidelity acoustic generation
- 3 **Podcast data pipeline** — segmentation + diarization + quality filtering; no speech enhancement preserves expressive cues
- 4 **90 min, 4 speakers** — zero-shot generation, far beyond MoonCast (2 speakers, <10 min)

7.5 Hz

Frame Rate

3.76

Best MOS

90 min

Max Duration

4

Speakers

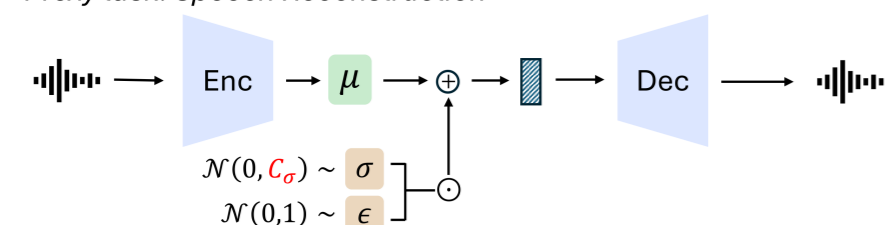
Decoupled Speech Tokenizers

At 7.5 Hz (3200× compression from 24 kHz), a single encoder cannot balance acoustic fidelity and semantic understanding. Coupled tokenizer improves WER (6.22 → 3.55) but speaker similarity collapses (0.68 → 0.45). **Decoupled hybrid achieves the best balance: WER 1.84, SIM-O 0.64.**

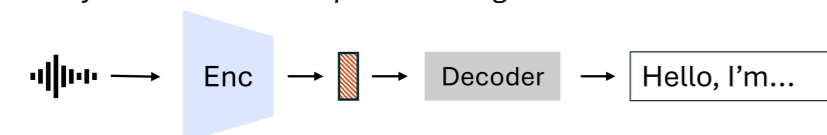
$$\mu = \text{Enc}_{\phi}(\mathbf{x}), \quad z = \mu + \sigma \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I), \quad \sigma \sim \mathcal{N}(0, C_{\sigma})$$

Unlike standard VAE where σ is learned, σ -VAE samples from a fixed prior — preventing variance collapse in AR generation.

Proxy task: Speech Reconstruction



Proxy task: Automatic Speech Recognition



Acoustic Tokenizer (upper): σ -VAE reconstructs waveform. Semantic Tokenizer (lower): ASR proxy task. Both at 7.5 Hz.

