

VOCALNET-M2: ADVANCING LOW-LATENCY SPOKEN LANGUAGE MODELING VIA INTEGRATED MULTI-CODEBOOK TOKENIZATION AND MULTI-TOKEN PREDICTION

Yuhao Wang^{1,2}, Ziyang Cheng¹, Heyang Liu^{1,2}, Ronghua Wu², Qunshan Gu², Yanfeng Wang¹, Yu Wang^{1†}

¹Shanghai Jiao Tong University

²Ant Group

ABSTRACT

Current end-to-end spoken language models (SLMs) have made notable progress, yet they still encounter considerable response latency. This delay primarily arises from the autoregressive generation of speech tokens and the reliance on complex flow-matching models for speech synthesis. To overcome this, we introduce VocalNet-M2, a novel low-latency SLM that integrates a multi-codebook tokenizer and a multi-token prediction (MTP) strategy. Our model directly generates multi-codebook speech tokens, thus eliminating the need for a latency-inducing flow-matching model. Furthermore, our MTP strategy enhances generation efficiency and improves overall performance. Extensive experiments demonstrate that VocalNet-M2 achieves a substantial reduction in first chunk latency (from approximately 725ms to 350ms) while maintaining competitive performance across mainstream SLMs. This work also provides a comprehensive comparison of single-codebook and multi-codebook strategies, offering valuable insights for developing efficient and high-performance SLMs for real-time interactive applications.

Index Terms— spoken language models, multi-token prediction, multi-codebook

1. INTRODUCTION

In recent years, end-to-end spoken language models (SLMs) have made rapid progress, marking a significant milestone in generative AI [1, 2, 3, 4]. These models learn to directly model discrete speech tokens derived from a speech tokenizer, endowing Large Language Models (LLMs) with the dual ability to comprehend and generate both text and speech, which is particularly valuable for applications requiring natural and fluent human-computer dialogue, such as voice assistants and real-time conversational agents.

Current approaches for multimodal modeling of text and speech in SLMs can be broadly categorized into two paradigms. The first is the speech-native multimodal approach, which directly extends an LLM’s vocabulary to include tokens from speech corpora [2, 3, 4, 5, 6]. The second is the modality-alignment approach [7, 8, 9]. This method leverages pre-trained LLMs and integrates separate speech input/output modules through efficient cross-modal alignment. This reduces the need for extensive training from scratch, building upon existing LLM capabilities.

Despite these advancements, a significant challenge for current open-source SLMs is high response latency, which severely degrades user experience in speech interactions. This latency primarily arises from two sources: the autoregressive generation of both text and speech tokens by the Transformer decoder, and the subsequent conversion of speech tokens into waveforms. Most SLMs are trained to

generate a single codebook of semantic speech tokens [8, 9, 4, 5, 10]. These tokens are then converted into a mel spectrogram by a flow-matching model, which a vocoder finally synthesizes into an audio waveform [11]. While semantic tokens effectively capture linguistic content, their limited acoustic information necessitates a flow-matching model for detailed acoustic reconstruction. Although this simplifies the speech modeling task for the LLM and can produce high-quality speech, it introduces a critical bottleneck. The heavy reliance on the flow-matching model for acoustic reconstruction leads to substantial computational overhead and, critically, significant inference latency. This poses a major challenge for real-time interactive applications where low-latency turn-taking is essential.

To mitigate this, a direct approach is to empower SLMs to generate multi-codebook speech tokens. These tokens inherently contain richer acoustic information, which could eliminate the need for a separate, latency-inducing flow-matching model. While Moshi [12] has explored multi-codebook strategies, its application has been confined to the speech-native multimodal paradigm. Furthermore, a systematic analysis of the architectural design and a comprehensive comparison between single-codebook and multi-codebook strategies remain largely unexplored.

Therefore, in this work, we propose VocalNet-M2, a novel low-latency modality-alignment SLM incorporating a multi-codebook tokenizer and a multi-token prediction (MTP) strategy [13]. Our main contributions are threefold:

- We propose a novel modality-aligned spoken language model architecture capable of directly generating **multi-codebook speech tokens**, eliminating the need for a flow-matching model and enabling more streamlined and efficient speech response generation.
- We design a specialized **multi-token prediction strategy** tailored for multi-codebook generation, which significantly improves the performance of VocalNet-M2 while further reducing inference latency.
- We conduct extensive experiments comparing single-codebook and multi-codebook strategies, revealing their respective strengths and trade-offs. Our findings provide actionable insights for future research on efficient and high-performance spoken language models.

