# Multi-Reward GRPO for Stable and Prosodic Single-Codebook TTS LLMs at Scale

Yicheng Zhong
ajaxzhong@tencent.com
Tencent Technology Co.Ltd
Shenzhen, Guangdong, China

Peiji Yang
Tencent Technology Co.Ltd
Shenzhen, Guangdong, China
peijiyang@tencent.com

Zhisheng Wang
Tencent Technology Co.Ltd
Shenzhen, Guangdong, China
plorywang@tencent.com

## Abstract

Recent advances in Large Language Models (LLMs) have transformed text-to-speech (TTS) synthesis, inspiring autoregressive frameworks that represent speech as sequences of discrete codec tokens. Among them, single-codebook TTS LLMs have emerged as compact and streamable architectures that jointly model semantic and acoustic integration. However, despite their efficiency, these models often exhibit unstable prosody, speaker drift, and degraded naturalness. To address these issues, we propose a multi-reward Group Relative Policy Optimization (GRPO) framework that directly optimizes the token generation policy of single-codebook TTS LLMs. Beyond standard intelligibility and speaker similarity objectives, our design integrates three rule-based rewards: a length penalty for duration consistency, an entropy regularization reward for decoding stability, and an LLM-annotated prosody alignment reward that explicitly supervises rhythm. In this prosody reward, an external reasoning LLM predicts multiple plausible pause structures via in-context learning, providing a human-preference-aligned supervisory signal for GRPO training. To assess universality, we further attach a flow-matching (FM) decoder on top of the GRPO-optimized AR backbone and observe consistent additional gains, indicating that our reinforcement optimization enhances the intrinsic AR policy. We further conduct a scalability analysis across data sizes and model scales, revealing that the proposed method consistently enhances prosodic stability, speaker similarity, and overall speech naturalness in single-codebook TTS LLMs.

## CCS Concepts

• **Computing methodologies → Natural language processing**.

## Keywords

Text-to-speech synthesis, Reinforecement learning, Large language model

## 1 Introduction

Recent advances in Large Language Models (LLMs) have reshaped text-to-speech (TTS) synthesis, enabling autoregressive (AR) decoding over discrete codec tokens. Among these, single-codebook TTS LLMs stand out for their compactness and native streaming capability. Currently, zero-shot TTS systems span three main families: (1) LLM-based acoustic-token models with strong linguistic–acoustic modeling [18, 19]; (2) diffusion-based architectures that implicitly learn text–speech alignment with fine-grained control [4]; and (3) coarse-to-fine pipelines where AR LLMs predict semantic tokens refined by diffusion or flow-matching modules [1, 5–7]. Within the single-codebook paradigm, semantic-focused systems rely on secondary models to recover acoustic detail, whereas joint semantic–acoustic systems directly generate tokens encoding both linguistic and paralinguistic cues [20]. While the latter offers more expressivity and lower latency, unified modeling often introduces suboptimal decoding policies, leading to prosody instability, speaker drift, and weakened temporal controllability. Reinforcement learning (RL) provides a natural mechanism to directly optimize these AR policies but remains underexplored.

Prior RL efforts include uncertainty aware learning [2] and synthetic positive construction via reverse inference optimization [12]. Differentiable reward frameworks such as DiffRO [8] further enable supervised optimization with programmatic objectives, though DPO-style methods are sensitive to preference noise and costly to scale. In contrast, Group Relative Policy Optimization (GRPO) [17] stabilizes learning through group-wise advantage normalization and avoids dependence on dense preference labels.