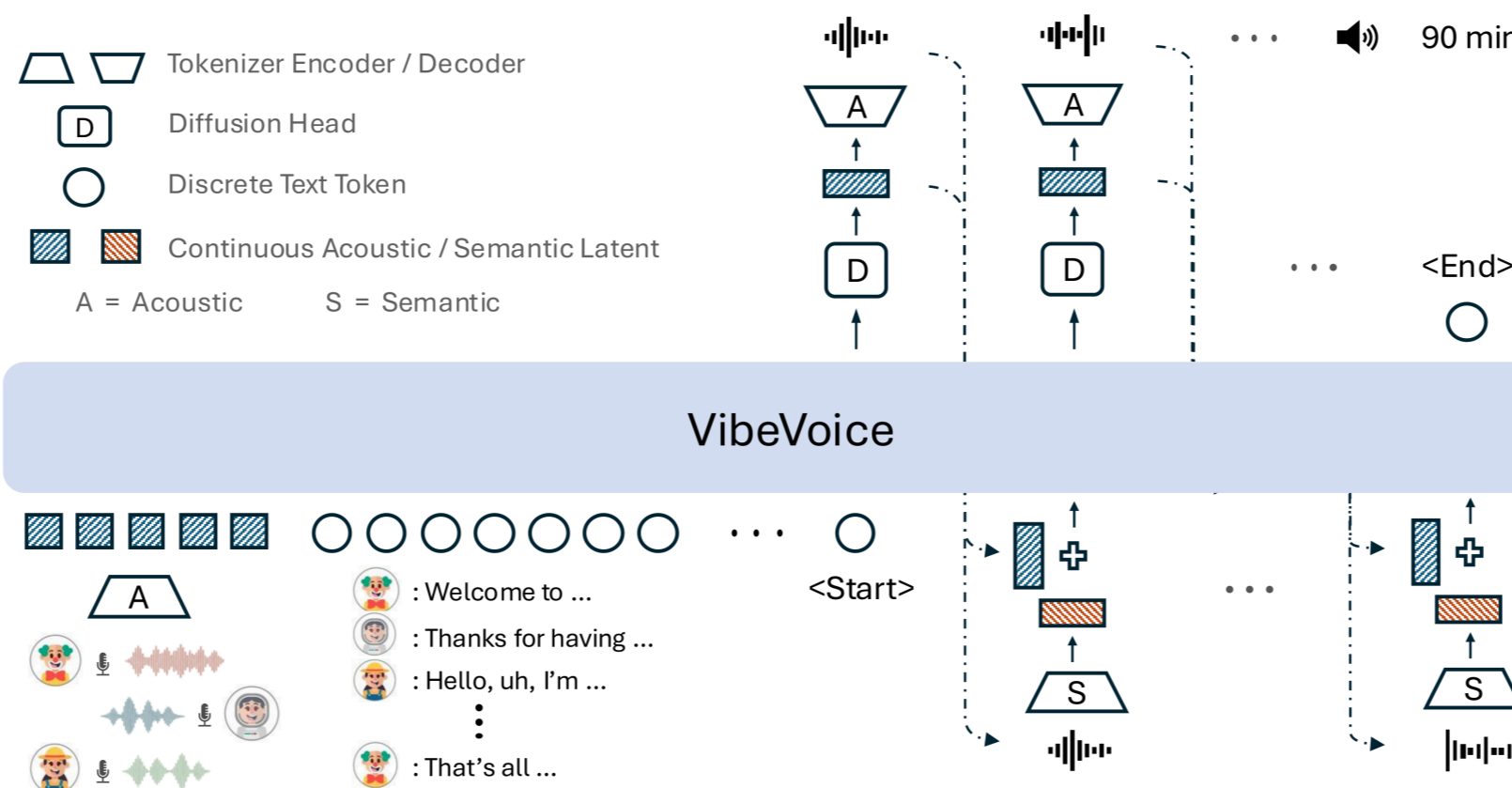
Acoustic (σ -VAE): encoder-decoder ~340M params, continuous latent space, trained with DAC discriminator. PESQ 3.068, UTMOS 4.181 at 7.5 Hz.

Semantic (ASR proxy): mirrors acoustic encoder architecture, trained on ASR task for content-centric features. Decoder discarded after training.

Why σ -VAE Enables Next-Token Diffusion

The Transformer runs **once** per token; only the tiny diffusion head iterates. Critically, **σ -VAE enables Next-Token Diffusion**: standard VAE's learned variance collapses to near-zero, leaving no room for generation error. σ -VAE enforces a fixed variance $\sigma \sim \mathcal{N}(0, C_{\sigma})$, creating a **tolerance zone** in latent space that absorbs the inevitable error accumulation from autoregressive generation — enabling stable synthesis over 90 minutes.

Architecture Overview



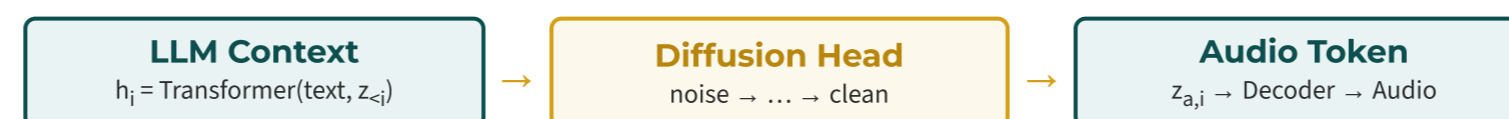
User Input: Voice & Text Scripts

Voice prompts + text scripts → LLM with hybrid context → Diffusion Head generates acoustic tokens → decoder synthesizes waveform.