2. VOCALNET-M2

2.1. Model Architecture

VocalNet-M2 adopts a Thinker-Talker architecture [1], as illustrated in Figure 1 (Left).

[†]Corresponding author

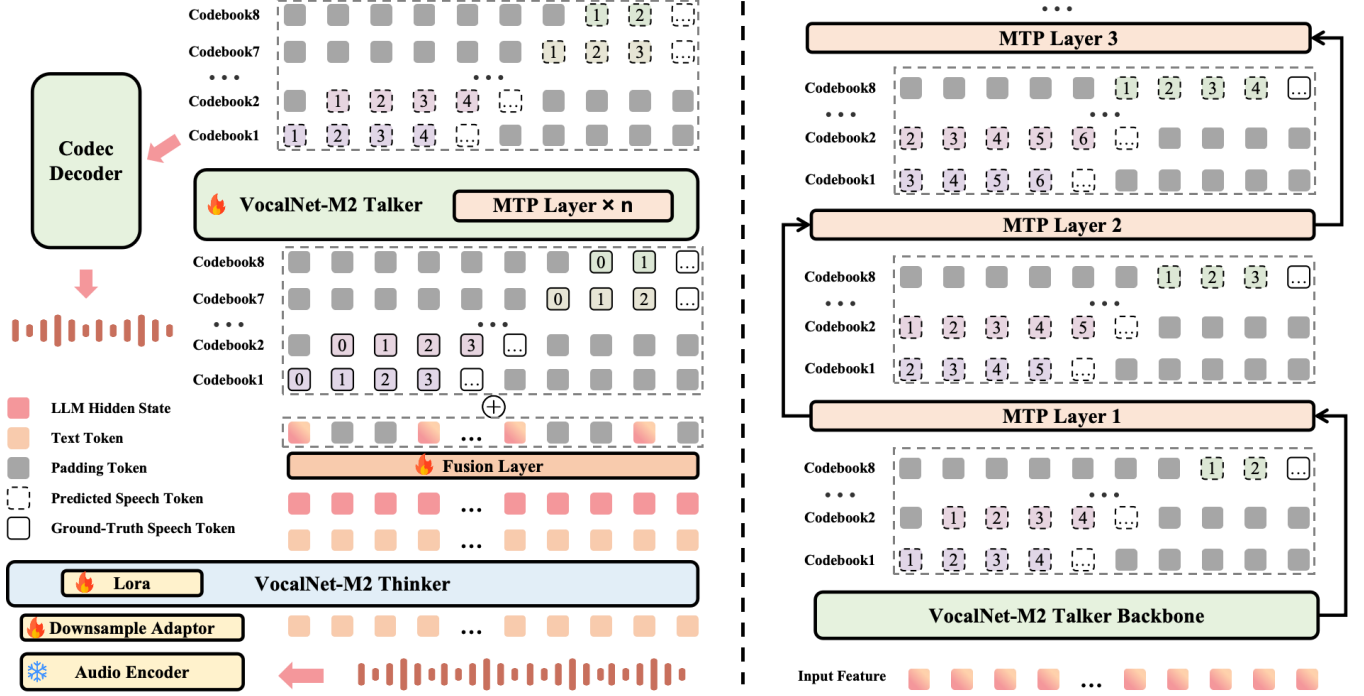


Fig. 1. Left: Overview of the VocalNet-M2 architecture. Right: The detailed architecture of VocalNet-M2 Talker.

Given a raw audio input x^a , it is first processed by an **Audio Encoder** and a **Downsample Adaptor**. This converts the raw audio into continuous representations $r_{1:T}^a$ that capture the high-quality features of the input. These representations are then fed into the **VocalNet-M2 Thinker**. The Thinker is responsible for autoregressively generating both the textual response tokens and their corresponding hidden states:

$$(t_{1:N}^{\text{text}}, h_{1:N}^{\text{text}}) = \mathcal{T}_{\text{thinker}}(r_{1:T}^a) \quad (1)$$

Here, $t_{1:N}^{\text{text}}$ represents the sequence of generated text tokens, and $h_{1:N}^{\text{text}}$ are their associated LLM hidden states, providing a rich semantic embedding for subsequent speech generation.