We introduce a multi-reward GRPO framework that integrates objective metrics with three rule-based rewards: a length penalty for duration consistency, an entropy reward for stable decoding, and an LLM-annotated Prosody Alignment Reward for explicit rhythm supervision. To obtain prosodic templates, a reasoning LLM (e.g., DeepSeek-R1 [9]) generates multiple plausible pause patterns via in-context learning offline; these templates then guide online optimization toward natural rhythm. We further conduct a systematic scaling study across 1K–1M data and 1B–8B models, revealing a strong correlation between RL effectiveness and data scale. Experiments show robust gains in prosodic stability, speaker similarity, and naturalness. Moreover, adding a Flow Matching refinement module after RL continues to yield improvements, confirming that our method strengthens the intrinsic AR policy in a manner complementary to acoustic refiners.
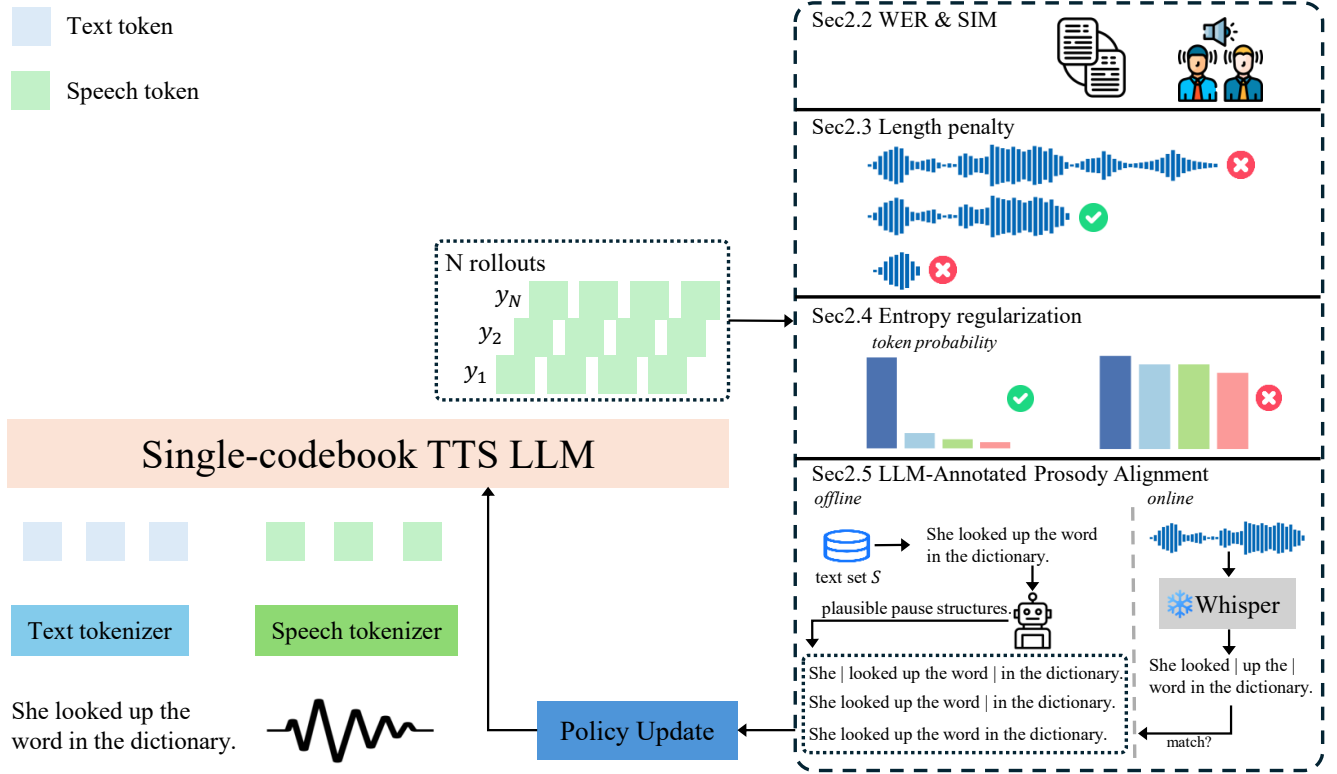
**Figure 1: Overview of our framework. Multiple rollouts are evaluated with WER/SIM, length penalty, entropy regularization, and an LLM-annotated prosody alignment reward to guide policy updates.**

## 2 Methodology

### 2.1 Overview

We propose a reinforcement learning framework based on Group Relative Policy Optimization (GRPO) to enhance the stability, prosody, and speaker similarity of single-codebook TTS large language models (LLMs). Starting from a single-codebook LLM as a policy $\pi_\theta(a_t|s_t)$, we employ GRPO to optimize the model with multiple complementary rewards that reflect both objective and rule-based criteria. The objective is to maximize the expected cumulative reward across the generated trajectory $\tau = (s_1, a_1, \ldots, s_T, a_T)$:

$$\mathcal{J}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=1}^{T} R(s_t, a_t) \right]. \tag{1}$$

The reward function $R$ is decomposed into several interpretable components that capture different aspects of synthesis quality:

$$R(s_t, a_t) = \alpha_{intl}R_{intl} + \alpha_{sim}R_{sim} + \alpha_{len}R_{len} + \alpha_{ent}R_{ent} + \alpha_{pro}R_{pro}, \tag{2}$$

where $\alpha_i$ are tunable scaling coefficients that balance the contribution of different reward terms. The policy parameters $\theta$ are updated via GRPO to maximize $\mathcal{J}(\theta)$ under the combined signal.

