

M3-TTS: MULTI-MODAL DIT ALIGNMENT & MEL-LATENT FOR ZERO-SHOT HIGH-FIDELITY SPEECH SYNTHESIS

Xiaopeng Wang^{1,†}, Chunyu Qiang^{2,†}, Ruibo Fu^{3,*}, Zhengqi Wen³, Xuefei Liu³, Yukun Liu³, Yuzhe Liang², Kang Yin², Yuankun Xie³, Heng Xie¹, Chenxing Li³, Chen Zhang², Changsheng Li¹

¹ Beijing Institute of Technology, Beijing, China ² Kuaishou Technology, Beijing, China

³ Institute of Automation, Chinese Academy of Sciences, Beijing, China

ABSTRACT

Non-autoregressive (NAR) text-to-speech synthesis relies on length alignment between text sequences and audio representations, constraining naturalness and expressiveness. Existing methods depend on duration modeling or pseudo-alignment strategies that severely limit naturalness and computational efficiency. We propose M3-TTS, a concise and efficient NAR TTS paradigm based on multi-modal diffusion transformer (MM-DiT) architecture. M3-TTS employs joint diffusion transformer layers for cross-modal alignment, achieving stable monotonic alignment between variable-length text-speech sequences without pseudo-alignment requirements. Single diffusion transformer layers further enhance acoustic detail modeling. The framework integrates a mel-vae codec that provides 3× training acceleration. Experimental results on Seed-TTS and AISHELL-3 benchmarks demonstrate that M3-TTS achieves state-of-the-art NAR performance with the lowest word error rates (1.36% English, 1.31% Chinese) while maintaining competitive naturalness scores. Code and demos will be available at <https://wwwwpxp.github.io/M3-TTS-Demo>.

Index Terms— Text-to-Speech, MMDiT, Mel-VAE

1. INTRODUCTION

In recent years, text-to-speech (TTS) has achieved substantial gains in fidelity and naturalness. Existing methods generally fall into two paradigms: autoregressive (AR) [1, 2, 3, 4, 5] and non-autoregressive (NAR) [6, 7, 8, 9, 10]. AR models generate speech frame by frame [11, 12] and thus avoid explicit duration modeling; without hard duration constraints [13, 14], they often produce more expressive prosody and higher naturalness. However, autoregressive decoding results in slow inference, and the teacher-forcing training paradigm induces train-inference mismatch (exposure bias), which can undermine stability. By contrast, NAR models typically model the entire acoustic sequence and synthesize speech in parallel, offering significantly faster inference. Nevertheless, most

NAR systems depend on text-speech alignment and duration constraints [15, 3, 16, 9, 6, 8] (e.g., duration predictors, uniform upsampling, or filler padding), which can impose over-regularized timing and pauses and thus average out prosody and expressive variation.

A central challenge for NAR TTS is reliable text-speech alignment. Early solutions employ duration-based aligners: FastSpeech series [17, 18] depends on forced alignments to provide explicit duration targets, whereas VITS [19, 20, 15] uses monotonic alignment search to infer durations without labeled supervision. Although effective, these strategies impose strong duration constraints that tend to average out prosody and limit expressiveness. Motivated by these limitations, recent Conditional Flow Matching (CFM) [21] NAR models remove phoneme-level duration modeling. For example, F5-TTS [6] and E2-TTS [9] pad the text sequence with filler tokens until it matches the mel length, while ZipVoice [8] applies uniform upsampling to assign equal duration to every token. However, both padding and uniform upsampling are surrogate (pseudo-alignment) mechanisms aimed at matching the audio sequence length; such proxy alignment can restrict the model’s ability to learn natural timing and rhythm, and it wastes computation by inflating sequences or introducing redundant operations.

In this paper, we introduce M3-TTS, a non-autoregressive TTS framework that couples a Multi-Modal Diffusion Transformer (MMDiT) architecture with a Mel-VAE latent acoustic target. The model establishes reliable text-speech correspondence without resorting to pseudo-alignment and supports efficient inference; its latent pathway compresses speech in time and dimension, reducing sequence length and memory, stabilizing optimization, and enabling zero-shot synthesis at 44.1 kHz. M3-TTS comprises three elements:

(1) Learnable cross-modal attention performs dynamic, variable-length correspondence between text and speech tokens, eliminating padding or uniform upsampling.