The **VocalNet-M2 Talker** is a multi-track autoregressive transformer decoder designed to generate speech tokens. VocalNet-M2 utilizes the XY-Tokenizer [14] to extract eight codebook audio tokens. Consequently, as shown in Figure 1, the Talker operates with eight distinct audio tracks (one for each codebook) and a semantic representation track as input. Its output consists of eight audio tokens, corresponding to the predicted tokens for each codebook.

Before being input to the Talker, the hidden states and text embeddings from the Thinker are processed by a **Fusion Layer**. This layer, composed of 2 linear layers, fuses the text embeddings and the Thinker's hidden states to create a unified semantic representation:

$$h_{1:N}^{\text{fused}} = \text{Linear}\left(\sigma\left(\text{Linear}(\text{Emb}(t_{1:N}^{\text{text}}) || h_{1:N}^{\text{text}})\right)\right) \quad (2)$$

Where $\text{Emb}(t_{1:N}^{\text{text}})$ denotes the embeddings of the text tokens, and $||$ signifies concatenation.

To address the inherent frequency mismatch between text and speech tokens, we first upsample the fused representation $h_{1:N}^{\text{fused}}$ to three times its original length. This upsampling aims to promote better temporal alignment between the semantic information and the

audio tokens. For the current decoding timestep t , this upsampled sequence is then either truncated or padded with zeros to match the current audio tokens length t , forming $h_{1:t}^{\text{up}}$:

$$h_{1:t}^{\text{up}} = \begin{cases} [h_1^{\text{fused}}, \mathbf{0}, \dots, h_N^{\text{fused}}, \mathbf{0}, \mathbf{0}]_{1:t} & \text{if } 3N \geq t \\ [h_1^{\text{fused}}, \mathbf{0}, \mathbf{0}, \dots, h_N^{\text{fused}}, \mathbf{0}, \mathbf{0}, \dots, \mathbf{0}]_{1:t} & \text{if } 3N < t \end{cases} \quad (3)$$

Then, the VocalNet-M2 Talker predicts the audio tokens at time $t + 1$. Its input consists of the upsampled hidden states ($h_{1:t}^{\text{up}}$) and the embeddings of all previously predicted audio tokens ($\sum_{j=1}^8 \text{Emb}(a_{1:t}^{\text{cbj}})$). The Talker then outputs the eight audio tokens for time $t + 1$ via eight linear layers:

$$\{a_{t+1}^{\text{cbj}}\}_{j=1}^8 = \mathcal{T}_{\text{talker}}(h_{1:t}^{\text{up}} + \sum_{j=1}^8 \text{Emb}(a_{1:t}^{\text{cbj}})) \quad (4)$$

This design allows VocalNet-M2 to generate audio tokens in a streaming manner by continuously feeding the h^{up} sequence into the Talker. During training, the cross-entropy loss is computed for each codebook at each time step t :

$$\mathcal{L}_{\text{talker}} = - \sum_{t=0}^{M-1} \sum_{j=1}^8 \log P(a_{t+1}^{\text{cbj}} | h_{1:t}^{\text{up}}, \{a_{1:t}^{\text{cbj}}\}_{i=1}^8) \quad (5)$$

Here, M is the length of audio token.

2.2. Multi-token Prediction (MTP)

To enhance generation efficiency and capture local dependencies more effectively [9], VocalNet-M2 incorporates a MTP mechanism. As shown in Figure 1 (Right), the VocalNet-M2 Talker consists of a Talker backbone followed by N_{mtp} sequential MTP layers. This

Table 1. Comparison between VocalNet-M2 and other mainstream SLMs. We evaluate their performance including text quality, speech quality, and first chunk latency. Note that Qwen2.5-Omni’s first chunk latency was not measured due to the absence of officially provided streaming inference code. The latency results are shown in the formation ‘mean±standard error’.