### 2.2 Intelligibility and Speaker Similarity Rewards

*Intelligibility Reward ($R_{intl}$).* To measure intelligibility and fidelity, we use the pre-trained Whisper ASR model[15] to transcribe the generated audio $A$ into text $\hat{S}$, and compute the Character Error Rate (CER) / Word Error Rate (WER) against the input text $S$:

$$R_{intl} = 1 - \frac{D_{lev}(\hat{S}, S)}{|S|}, \tag{3}$$

where $D_{lev}$ denotes the Levenshtein distance. This reward encourages the model to produce speech that is semantically consistent with the input text.

*Speaker Similarity Reward.* To evaluate speaker consistency, we adopt the WavLM-large model[3] fine-tuned for speaker verification to extract speaker embeddings $E(A)$ and $E(A_{\text{ref}})$ from the generated and reference audios. We compute the cosine similarity between embeddings:

$$R_{sim} = \cos(E(A), E(A_{\text{ref}})) = \frac{E(A) \cdot E(A_{\text{ref}})}{\|E(A)\|\|E(A_{\text{ref}})\|}. \tag{4}$$

### 2.3 Length Penalty Reward

To prevent premature stopping or overly long generations—a common instability issue in AR TTS—we constrain the generated speech duration $T$ relative to a target $T_{\text{target}}$. The target is estimated from

the reference text length and the reference speech speed ratio. Given a tolerance range $[a, b]$, the reward is defined as:

$$R_{len} = \begin{cases} 1, & \text{if } \frac{T_{\text{text}}/T}{r_{\text{ref}}} \in [a, b], \\ 0, & \text{otherwise,} \end{cases} \quad (5)$$

where $r_{\text{ref}}$ denotes the reference speaking rate computed from the paired reference input, which stabilizes the generation length distribution and mitigates truncation or excessive elongation.

## 2.4 Entropy Regularization Reward

To encourage stable and deterministic token generation, we regularize the policy entropy along the generation trajectory. Let $\bar{H}$ denote the average token-level entropy across the sequence, and $H_{\text{target}}$ be the entropy target estimated from high-quality samples. The entropy reward is formulated as:

$$R_{ent} = -\lambda_{ent} \cdot \max(0, \bar{H} - H_{\text{target}}). \quad (6)$$

This reward penalizes excessively high entropy that leads to erratic prosodic variations, encouraging smoother generation paths.

## 2.5 LLM-Annotated Prosody Alignment Reward

*Offline Annotation.* To capture human-like prosody, we leverage an auxiliary reasoning LLM (*DeepSeek-R1*) to annotate a set of input texts $S$ with appropriate pause structures. Few-shot examples of prosodic annotation are provided to the LLM for both Chinese and English. For Chinese, we adopt discrete pause markers (#1–#4) corresponding to increasing pause durations. For English, we employ *Prosodic Word (PW)* and *Prosodic Phrase (PPH)* labels, following the linguistic hierarchy in [16]. The resulting pseudo-labels serve as target prosodic structures.

*Online Comparison.* During GRPO training, the generated waveform is decoded and timestamped using Whisper. We convert silence durations into discrete pause symbols through a handcrafted rule-based mapping, and compare them against the pseudo-labels. If the predicted pause sequence matches the annotated pattern, a binary reward is assigned:

$$R_{pro} = \begin{cases} 1, & \text{if } \hat{P}(A) \in \{P(S)\}, \\ 0, & \text{otherwise,} \end{cases} \quad (7)$$

where $\hat{P}(A)$ denotes the predicted pause structure from generated audio, and $\{P(S)\}$ is the annotated reference pattern set. This binary reward encourages alignment with human-preferred pause structures and suppresses unnatural rhythm patterns.