(2) The VAE latent provides low-dimensional continuous modeling, shortens both temporal span and feature dimension, lowers GPU memory, improves optimization stability, and natively supports 44.1 kHz synthesis.

[†] denotes equal contribution. * denotes corresponding author.

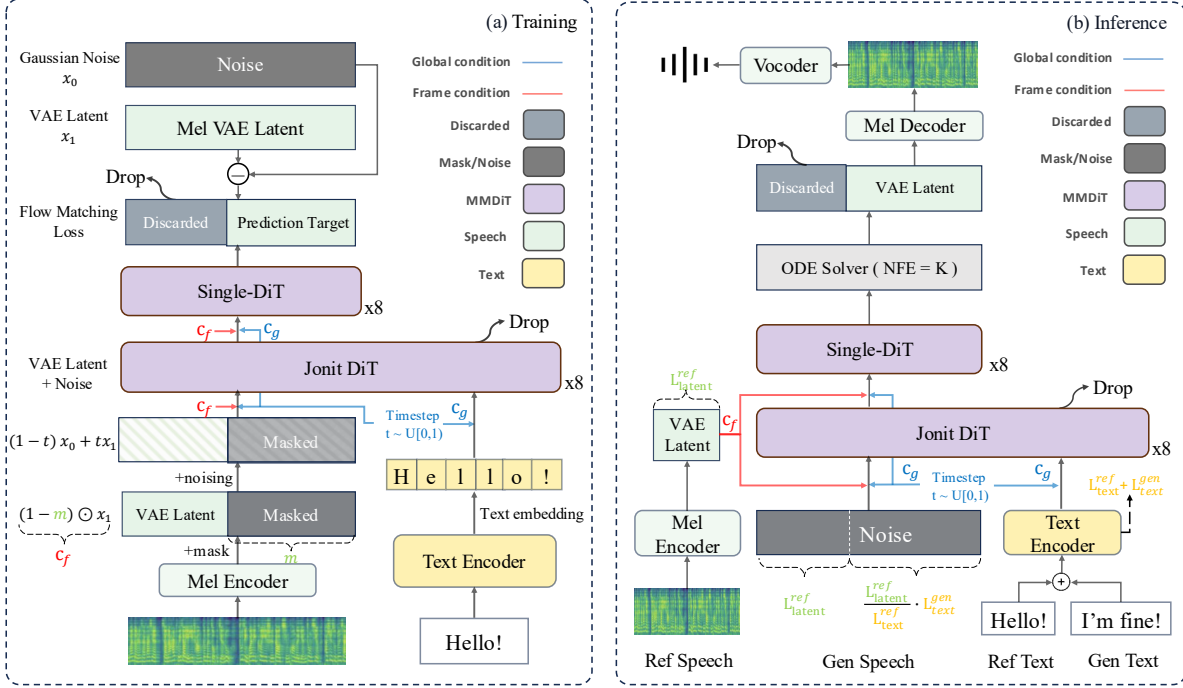


Fig. 1. Overview of M3-TTS: training (left) and inference (right). Text is encoded into T , and the reference speech is encoded by a Mel-VAE into latents x_1 . Noise x_0 and x_1 are linearly interpolated to form x_t , and conditioning is applied with a global token c_g and a frame-level token $c_f = c_g + (1 - m) \odot x_1$. Joint-DiT aligns $[x_t; T]$ in a unified attention space and splits the output into (H^a, H^t) . Single-DiT refines only the speech branch H^a . During inference, the output length using the reference speech-to-text ratio; an ODE solver integrates from noise to a latent; the Mel decoder then decodes it into a spectrogram.

(3) Experimental results on Seed-TTS, M3-TTS attains the lowest WER (EN 1.36%, ZH 1.31%) among compared systems while maintaining competitive naturalness and speaker similarity.

2. M3-TTS

2.1. Overview

M3-TTS is designed to overcome the limitations of cross-modal alignment and the inefficiency of high-dimensional mel features. To this end, Mel-VAE compresses speech into a low-dimensional latent space, and a two-stage DiT is adopted: Joint-DiT aligns text representations T with speech latents x_t , while Single-DiT refines the speech branch to predict the vector field v_θ for CFM. During inference, an ODE solver integrates noise into a latent, which is then decoded by the Mel decoder to reconstruct speech. This design avoids padding/upsampling, improves alignment stability and prosody naturalness.