Model	text				speech		First chunk latency (ms)
	AlpacaEval	Llama Questions	TriviaQA	Web Questions	wer	utmos	
SLAM-Omni [5]	3.50	2.94	0.39	0.84	5.78	4.46	702.41 ± 30.30
VocalNet-8B [9]	7.12	7.95	6.24	6.48	3.64	4.49	556.00 ± 8.29
GLM-4-Voice [4]	5.86	7.74	4.95	5.56	11.90	4.23	1060.36 ± 2.36
MiniCPM-o [15]	6.13	7.72	6.43	7.16	9.52	4.14	893.82 ± 81.80
kimio-audio [16]	6.49	8.10	6.15	7.10	14.71	2.87	1744.80 ± 139.99
Qwen2.5-Omni [1]	6.01	7.90	5.89	6.88	2.31	4.34	\
VocalNet-M2	7.29	8.33	6.13	6.65	6.07	4.31	348.86 ± 2.86

design allows the model to predict $N_{\text{mtp}} + 1$ audio tokens for each codebook in a single inference step while preserving the essential temporal relationships between these tokens.

Specifically, for an MTP Layer n (where $n \in \{1, \dots, N_{\text{mtp}}\}$), it is designed to predict the audio tokens at time step $t + n + 1$, leveraging information available up to time t . The output of an MTP layer can be formally expressed as:

$$\{a_{t+n+1}^{\text{cbj}}\}_{j=1}^8 = \mathcal{T}_{\text{MTP}_n} \cdots \mathcal{T}_{\text{MTP}_1} \mathcal{T}_{\text{talker}}(h_{1:t}^{\text{up}} + \sum_{j=1}^8 \text{Emb}(a_{1:t}^{\text{cbj}})) \quad (6)$$

As analyzed in VocalNet [9], this approach helps the model to more efficiently leverage limited training data and better capture local dependencies between speech tokens. Consequently, the overall training objective integrates the standard Talker loss with losses from each MTP layer:

$$\mathcal{L}_{\text{mtp}} = - \sum_{n=0}^{N_{\text{mtp}}} \sum_{t=0}^{M-1} \sum_{j=1}^8 \log P(a_{t+n+1}^{\text{cbj}} | h_{1:t}^{\text{up}}, \{a_{1:t}^{\text{cbi}}\}_{i=1}^8) \quad (7)$$

2.3. Training Strategy

The training strategy for VocalNet-M2 is structured in three sequential stages. The initial stage focuses on pre-training the VocalNet-M2 Talker using TTS data, where the model learns to synthesize high-quality speech solely from text tokens, thereby establishing its fundamental speech generation capabilities. Following this, the second stage is dedicated to training the downsample adaptor and VocalNet-M2 Thinker with Lora, enabling it to comprehend and process raw audio inputs and generate text responses. The final stage integrates both the Thinker and Talker modules for end-to-end fine-tuning on speech-to-speech dialogue data. Distinct from the initial TTS pre-training, here the Talker module receives both the hidden states and text embeddings generated by the Thinker. This comprehensive final phase allows the entire VocalNet-M2 model to process audio input, produce a pertinent textual response, and concurrently generate the corresponding speech output.

3. EXPERIMENTAL SETTINGS

3.1. Training Data

For TTS pre-training, we used approximately 10k hours of randomly sampled audio from the Emilia dataset [17]. Our speech dialogue

training dataset totals about 800K samples (approximately 7k hours of audio). This includes 400K dialogues from VoiceAssistant [6], 300K from Ultrachat [9], and an additional 100K English multi-turn speech dialogues synthesized from tulu-3-sft-mixture [18] using Cosyvoice2 [11].

Regarding model initialization, the audio encoder is based on the Whisper-large-v3 [19]. The VocalNet-M2 Thinker is initialized from Qwen3-8B [20]. The VocalNet-M2 Talker shares a similar architectural design with the Thinker but employs a reduced number of transformer layers and is trained separately from scratch. As for audio labels, we utilize XY-tokenizer [14] to extract tokens.

3.2. Evaluation Metrics

To thoroughly assess VocalNet-M2’s voice interaction capabilities, we utilize the English subsets of OpenAudioBench [3]. The quality of textual responses generated by the model is evaluated using Qwen-max, which scores correctness and relevance on a normalized scale of 0 to 10. For evaluating speech quality, we employ two distinct metrics. UTMOS [21] predicts Mean Opinion Scores (MOS) for objective measurement of perceived naturalness. To quantify the alignment between synthesized speech and its text, we transcribe the speech using Whisper-large-v3 [19] and then compute the Word Error Rate (WER) against the ground-truth text.