## 3 Experiments

## 3.1 Implementation Details

All experiments are conducted on **8×H20 GPUs** with mixed-precision training. The GRPO optimization adopts a batch size of **16**, a learning rate of **1e−6**, and a group size of **12**. During generation, we employ the vLLM decoding configuration with *top-k = 75*, *top-p = 0.9*, *temperature = 1.1*, and *repetition penalty = 1.1*. The reward coefficients are empirically set as:

$$\alpha_{\text{intl}} = \alpha_{\text{sim}} = \alpha_{\text{ent}} = \alpha_{\text{pro}} = 1.0, \quad \alpha_{\text{len}} = 0.1.$$

**Table 1: Comparison of different methods on SEED test sets (test-zh, test-en). For each subset we report WER (↓ lower is better), SIM (↑ higher is better).**

| Method | test-zh | | test-en | | test-hard | | |
|---|---|---|---|---|---|---|---|
| | CER↓ | SIM↑ | WER↓ | SIM↑ | CER↓ | SIM↑ | MOS↑ |
| Seed-TTS | 1.12 | **0.796** | 2.62 | <u>0.714</u> | 7.59 | **0.776** | - |
| FireRedTTS | 1.51 | 0.635 | 3.82 | 0.460 | 17.45 | 0.621 | 3.53 |
| MaskGCT | 2.27 | 0.774 | 2.62 | <u>0.714</u> | 10.27 | 0.748 | - |
| F5-TTS | 1.56 | 0.741 | **1.83** | 0.615 | 8.67 | 0.713 | 3.94 |
| Spark-TTS | 1.20 | 0.672 | <u>1.98</u> | 0.584 | - | - | 4.01 |
| CosyVoice | 3.63 | 0.723 | 4.29 | 0.609 | 11.75 | 0.709 | 3.89 |
| CosyVoice2 | 1.45 | 0.748 | 2.57 | 0.652 | 6.83 | 0.724 | 3.98 |
| CosyVoice3 | 1.12 | 0.781 | 2.21 | 0.720 | **5.83** | 0.758 | 4.07 |
| LLaSA-8B | 1.59 | 0.684 | 2.97 | 0.574 | 11.09 | 0.660 | 3.67 |
| LLaSA+SFT | 1.51 | 0.688 | 2.89 | 0.582 | 10.63 | 0.674 | 3.76 |
| LLaSA+RL | <u>1.10</u> | 0.758 | 2.12 | 0.672 | 6.04 | 0.731 | <u>4.12</u> |
| LLaSA+RL+FM | **1.08** | <u>0.790</u> | 2.08 | **0.733** | <u>5.98</u> | <u>0.775</u> | **4.21** |

For training data, we construct a bilingual corpus by sampling from Emilia[11] and libriheavy[13], selecting 1 million text samples (10–100 words) and 1 million speech samples balanced across Chinese and English (1:1 ratio). These are paired into approximately 1 million (ref_text, ref_speech, target_text) triplets for online GRPO training, totaling about 5115 hours of audio.

## 3.2 Main Results

Table 1 reports results on the SEED[1] benchmark, evaluating both objective metrics (CER/WER, SIM) and subjective naturalness through MOS, where 100 randomly sampled utterances were rated by 10 participants. We compare against several systems including Seed-TTS [1], FireRedTTS [10], MaskGCT [19], F5-TTS[4], Spark-TTS[18], CosyVoice[5–7], and the LLaSA[20] baseline, an SFT-only variant, and a post-RL Flow Matching (FM)[14] refinement model. Our GRPO-optimized LLaSA achieves the best CER on *test-zh* and competitive results across languages, outperforming all open-source single-codebook TTS LLMs (e.g., Spark-TTS, CosyVoice/2) and the SFT-only counterpart, showing the higher sample efficiency of RL relative to supervised fine-tuning. Despite CosyVoice3 benefiting from a hybrid architecture and substantially larger training data (1M h vs. our 250k h), our model still attains lower CER and comparable SIM.

On the challenging *test-hard* split, GRPO delivers notable gains in both CER and SIM, demonstrating improved robustness under difficult scenarios. Importantly, our method also achieves the highest MOS, indicating strong alignment with human preference. Adding FM after RL further improves performance, especially SIM and MOS, confirming that our RL optimization strengthens the intrinsic AR policy in a manner complementary to acoustic refinement.

## 3.3 Scalability Analysis

To investigate the effect of model and data scale, we apply the GRPO framework to models with 1B, 3B, and 8B parameters, trained with data scales ranging from 1K → 10K → 100K → 1M samples. Figure 2 visualizes the results, plotting *CER* and *SIM* against data scale.