2.2. Mel-VAE Codec

The Mel-VAE codec, based on the VQ-CTAP [22] design and comprising a Mel encoder-decoder pair, normalizes

speech to 44.1 kHz and produces a latent sequence at ~ 43 Hz. Compared to mainstream NAR systems, it yields roughly $2\times$ temporal compression and $2.5\times$ dimensional compression ($100 \rightarrow 40$), thereby reducing training/inference memory and compute while maintaining fidelity. Moreover, because the predictor outputs a distribution over latents rather than deterministic log-mel amplitudes, it exhibits improved robustness.

2.3. Multi-Modal Diffusion Transformers

MMDiT architecture comprises two modules: a Joint-DiT for cross-modal alignment and a Single-DiT for speech-only refinement.

Joint-DiT. Let speech latents $A \in \mathbb{R}^{B \times T_s \times D}$ and text features $T \in \mathbb{R}^{B \times T_t \times D}$. We concatenate them along the temporal axis and model the unified sequence with shared attention:

$$Z = \text{Concat}_{\text{time}}(A, T) \in \mathbb{R}^{B \times (T_s + T_t) \times D}. \quad (1)$$

Joint-DiT stacks pre-normalized Transformer blocks with *shared* scaled dot-product self-attention and a feed-forward sublayer, aided by modality tags and positional/time embeddings; RoPE is applied to Q, K before attention. Conditions are injected via AdaLN: the text stream uses the global time condition $c_g = \text{Emb}(t)$, while the speech stream uses the

Table 1. Objective (SIM-o \uparrow , WER \downarrow , UTMOS \uparrow) and subjective (NMOS \uparrow , QMOS \uparrow) results on AISHELL3-test (44.1 kHz) and Seed-TTS test-en/zh (24 kHz). Compared are AR, NAR, and VAE Reconstruction as the codec upper bound. Bold indicates the best result; underlining denotes the second best.

Model	Data (hrs)	Params	AISHELL3-test (44.1k)			Seed-TTS test-en (24k)			Seed-TTS test-zh (24k)			Subjective Metrics	
			SIM-o \uparrow	WER \downarrow	UTMOS \uparrow	SIM-o \uparrow	WER \downarrow	UTMOS \uparrow	SIM-o \uparrow	WER \downarrow	UTMOS \uparrow	NMOS \uparrow	QMOS \uparrow
Ground Truth	—	—	0.631	6.20	2.50	0.734	2.14	3.52	0.755	1.25	2.78	3.78 \pm 0.09	3.95 \pm 0.08
VAE Reconstruction	—	—	0.584	5.19	1.85	0.631	1.85	2.34	0.699	1.35	1.89	—	—
CosyVoice	170K Multi.	416M	—	—	—	0.609	4.29	—	0.723	3.63	—	—	—
CosyVoice 2	167K Multi.	618M	—	—	—	0.652	2.57	—	0.748	1.45	—	—	—
Spark-TTS	102K Multi.	507M	—	—	—	0.584	1.98	—	0.672	1.20	—	—	—
E2-TTS (32 NFE)	100K Emilia	333M	—	—	—	0.706	2.32	3.21	0.713	1.91	2.26	—	—
F5-TTS (32 NFE)	100K Emilia	336M	—	—	—	0.664	1.85	3.72	0.750	1.53	2.93	3.76 \pm 0.20	3.90 \pm 0.18
F5-TTS (32 NFE)	100K Emilia	155M	—	—	—	0.628	1.96	3.66	0.733	1.57	2.93	—	—
ZipVoice (16 NFE)	100K Emilia	123M	—	—	—	<u>0.697</u>	1.70	<u>3.82</u>	<u>0.751</u>	1.40	<u>3.15</u>	3.78 \pm 0.15	3.95 \pm 0.13
M3-TTS-Fbank (32 NFE)	100K Emilia	355M	—	—	—	0.681	<u>1.48</u>	3.88	0.762	1.36	3.18	3.80 \pm 0.12	3.99 \pm 0.11
M3-TTS-VAE (32 NFE)	100K Emilia	355M	0.540	10.7	1.78	0.604	1.36	2.80	0.621	<u>1.31</u>	2.18	3.62 \pm 0.19	3.75 \pm 0.17

frame-level condition:

$$c_f = c_g + (1 - m) \odot A, \quad (2)$$

where $m \in \{0, 1\}^{B \times T_s}$ is a frame-level binary mask (broadcast along the feature dimension). After cross-modal fusion, the output preserves its shape and is split back into (H^a, H^t) , yielding explicit text–speech alignment.