4. EXPERIMENTAL RESULTS

4.1. Main Results

Table 1 presents a comparative analysis between VocalNet-M2 and other mainstream SLMs. VocalNet-M2 demonstrates strong performance in text quality, retaining the knowledge and reasoning capabilities of Qwen3-8B. It achieves the highest scores in AlpacaEval and Llama Questions, alongside competitive results in TriviaQA and Web Questions. For speech quality, evaluated using UTMOS and WER, VocalNet-M2’s performance is in the mid-range among current mainstream models, as shown in Table 1. This is an expected outcome given the limited training data and the inherent complexity of learning multi-codebook speech tokens compared to single-codebook approaches.

Furthermore, we conducted a latency analysis, measuring the first audio chunk generation time for all models. To ensure a fair comparison and minimize variability arising from custom implementations, we prioritized models with officially provided streaming inference code. Consequently, Qwen2.5-Omni was excluded due to the absence of such code. Most models were evaluated using

Table 2. Impact of tokenizer type (single vs. multi-codebook) and training data quality on speech generation performance.

tokenizer	Training data	WER	utmos	First chunk latency
\mathcal{S}^3 tokenizer(Single-codebook)	v1	10.66	4.34	725.90 \pm 9.17
	emilia + v1	3.73	4.35	
	emilia + v2	3.68	4.37	
XY-tokenizer(Multi-codebook)	v1	20.49	3.89	405.23 \pm 6.29
	emilia + v1	10.43	4.08	
	emilia + v2	8.56	4.24	

a fixed first chunk duration of 0.8 seconds; the sole exception was MiniCPM-o, whose official implementation has a fixed chunk size of 0.533 seconds. All tests were performed on a single L20 GPU without acceleration frameworks (e.g., vLLM), as these are not universally supported by the baseline models. As detailed in Table 1, VocalNet-M2 exhibits significantly lower first chunk latency while maintaining competitive performance. This efficiency is attributed to its direct modeling of multi-codebook speech tokens and the incorporation of the MTP approach.

4.2. Comparison Between Single and Multi-Codebook Tokens

This section explores the differences in modeling single and multi-codebook speech tokens. For the tokenizer with single-codebook, we utilize the \mathcal{S}^3 tokenizer from Cosyvoice2 [11]. Conversely, the XY-tokenizer with 8 codebooks is utilized for extracting multi-codebook tokens. Two versions of training data were constructed:

- **v1 data:** The speech dialogue data described in Section 3.1.
- **v2 data:** A high-quality dialogue dataset of approximately 400K samples, derived from v1. It exhibits higher UTMOS and lower WER, achieved by filtering samples with high WER and re-synthesizing audio with high-quality prompts.

The results of this ablation study are presented in Table 2.

Firstly, we observe that learning to generate multi-codebook speech tokens requires more training data than single-codebook tokens to achieve comparable performance. Without TTS pretraining, the model struggles to generate the multi-codebook tokens, resulting in very high WER and low UTMOS scores. In contrast, learning speech tokens extracted by the \mathcal{S}^3 tokenizer is easier; even without pretraining, the model can achieve respectable performance.

Secondly, high-quality training data is crucial for learning to generate multi-codebook speech tokens. Models trained with ‘emilia + v2’ data demonstrate significantly better WER and UTMOS scores compared to those trained with ‘emilia + v1’ data. This improvement is not observed in models leveraging single codebook speech tokens. The reason for this disparity is that single codebook speech tokens inherently lack acoustic information, which must be reconstructed through a flow-matching model to convert them back into speech. Consequently, the quality of speech generated from single codebook tokens heavily relies on the performance of the flow-matching model. In this work, we utilize the CosyVoice2 flow-matching model, enabling single-codebook-based models to achieve high performance in both WER and UTMOS, even when trained on noisier, lower-quality ‘v1’ data. In contrast, models based on multi-codebook tokens directly learn acoustic features from the dialogue data, making them more sensitive to the quality of the training data.

4.3. Ablation study on MTP

This section investigates the impact of MTP on the model’s performance. In this experiment, we varied the number of MTP layers used

Table 3. Ablation study on the impact of MTP on model’s performance.

