

THINKING IN COCKTAIL PARTY: CHAIN-OF-THOUGHT AND REINFORCEMENT LEARNING FOR TARGET SPEAKER AUTOMATIC SPEECH RECOGNITION

Yiru Zhang[†], Hang Su[†], Lichun Fan^{*}, Zhenbo Luo, Jian Luan

MiLM Plus, Xiaomi Inc., China

ABSTRACT

Target Speaker Automatic Speech Recognition (TS-ASR) aims to transcribe the speech of a specified target speaker from multi-speaker mixtures in cocktail party scenarios. Recent advancement of Large Audio-Language Models (LALMs) has already brought some new insights to TS-ASR. However, significant room for optimization remains for the TS-ASR task within the LALMs architecture. While Chain of Thoughts (CoT) and Reinforcement Learning (RL) have proven effective in certain speech tasks, TS-ASR, which requires the model to deeply comprehend speech signals, differentiate various speakers, and handle overlapping utterances is particularly well-suited to a reasoning-guided approach. Therefore, we propose a novel framework that incorporates CoT and RL training into TS-ASR for performance improvement. A novel CoT dataset of TS-ASR is constructed, and the TS-ASR model is first trained on regular data and then fine-tuned on CoT data. Finally, the model is further trained with RL using selected data to enhance generalized reasoning capabilities. Experiment results demonstrate a significant improvement of TS-ASR performance with CoT and RL training, establishing a state-of-the-art performance compared with previous works of TS-ASR on comparable datasets.

Index Terms— Target Speaker Automatic Speech Recognition, Chain-of-Thought, Reinforcement Learning

1. INTRODUCTION

Large Language Models (LLMs), which originally developed for text understanding and generation, have achieved significant advances by scaling the parameters and the data exponentially [1, 2]. In the field of speech and audio processing, Large Audio-Language Models (LALMs) have developed rapidly in recent years, expanding the capabilities of LLMs to speech-related tasks through pre-trained speech encoders [3, 4]. When dealing with Automatic Speech Recognition (ASR) tasks, for example, prevailing architectures of LLM-based ASR typically employ a pre-trained speech encoder to transform acoustic signals into feature representations, which are then injected into LLMs as prompts for end-to-end ASR training [5, 6]. Furthermore, LALMs architecture holds the potential to extend capabilities to more complex acoustic environments, such as TS-ASR in cocktail party scenarios, which remains a challenging research frontier.

Given a reference speech from the target speaker, TS-ASR is capable of selectively transcribing the speech of the target speaker while suppressing interference from other speakers [7]. Some previous works designed a neural network architecture to jointly train Target Speaker Extraction (TSE) and ASR in order to transcribe the speech of the target speaker [8, 9]. There were also some works integrating speaker embeddings into the training of conventional end-

to-end ASR models to directly train the TS-ASR model [10, 11, 12]. With the recent advancement of LALMs, MT-LLM [13], which utilized both an ASR encoder and a speaker encoder to extract features and then fed them into the LLM backbone to get TS-ASR output, has also demonstrated a good performance on the TS-ASR task. Although the development of LALMs has already brought new insights into TS-ASR, the optimization potential remains far from exhausted, as TS-ASR is a complex task that requires deep audio comprehension to achieve higher accuracy.

The breakthrough of DeepSeek-R1 [14] shows that combining CoT with RL significantly improves reasoning in complex tasks. This approach has also proven to be effective on many speech tasks, such as Audio Question and Answering (AQA) [15, 16] and ASR [17]. TS-ASR is a complex task that requires models to deeply understand speech and distinguish between different speakers, which involves logical reasoning to determine the number of speakers, identify the target speaker, and select the relevant speech segments for recognition. Therefore, we believe that CoT and RL are able to improve the performance of TS-ASR task by guiding the model to produce explicit intermediate reasoning steps before arriving at the final answer.

In this work, we introduce CoT and RL into the TS-ASR task to enhance reasoning capabilities in cocktail party scenarios. First, a TS-ASR *BASE* model is trained based on the LALMs architecture. Then, a novel method for constructing CoT training data is proposed, in which key information such as speaker count, overlap duration, speaker gender, and similarity to the reference is extracted from mixed speech. These attributes are subsequently structured through logical organization to form reasoning data for the CoT training, which is fine-tuned on the TS-ASR *BASE* model. Finally, based on the model trained on CoT data, part of the incorrectly predicted training data are randomly selected to do the further RL training. Experiment results show that our proposed method achieves an averaged 8.33% Word Error Rate (WER) of LibriSpeech, Libri2Mix and Libri3Mix. As far as we know, our method achieves a State-Of-The-Art (SOTA) performance compared with previous works of TS-ASR [9, 12, 13] on comparable datasets.