Single-DiT. We refine only the speech branch H^a , still conditioned on c_f , and finally output the vector field $v_\theta(\cdot)$ for CFM; the text branch is dropped at this stage.

2.4. Training and Inference

Training. We sample $x_0 \sim p_0$ and $t \sim \mathcal{U}(0, 1)$, and form the interpolant $x_t = (1 - t)x_0 + tx_1$. A binary mask m yields a masked view of x_1 ; the global time condition is $c_g = \text{Emb}(t)$ and the fused condition is $c_f = c_g + (1 - m) \odot x_1$. With target velocity $u_t = \dot{\alpha}(t)(x_1 - x_0)$ (linear schedule $\alpha(t) = t$ reduces to $\dot{\alpha}(t)(x_1 - x_0)$), we minimize:

$$\mathcal{L}(\theta) = \mathbb{E} \left[\left\| v_\theta(x_t, t, c_f) - \dot{\alpha}(t)(x_1 - x_0) \right\|_2^2 \right]. \quad (3)$$

Inference. The generated latent length is set by:

$$L_{\text{gen}} = \text{round} \left(\frac{L_{\text{speech}}^{\text{ref}}}{L_{\text{text}}^{\text{ref}}} \cdot L_{\text{text}}^{\text{tar}} \right). \quad (4)$$

Starting from $x_0 \sim p_0$, we integrate $\dot{x}_t = v_\theta(x_t, t, c_f)$ for $t \in [0, 1]$ to obtain x_1 , then decode it to Mel and waveform.

3. EXPERIMENTAL SETUP

Training configuration. We use a 16-layer MMDiT acoustic model—8 Joint-DiT layers and 8 Single-DiT layers—with model dimension 640 and 10 attention heads ($\approx 355\text{M}$ parameters). The text encoder follows the ZipVoice [8] design with a 4-layer Zipformer [23]. Training is performed on $8 \times \text{A100}$ GPUs with batch size 192 and learning rate 7.5×10^{-5} . For

the infilling objective, we randomly mask 70–100% of Mel frames. For CFG [24] training, masked speech and text inputs are independently dropped with probability 0.2. To verify the advantages of the M3-TTS architecture and keep comparability with mainstream NAR pipelines, we train two acoustic targets under the *same* architecture and schedule: (i) an Fbank variant using 100-dimensional log-mel filterbanks at 24 kHz with hop length 256, decoded by the Vocos vocoder [25]; and (ii) a Mel-VAE latent variant decoded by BigVGAN [26].

Datasets. We train on the Emilia corpus [27] (approximately $\sim 95\text{k}$ hours of English and Chinese) after filtering transcription errors and other anomalies. Zero-shot TTS is evaluated on three benchmarks: Seed-TTS [14] test-en (1,088 English utterances from Common Voice), Seed-TTS test-zh (2,020 Chinese utterances from DiDiSpeech), and a 44.1 kHz test set of 1,000 samples constructed from AISHELL-3 [28], following the Seed-TTS evaluation protocol.

Baselines. We compare against representative AR and NAR systems. *AR models:* CosyVoice [1], CosyVoice2 [2], Spark-TTS [4]. *NAR models:* MaskGCT [3], E2-TTS [9], F5-TTS [6], ZipVoice [8].

Metrics. We adopt a cross-sentence, zero-shot setting with both reproducible, model-based metrics and human perception scores. Intelligibility is measured by WER using ASR backends: Whisper-large-v3 [29] for English and Paraformer-zh [30] for Chinese. Speaker similarity (SIM-o) is computed as the cosine similarity between WavLM-based ECAPA-TDNN embeddings [31] extracted from the prompt and synthesized speech. Naturalness is estimated by UTMOS [32]. For human evaluation, we report NMOS (naturalness MOS) and QMOS (quality MOS).