Metrics	w/o MTP Layer	w/ n MTP layer				
		n=1	n=2	n=3	n=4	n=5
WER	8.56	7.64	6.53	6.64	6.07	6.33
UTMOS	4.24	4.28	4.29	4.28	4.31	4.29

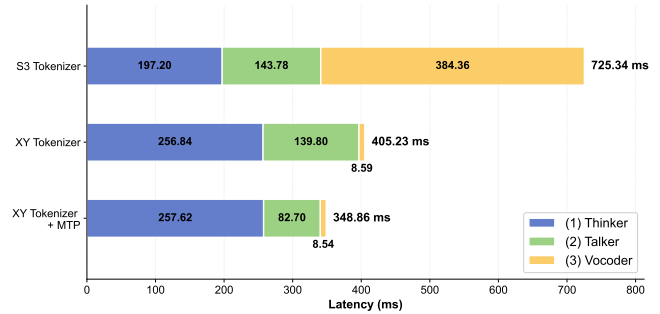


Fig. 2. Breakdown of first chunk latency across different model components for various configurations.

during training, while MTP layers were not employed during inference. As shown in Table 3, the introduction of MTP significantly reduced the WER from 8.56 (without MTP layers) to an optimal 6.07 with four MTP layers. UTMOS scores remained consistently high, with the best score of 4.31 also achieved with four MTP layers. Therefore, a configuration with four MTP layers was adopted for VocalNet-M2.

4.4. Latency Analysis

We categorize the first chunk latency into three distinct parts: (1) the Thinker’s generation of text tokens and hidden states, (2) the Talker’s generation of speech tokens, and (3) the Vocoder’s conversion of these tokens into speech. Figure 2 illustrates the latency for various configurations. Notably, the integration of multi-codebook speech tokens and the MTP method significantly reduced the latency in both the (2) Talker and (3) Vocoder stages. These advancements collectively enabled VocalNet-M2 to achieve a first chunk latency reduction from 725 ms to 348 ms, resulting in an inference speedup of approximately 2 \times .

5. CONCLUSION

In this work, we introduced VocalNet-M2, a novel low-latency modality-alignment SLM. Our key contributions include a new model architecture designed to directly generate multi-codebook speech tokens, thereby eliminating the need for a computationally intensive flow-matching model for speech synthesis and significantly reducing latency. Additionally, we developed a MTP strategy that not only enhances overall performance but also further reduces inference latency. Our experimental results demonstrate that VocalNet-M2 effectively reduces first chunk latency by approximately 50%, from 725ms to 348ms, showcasing a substantial improvement in responsiveness. We also provided a thorough comparison between single and multi-codebook approaches, highlighting their respective advantages and limitations. These advancements pave the way for more efficient and responsive spoken dialogue systems.