We observe a clear monotonic trend: larger model capacities yield better prosodic stability and speaker preservation, while increasing GRPO data size further refines rhythm and decoding stability. Even
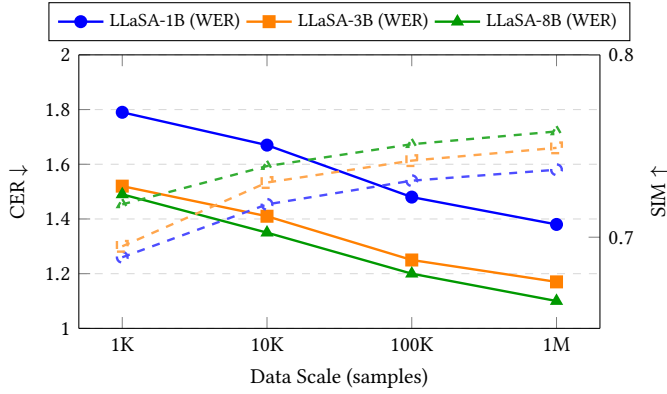
Figure 2: Scalability analysis. WER (solid, left axis) decreases while SIM (dashed, right axis) increases as both data scale and model size grow.

Table 2: Ablation study on the contribution of different reward components. Lower CER/WER and higher SIM/MOS indicate better performance.

| Method | zh | | | en | | |
|---|---|---|---|---|---|---|
| | CER↓ | SIM↑ | MOS↑ | WER↑ | SIM↓ | MOS↑ |
| LLaSA | 1.59 | 0.684 | 3.68 | 2.97 | 0.574 | 3.57 |
| + $R_{intl}$ & $R_{sim}$ | 1.31 | 0.719 | 3.77 | 2.66 | 0.623 | 3.69 |
| + $R_{len}$ | 1.23 | 0.738 | 3.81 | 2.48 | 0.647 | 3.75 |
| + $R_{ent}$ | 1.12 | 0.751 | 4.01 | 2.25 | 0.668 | 3.90 |
| + $R_{pro}$ (Full) | 1.10 | 0.758 | 4.25 | 2.12 | 0.672 | 4.12 |

small-scale RL (e.g., 10K samples) delivers measurable gains over the supervised baseline, highlighting the data efficiency of our proposed optimization scheme.

## 3.4 Ablation Study

Table 2 reports the incremental effect of each reward term across both objective (CER/WER, SIM) and subjective (MOS) metrics. $R_{intl}$ & $R_{sim}$ establish strong initial gains, while the length penalty further reduces duration mismatch. $R_{ent}$ yields notable improvements in both stability and MOS, indicating smoother and more natural token dynamics. $R_{pro}$ delivers the largest additional boost across all metrics—especially in MOS—showing that explicit rhythmic supervision not only improves prosodic structure but also better aligns the model's outputs with human perceptual preferences.

## 4 Conclusion

We introduced a GRPO-based RL framework that directly improves the intrinsic autoregressive policy of single-codebook TTS LLMs. By integrating objective, rule-based, and LLM-assisted prosody rewards, our method enhances prosodic stability, speaker similarity, and naturalness across model and data scales. These gains further compound with flow-matching refinement, indicating that our RL optimization strengthens the core AR decoder rather than overlapping with downstream token-refinement methods.

## References

[1] Philip Anastassiou, Jiawei Chen, Jitong Chen, Yuanzhe Chen, Zhuo Chen, Ziyi Chen, Jian Cong, Lelai Deng, Chuang Ding, Lu Gao, et al. 2024. Seed-tts: A family of high-quality versatile speech generation models. *arXiv preprint arXiv:2406.02430* (2024).

[2] Chen Chen, Yuchen Hu, Wen Wu, Helin Wang, Eng Siong Chng, and Chao Zhang. 2024. Enhancing zero-shot text-to-speech synthesis with human feedback. *arXiv preprint arXiv:2406.00654* (2024).

[3] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. 2022. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing* 16, 6 (2022), 1505–1518.

[4] Yushen Chen, Zhikang Niu, Ziyang Ma, Keqi Deng, Chunhui Wang, JianZhao JianZhao, Kai Yu, and Xie Chen. 2025. F5-tts: A fairytaler that fakes fluent and faithful speech with flow matching. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 6255–6271.