4. EXPERIMENTAL RESULTS

4.1. Overall Comparison

Table 1 reports objective and subjective results across datasets. Compared with representative NAR systems (E2-TTS, F5-TTS, ZipVoice, MaskGCT) and AR systems (CosyVoice,

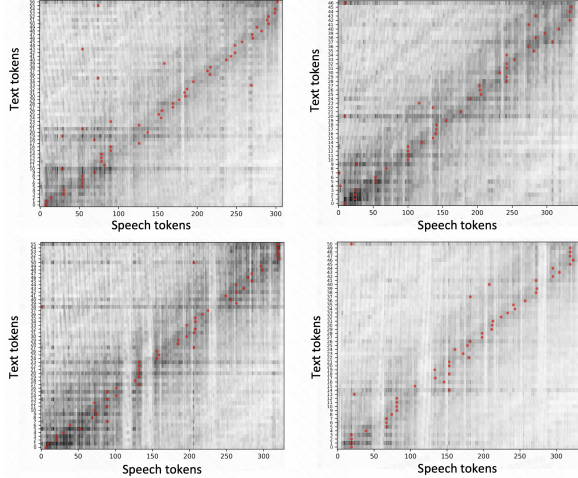


Fig. 2. Joint-DiT cross-modal attention visualization for four test samples. Each heatmap shows attention from speech tokens (rows) to text tokens (columns); red dots indicate the row-wise argmax.

Table 2. Training time on the Emilia corpus ($8 \times A100$, batch size 192), comparing VAE and Fbank features under the MMDiT architecture.

Variant	Params (M)	Time (h)	Speedup (\times)
M3-TTS-Fbank	355	90	1.0
M3-TTS-VAE	355	31	2.9

CosyVoice2, Spark-TTS), M3-TTS-Fbank attains strong overall performance: on Seed-TTS en, it achieves the lowest WER among NAR models (1.48) and the highest UTMOS (3.88) with SIM-o 0.681; on Seed-TTS zh, it yields the best SIM-o (0.762), the best UTMOS (3.18), and a competitive WER (1.36, second overall). In double-blind listening tests, M3-TTS-Fbank obtains NMOS/QMOS of 3.80/3.99, surpassing strong NAR baselines (e.g., ZipVoice 3.78/3.95) and closely matching or slightly exceeding ground truth (3.78/3.95). M3-TTS-VAE further attains the lowest WER on Seed-TTS en (1.36) and the lowest NAR WER on Seed-TTS zh (1.31), albeit with lower SIM-o and UTMOS (2.80/2.18). We attribute the gains primarily to Joint-DiT, which aligns text and speech in a unified attention space and avoids alignment bias from filler-token padding and uniform upsampling, thereby improving intelligibility and prosody.

To assess the impact of representation choice, we evaluate M3-TTS-VAE on AISHELL3-test (44.1 kHz) and Seed-TTS (24 kHz), and include VAE Reconstruction to indicate the codec upper bound. At 44.1 kHz, M3-TTS-VAE performs below ground truth and the reconstruction ceiling: WER 10.7, SIM-o 0.540, UTMOS 1.78, compared with ground truth and VAE Reconstruction, suggesting constraints from accumulated generation errors and codec capacity under a high sampling rate and cross-corpus evaluation. At 24 kHz, M3-

TTS-VAE often attains lower WER than M3-TTS-Fbank (en 1.36, zh 1.31) but lags on SIM-o/UTMOS. These results indicate that the VAE latent’s roughly $2\times$ temporal compression and distributional prediction ease alignment and regression (benefiting WER), whereas naturalness and timbral detail are limited by codec bandwidth and the reconstruction ceiling; domain shift across corpora and sampling rates further amplifies the gap.

4.2. Joint Attention Visualization

Figure 2 presents four examples of Joint-DiT cross-modal attention, with speech tokens on the x -axis and text tokens on the y -axis; brighter intensities indicate larger attention weights, and red dots mark the row-wise argmax (the best speech alignment for each text token). The maps exhibit an approximately monotonic diagonal structure with minimal off-diagonal drift. These qualitative observations suggest that the unified attention space learns stable text-speech alignment.