In summary, our contributions are as follows:

- We pioneer the integration of CoT and RL training into the TS-ASR task, establishing a novel framework that enhances reasoning capabilities in cocktail party scenarios.
- We construct a novel CoT training dataset for TS-ASR by extracting task-relevant information in logic from the mixed speech. This CoT dataset will be open-sourced to support further research at <https://github.com/Ease-Z/TS-ASR-CoT>
- Our approach demonstrates a significant improvement of TS-ASR performance with CoT and RL training, establishing a SOTA performance compared with previous works of TS-ASR on comparable datasets.

[†]Equal contribution, ^{*}Corresponding author: fanlichun1@xiaomi.com

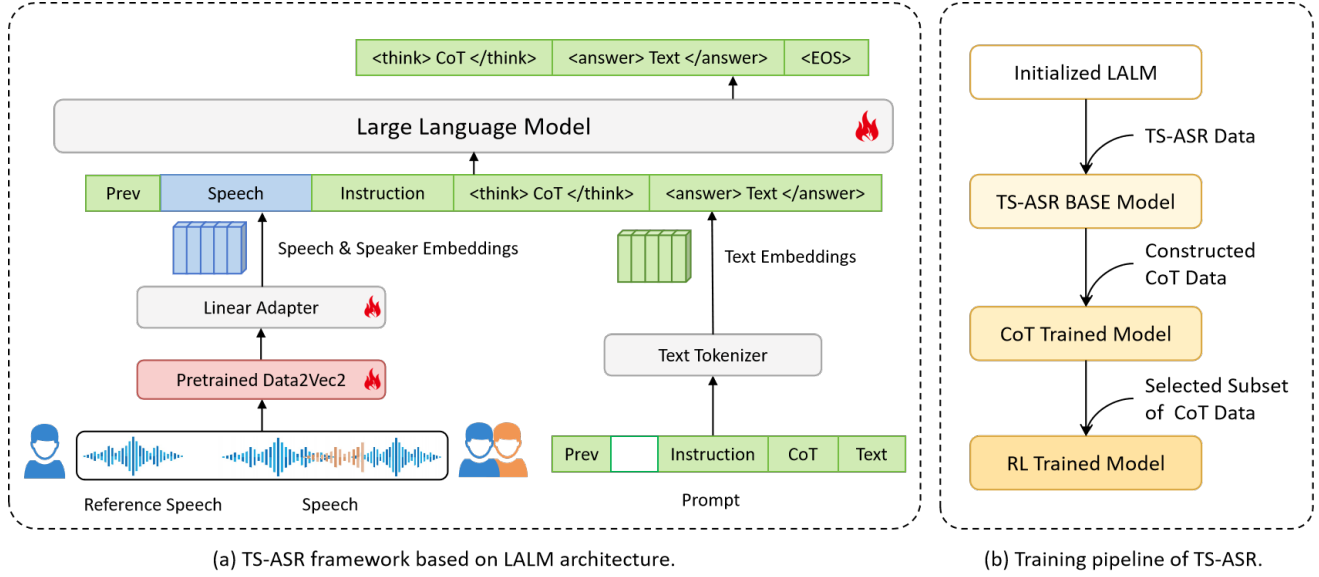


Fig. 1. Overview of the proposed method.

2. METHODOLOGY

2.1. LALM-based TS-ASR

2.1.1. Model structure

The proposed framework of the TS-ASR model is shown in Figure 1 (a). It consists of speech encoder, linear adapter, text tokenizer, and an LLM backbone network. We select a 3s segment from the target speaker utterance as reference, concatenate it with the mixed input (separated by a 3s silence), and feed the concatenated audio into the speech encoder. The pre-trained speech encoder Data2Vec2 [18] converts the input speech into frame-level embedding vectors that contain semantic and speaker information. To achieve modality alignment, a linear adapter maps the output from the speech encoder to the feature space of LLM. Meanwhile, the text tokenizer transforms the textual prompt into text embeddings. The speech and text embeddings are then combined and fed into the LLM. After fine-tuning, the model can focus on the target speaker within the mixed speech and generate accurate transcriptions.

In our framework, reference speech is concatenated with mixed speech as an audio prompt [19]. This approach effectively leverages the inherent capabilities of LLMs in generating targeted responses from prompts for TS-ASR. Moreover, the architecture employs only the Data2Vec2 speech encoder without introducing a separate speaker encoder. As speaker and semantic information are both inherently captured by the Self-Supervised Learning (SSL) model [20], their fusion is naturally achieved within the encoder.