6. REFERENCES

- [1] Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, et al., “Qwen2. 5-omni technical report,” *arXiv preprint arXiv:2503.20215*, 2025.
- [2] Boyong Wu, Chao Yan, Chen Hu, Cheng Yi, Chengli Feng, Fei Tian, Feiyu Shen, Gang Yu, Haoyang Zhang, Jingbei Li, et al., “Step-audio 2 technical report,” *arXiv preprint arXiv:2507.16632*, 2025.
- [3] Yadong Li, Jun Liu, Tao Zhang, Song Chen, Tianpeng Li, Zehuan Li, Lijun Liu, Lingfeng Ming, Guosheng Dong, Da Pan, et al., “Baichuan-omni-1.5 technical report,” *arXiv preprint arXiv:2501.15368*, 2025.
- [4] Aohan Zeng, Zhengxiao Du, Mingdao Liu, Kedong Wang, Shengmin Jiang, Lei Zhao, Yuxiao Dong, and Jie Tang, “Glm-4-voice: Towards intelligent and human-like end-to-end spoken chatbot,” *arXiv preprint arXiv:2412.02612*, 2024.
- [5] Wenxi Chen, Ziyang Ma, Ruiqi Yan, Yuzhe Liang, Xiquan Li, Ruiyang Xu, Zhikang Niu, Yanqiao Zhu, Yifan Yang, Zhanxun Liu, et al., “Slam-omni: Timbre-controllable voice interaction system with single-stage training,” *arXiv preprint arXiv:2412.15649*, 2024.
- [6] Zhifei Xie and Changqiao Wu, “Mini-omni: Language models can hear, talk while thinking in streaming,” *arXiv preprint arXiv:2408.16725*, 2024.
- [7] Wenyi Yu, Siyin Wang, Xiaoyu Yang, Xianzhao Chen, Xiaohai Tian, Jun Zhang, Guangzhi Sun, Lu Lu, Yuxuan Wang, and Chao Zhang, “Salmonn-omni: A standalone speech llm without codec injection for full-duplex conversation,” *arXiv preprint arXiv:2505.17060*, 2025.
- [8] Qian Chen, Yafeng Chen, Yanni Chen, Mengzhe Chen, Yingda Chen, Chong Deng, Zhihao Du, Ruize Gao, Changfeng Gao, Zhifu Gao, et al., “Minmo: A multimodal large language model for seamless voice interaction,” *arXiv preprint arXiv:2501.06282*, 2025.
- [9] Yuhao Wang, Heyang Liu, Ziyang Cheng, Ronghua Wu, Qunshan Gu, Yanfeng Wang, and Yu Wang, “Vocalnet: Speech llm with multi-token prediction for faster and high-quality generation,” *arXiv preprint arXiv:2504.04060*, 2025.
- [10] Qingkai Fang, Yan Zhou, Shoutao Guo, Shaolei Zhang, and Yang Feng, “Llama-omni2: Llm-based real-time spoken chatbot with autoregressive streaming speech synthesis,” *arXiv preprint arXiv:2505.02625*, 2025.
- [11] Zhihao Du, Yuxuan Wang, Qian Chen, Xian Shi, Xiang Lv, Tianyu Zhao, Zhifu Gao, Yexin Yang, Changfeng Gao, Hui Wang, et al., “Cosyvoice 2: Scalable streaming speech synthesis with large language models,” *arXiv preprint arXiv:2412.10117*, 2024.
- [12] Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard Grave, and Neil Zeghidour, “Moshi: a speech-text foundation model for real-time dialogue,” *arXiv preprint arXiv:2410.00037*, 2024.
- [13] Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve, “Better & faster large language models via multi-token prediction,” *arXiv preprint arXiv:2404.19737*, 2024.
- [14] Yitian Gong, Luo Zhijie Jin, Ruifan Deng, Dong Zhang, Xin Zhang, Qinyuan Cheng, Zhaoye Fei, Shimin Li, and Xipeng Qiu, “Xy-tokenizer: Mitigating the semantic-acoustic conflict in low-bitrate speech codecs,” *arXiv preprint arXiv:2506.23325*, 2025.
- [15] OpenBMB, “Minicpm-o 2.6: A gpt-4o level mllm for vision, speech, and multimodal live streaming on your phone,” <https://openbmb.notion.site/185ede1b7a558042b5d5e45e6b237da9>, Accessed: 2025-03-28.
- [16] Ding Ding, Zeqian Ju, Yichong Leng, Songxiang Liu, Tong Liu, Zeyu Shang, Kai Shen, Wei Song, Xu Tan, Heyi Tang, et al., “Kimi-audio technical report,” *arXiv preprint arXiv:2504.18425*, 2025.
- [17] Haorui He, Zengqiang Shang, Chaoren Wang, Xuyuan Li, Yicheng Gu, Hua Hua, Liwei Liu, Chen Yang, Jiaqi Li, Peiyang Shi, et al., “Emilia: An extensive, multilingual, and diverse speech dataset for large-scale speech generation,” in *2024 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2024, pp. 885–890.
- [18] Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al., “Tulu 3: Pushing frontiers in open language model post-training,” *arXiv preprint arXiv:2411.15124*, 2024.
- [19] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever, “Robust speech recognition via large-scale weak supervision,” in *International conference on machine learning*. PMLR, 2023, pp. 28492–28518.
- [20] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu, “Qwen3 technical report,” *arXiv preprint arXiv:2505.09388*, 2025.
- [21] Takaaki Saeki, Detai Xin, Wataru Nakata, Tomoki Koriyama, Shinnosuke Takamichi, and Hiroshi Saruwatari, “Utmos: Utokyo-sarulab system for voicemos challenge 2022,” *arXiv preprint arXiv:2204.02152*, 2022.