4.3. Training Efficiency

We benchmark training wall-clock time on the Emilia corpus under controlled conditions: $8 \times A100$ GPUs, batch size 192, and an identical training schedule. As shown in Table 2, M3-TTS-Fbank completes in 90 h, whereas M3-TTS-VAE completes in 31 h, yielding a $\sim 3\times$ speedup. We attribute this gain primarily to the reduced sequence length and dimensionality in the Mel VAE pathway, which increases throughput without altering optimization hyperparameters.

4.4. Discussion

Despite strong empirical results, M3-TTS has two main limitations. First, Mel-VAE was trained on a small, mixed-sampling-rate corpus, which likely limits latent expressiveness; scaling to larger, better-balanced data with more 44.1 kHz audio should improve fidelity and robustness. Second, inference concatenates the prompt text and reference audio before the target sequence (as in ZipVoice and F5-TTS), which adds latency and limits zero-shot flexibility.

5. CONCLUSION

In this work, we introduced M3-TTS, a multimodal text-to-speech system that combines Mel-VAE latent representations with MMDiT alignment for TTS. Our approach features a Joint-DiT mechanism that enables dynamic cross-modal attention between text and speech, eliminating traditional padding and upsampling limitations. The Mel-VAE codec provides low-dimensional continuous modeling with significant memory reduction. On Seed-TTS and AISHELL-3, M3-TTS attains state-of-the-art non-autoregressive results. In the future, we plan to extend this framework to dialogue.