VibeVoice uses an LLM backbone with **Next-Token Diffusion** for audio generation. At each autoregressive step, the LLM reads text + previous audio context, then a lightweight diffusion head generates the next continuous audio token — decoded into waveform by the acoustic decoder.

Next-Token Diffusion

Standard LLMs predict the next *word* via softmax over a vocabulary. But audio lives in **continuous space** — there is no dictionary to choose from. VibeVoice replaces softmax with a **diffusion head**: a small network (4 layers, ~123M params) that iteratively shapes random noise into an audio token, guided by the LLM's hidden state.



Standard LLM: $h_i \rightarrow \text{softmax} \rightarrow \text{discrete word}$ | VibeVoice: $h_i \rightarrow \text{diffusion} \rightarrow \text{continuous audio vector}$

$$\text{Train: } \mathcal{L}_{\text{Diff}} = \mathbb{E} \left\| \epsilon - \epsilon_{\theta}(z_a^{(t)}, t, h_i) \right\|^2 \quad \text{— predict the noise added to clean token}$$

$$\text{Infer (CFG): } \hat{\epsilon} = \epsilon_{\text{uncond}} + w \cdot (\epsilon_{\text{cond}} - \epsilon_{\text{uncond}}) \quad \text{— amplify LLM conditioning (w=1.3)}$$

Main Results

VibeVoice-Eval benchmark: 108 podcast samples (1–30 min). Subjective: 24 annotators, ~6h audio.

Realism		Richness		Preference	
SesameAI-CSM	2.89	SesameAI-CSM	3.03	SesameAI-CSM	2.75
Higgs Audio V2	2.95	Higgs Audio V2	3.19	Higgs Audio V2	2.83
ElevenLabs v3a	3.34	ElevenLabs v3a	3.48	ElevenLabs v3a	3.38
Gemini 2.5 Pro	3.55	Gemini 2.5 Pro	3.78	Gemini 2.5 Pro	3.65
VibeVoice-1.5B	3.59	VibeVoice-1.5B	3.59	VibeVoice-1.5B	3.44
VibeVoice-7B	3.71	VibeVoice-7B	3.81	VibeVoice-7B	3.75

Objective — Best WER-W **1.11** (VibeVoice-1.5B) | Best SIM-O **0.692** (VibeVoice-7B)

Tokenizer Reconstruction

7.5 Hz with 3200× compression. Continuous σ -VAE achieves best PESQ and UTMOS on LibriTTS test-clean.

Model	N_q	Rate	PESQ ↑	UTMOS ↑
Codec	8	600	2.72	3.04
DAC	4	400	2.738	3.433
WavTokenizer	1	75	2.373	4.049
Ours (Acoustic)	N/A	7.5	3.068	4.181

Scalability

Model	Short (0~12 min)		Long (12~30 min)	
	WER ↓	SIM-O ↑	WER ↓	SIM-O ↑
CosyVoice2 (concat)	4.27	0.73	4.95	0.74
MoonCast	10.4	0.55	crashes †	–
VibeVoice-1.5B	2.11	0.59	1.55	0.59
VibeVoice-7B	0.66	0.75	1.24	0.75

VibeVoice supports up to **90 min** duration with **4** distinct speakers — MoonCast is limited to 2 speakers and crashes on long-form (>12 min) generation.

Conclusion

VibeVoice: zero-shot, long-form, multi-speaker podcast generation. 7.5 Hz tokenizers + LLM next-token diffusion → SOTA on subjective and objective metrics. **90 min, 4 speakers.**



Code & Checkpoints

GitHub Star History