[5] Zhihao Du, Qian Chen, Shiliang Zhang, Kai Hu, Heng Lu, Yexin Yang, Hangrui Hu, Siqi Zheng, Yue Gu, Ziyang Ma, et al. 2024. Cosyvoice: A scalable multilingual zero-shot text-to-speech synthesizer based on supervised semantic tokens. *arXiv preprint arXiv:2407.05407* (2024).

[6] Zhihao Du, Changfeng Gao, Yuxuan Wang, Fan Yu, Tianyu Zhao, Hao Wang, Xiang Lv, Hui Wang, Chongjia Ni, Xian Shi, et al. 2025. Cosyvoice 3: Towards in-the-wild speech generation via scaling-up and post-training. *arXiv preprint arXiv:2505.17589* (2025).

[7] Zhihao Du, Yuxuan Wang, Qian Chen, Xian Shi, Xiang Lv, Tianyu Zhao, Zhifu Gao, Yexin Yang, Changfeng Gao, Hui Wang, et al. 2024. Cosyvoice 2: Scalable streaming speech synthesis with large language models. *arXiv preprint arXiv:2412.10117* (2024).

[8] Changfeng Gao, Zhihao Du, and Shiliang Zhang. 2025. Differentiable Reward Optimization for LLM based TTS system. In *Proc. Interspeech 2025*. 2450–2454.

[9] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).

[10] Hao-Han Guo, Yao Hu, Kun Liu, Fei-Yu Shen, Xu Tang, Yi-Chen Wu, Feng-Long Xie, Kun Xie, and Kai-Tuo Xu. 2024. Fireredtts: A foundation text-to-speech framework for industry-level generative speech applications. *arXiv preprint arXiv:2409.03283* (2024).

[11] Haorui He, Zengqiang Shang, Chaoren Wang, Xuyuan Li, Yicheng Gu, Hua Hua, Liwei Liu, Chen Yang, Jiaqi Li, Peiyang Shi, et al. 2024. Emilia: An extensive, multilingual, and diverse speech dataset for large-scale speech generation. In *2024 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 885–890.

[12] Yuchen Hu, Chen Chen, Siyin Wang, Eng Siong Chng, and Chao Zhang. 2024. Robust zero-shot text-to-speech synthesis with reverse inference optimization. *arXiv preprint arXiv:2407.02243* (2024).

[13] Wei Kang, Xiaoyu Yang, Zengwei Yao, Fangjun Kuang, Yifan Yang, Liyong Guo, Long Lin, and Daniel Povey. 2024. Libriheavy: A 50,000 hours ASR corpus with punctuation casing and context. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 10991–10995.

[14] Shivam Mehta, Ruibo Tu, Jonas Beskow, Éva Székely, and Gustav Eje Henter. 2024. Matcha-TTS: A fast TTS architecture with conditional flow matching. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 11341–11345.

[15] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust Speech Recognition via Large-Scale Weak Supervision. doi:10.48550/ARXIV.2212.04356

[16] Elisabeth O Selkirk. 1980. On prosodic structure and its relation to syntactic structure. *(No Title)* (1980).

[17] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300* (2024).

[18] Xinsheng Wang, Mingqi Jiang, Ziyang Ma, Ziyu Zhang, Songxiang Liu, Linqin Li, Zheng Liang, Qixi Zheng, Rui Wang, Xiaoqin Feng, et al. 2025. Spark-tts: An efficient llm-based text-to-speech model with single-stream decoupled speech tokens. *arXiv preprint arXiv:2503.01710* (2025).

[19] Yuancheng Wang, Haoyue Zhan, Liwei Liu, Ruihong Zeng, Haotian Guo, Ji-achen Zheng, Qiang Zhang, Xueyao Zhang, Shunsi Zhang, and Zhizheng Wu. 2025. MaskGCT: Zero-Shot Text-to-Speech with Masked Generative Codec Transformer. In *ICLR*.

[20] Zhen Ye, Xinfa Zhu, Chi-Min Chan, Xinsheng Wang, Xu Tan, Jiahe Lei, Yi Peng, Haohe Liu, Yizhu Jin, Zheqi Dai, et al. 2025. Llasa: Scaling train-time and inference-time compute for llama-based speech synthesis. *arXiv preprint arXiv:2502.04128* (2025).