6. REFERENCES

- [1] Zhihao Du, Qian Chen, et al., “Cosyvoice: A scalable multilingual zero-shot text-to-speech synthesizer based on supervised semantic tokens,” *arXiv preprint arXiv:2407.05407*, 2024.
- [2] Zhihao Du, Yuxuan Wang, et al., “Cosyvoice 2: Scalable streaming speech synthesis with large language models,” *arXiv preprint arXiv:2412.10117*, 2024.
- [3] Yuancheng Wang, Haoyue Zhan, Liwei Liu, and et al, “MaskGCT: Zero-shot text-to-speech with masked generative codec transformer,” *International Conference on Representation Learning*, 2025.
- [4] Xinsheng Wang, Mingqi Jiang, and et al, “Spark-TTS: An efficient llm-based text-to-speech model with single-stream decoupled speech tokens,” *arXiv preprint arXiv:2503.01710*, 2025.
- [5] Chunyu Qiang, Hao Li, Hao Ni, He Qu, Ruibo Fu, Tao Wang, Longbiao Wang, and Jianwu Dang, “Minimally-supervised speech synthesis with conditional diffusion model and language model: A comparative study of semantic coding,” in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 10186–10190.
- [6] Yushen Chen, Zhikang Niu, and et al, “F5-TTS: A fairytale that fakes fluent and faithful speech with flow matching,” *Proc. ACL*, 2025.
- [7] Qixi Zheng, Yushen Chen, Zhikang Niu, and et al, “Accelerating Flow-Matching-Based Text-to-Speech via Empirically Pruned Step Sampling,” in *Interspeech 2025*, 2025, pp. 2445–2449.
- [8] Han Zhu, Wei Kang, and et al, “Zipvoice: Fast and high-quality zero-shot text-to-speech with flow matching,” *ASRU*, 2025.
- [9] Sefik Emre Eskimez, Xiaofei Wang, et al., “E2 tts: Embarrassingly easy fully non-autoregressive zero-shot tts,” in *SLT*, 2024.
- [10] Chunyu Qiang, Kang Yin, Xiaopeng Wang, Yuzhe Liang, Jiahui Zhao, Ruibo Fu, Tianrui Wang, Cheng Gong, Chen Zhang, Longbiao Wang, et al., “Instructaudio: Unified speech and music generation with natural language instruction,” *arXiv preprint arXiv:2511.18487*, 2025.
- [11] Chunyu Qiang, Haoyu Wang, Cheng Gong, Tianrui Wang, Ruibo Fu, Tao Wang, Ruilong Chen, Jiangyan Yi, Zhengqi Wen, Chen Zhang, et al., “Secousticodec: Cross-modal aligned streaming single-codecbook speech codec,” *arXiv preprint arXiv:2508.02849*, 2025.
- [12] Chunyu Qiang, Wang Geng, Yi Zhao, Ruibo Fu, Tao Wang, Cheng Gong, Tianrui Wang, Qiuyu Liu, Jiangyan Yi, Zhengqi Wen, et al., “Vq-ctap: Cross-modal fine-grained sequence representation learning for speech processing,” *IEEE Transactions on Audio, Speech and Language Processing*, 2025.
- [13] Sanyuan Chen, Chengyi Wang, and et al, “Neural codec language models are zero-shot text to speech synthesizers,” *IEEE Transactions on Audio, Speech and Language Processing*, 2025.
- [14] Philip Anastassiou, Jiawei Chen, and et al, “Seed-TTS: A family of high-quality versatile speech generation models,” *arXiv preprint arXiv:2406.02430*, 2024.
- [15] Edresson Casanova et al., “Yourtts: Towards zero-shot multi-speaker tts and zero-shot voice conversion for everyone,” in *ICLR*, 2022.
- [16] Zeqian Ju, Yuancheng Wang, et al., “Naturalspeech 3: Zero-shot speech synthesis with factorized codec and diffusion models,” in *ICLR*, 2024.
- [17] Yi Ren, Yangjun Ruan, et al., “Fastspeech: Fast, robust and controllable text to speech,” *Advances in neural information processing systems*, vol. 32, 2019.
- [18] Yi Ren, Chenxu Hu, Xu Tan, et al., “Fastspeech 2: Fast and high-quality end-to-end text to speech,” in *International Conference on Learning Representations*, 2020.
- [19] Jaehyeon Kim, Jungil Kong, and Juhee Son, “Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech,” in *ICLR*.
- [20] Jungil Kong, Jihoon Park, et al., “Vits2: Improving quality and efficiency of single-stage text-to-speech with adversarial learning and architecture design,” in *Interspeech*, 2023.
- [21] Shivam Mehta, Ruibo Tu, Jonas Beskow, and et al, “Matcha-TTS: A fast TTS architecture with conditional flow matching,” in *Proc. ICASSP*, 2024.
- [22] Chunyu Qiang and Wang et al Geng, “Vq-ctap: Cross-modal fine-grained sequence representation learning for speech processing,” *IEEE Transactions on Audio, Speech and Language Processing*, vol. 33, 2025.
- [23] Zengwei Yao, Liyong Guo, et al., “Zipformer: A faster and better encoder for automatic speech recognition,” in *The Twelfth International Conference on Learning Representations*, 2024.
- [24] Jonathan Ho and et al, “Classifier-free diffusion guidance,” *arXiv preprint arXiv:2207.12598*, 2022.
- [25] Hubert Siuzdak, “Vocos: Closing the gap between time-domain and fourier-based neural vocoders for high-quality audio synthesis,” in *ICLR*, 2023.
- [26] Sang Gil Lee et al., “Bigvgan: A universal neural vocoder with large-scale training,” in *ICLR*, 2023.
- [27] Haorui He, Zengqiang Shang, and et al, “Emilia: An extensive, multilingual, and diverse speech dataset for large-scale speech generation,” in *SLT*, 2024.
- [28] Yao Shi, Hui Bu, et al., “Aishell-3: A multi-speaker mandarin tts corpus,” in *Proc. Interspeech*, 2021.
- [29] Alec Radford et al., “Robust speech recognition via large-scale weak supervision,” in *ICLR*, 2023.
- [30] Zhifu Gao, ShiLiang Zhang, et al., “Paraformer: Fast and accurate parallel transformer for non-autoregressive end-to-end speech recognition,” in *Proc. Interspeech 2022*, 2022, pp. 2063–2067.
- [31] Brecht Desplanques et al., “Ecapa-tdnn: Emphasized channel attention, propagation and aggregation in tdnn based speaker verification,” in *Interspeech*, 2020.
- [32] Takaaki Saeki, Detai Xin, Wataru Nakata, and et al, “UTMOS: Utokyo-sarulab system for voicemos challenge 2022,” *Proc. Interspeech 2022*, 2022.