2.1.2. Training pipeline

Our TS-ASR training follows a three-stage paradigm as shown in Figure 1 (b):

TS-ASR BASE A base TS-ASR model is trained via full parameter Supervised Fine-Tuning (SFT) on the initialized LALM, where all components are jointly updated, including the pre-trained LLM, the pre-trained speech encoder, and the randomly initialized adapter. This step aims to align multi-modal information and enable the model to perform TS-ASR. It effectively addresses the cold start problem and provides a foundation for subsequent CoT training.

CoT Training The TS-ASR BASE model is fine-tuned with constructed CoT data to guide step-by-step reasoning to improve TS-ASR performance. To focus model capacity on challenging cases, a random reasoning strategy is applied, providing full CoT procedure for hard samples while keeping empty <think> for easy samples.

RL Training Furthermore, based on the CoT trained model, RL refines the reasoning steps and improves the final performance using Generalized Reward Policy Optimization (GRPO) [21] with WER and CoT format rewards. In order to further focus learning on challenging cases, the RL training set is constructed by selecting a subset of constructed CoT data with wrong format or ASR transcription.

2.2. CoT training

2.2.1. CoT data construction

Our CoT dataset is constructed based on LibriSpeech [22] dataset, which contains gender information, speaker information, and ASR transcription for each single utterance. As shown in Figure 2, since the reasoning structure required for TS-ASR is relatively fixed, a standardized data template is designed that includes five components. First, audio information describes the composition of the input speech, the number of speakers, and the total duration. Second, reference speech primarily records the gender of the target speaker. Third, the corresponding speech of each speaker is described by timestamps, gender, and a similarity level to the reference. Since the mixed speech is generated from single-speaker utterances, these information can be obtained from the source single-speaker speech. Fourth, the target speaker is then inferred based on the above information. In single-speaker scenarios, identification is straightforward, while in multi-speaker scenarios, gender and similarity level are jointly considered to reach a decision. Finally, the answer is the transcription corresponding to the target speaker.

To compute similarity levels, we use CAM++ [23] to extract speaker embeddings and calculate the cosine similarity scores between embeddings of each source speech and the reference speech. Then, these continuous scores in $[0, 1]$ are mapped to five discrete levels $[1, 5]$ via uniform quantization. This discretization prevents the model from focusing on insignificant numerical differences, thereby improving training stability and convergence efficiency.

Single-speaker CoT Sample	2-Speaker CoT Sample
<p><think> Audio information: 0-3s is enrollment speech; 3-6s is silence; 6-9.04s is single-speaker audio; total duration 9.04s. Enrollment speech: male. Speaker1 information: male; from 6 to 9.04s; similarity to the enrollment speech is 4. Target speaker: Since this is a single-speaker audio, the Speaker1 must be the target speaker. Final output:\n </think> <answer>IT'S GONE TO RUI N AND DECAY THE LAST FEW YEARS</answer></p>	<p><think> Audio information: 0-3s is enrollment speech; 3-6s is silence; 6-20.95s is 2-speaker mixture audio; total duration 20.95s. Enrollment speech: male. Speaker1 information: male; from 6.0 to 9.04s; similarity to the enrollment speech is 4. Speaker2 information: female; from 6.0 to 20.95s; similarity to the enrollment speech is 2. Target speaker: Speaker1 and the enrollment speech are both male; 4(Speaker1) > 2(Speaker2); Speaker1 has the highest similarity score to the enrollment speech and is the target speaker. Final output: </think> <answer>THIRTY OR FORTY MILES EASTWARD FROM SEATTLE</answer></p>

Fig. 2. Sample of CoT data for single-speaker and 2-speaker mixed speech. Text in different colors represents different types of information.

2.2.2. Training with CoT data

The constructed CoT data are employed to further fine-tune the TS-ASR *BASE* model via SFT. Since LLMs often struggle with complex tasks, CoT guides them with explicit reasoning steps like humans. In this stage, the model is trained to analyze the speech content, estimate the number of speakers, describe the attributes of each speaker, and reason step by step to identify the target speaker. Finally, the model generates transcription from the identified target speaker.

Moreover, we introduce a random reasoning strategy for CoT training data. CoT data with the complete thinking is provided for hard samples, which is defined as those mispredicted by the *BASE* model. While the <think>content is left empty with 50% probability for easy samples that are correctly predicted. This strategy encourages the model to focus more on reasoning for challenging samples while allowing it to output answers directly for easier cases.

2.3. RL training

2.3.1. Data selection

Since RL fine-tuning can only use a subset of the training data, a data selection strategy is designed for TS-ASR, focusing on samples most beneficial for performance improvement. Specifically, the full training set is decoded using the model trained by CoT data, selecting all samples with incorrectly predicted CoT format and a random subset of samples with correct format but wrong ASR prediction. This approach is motivated by the fact that format errors are more fundamental and impactful than content errors, and including them in training helps the model further refine its reasoning process.

2.3.2. Generalized reward policy optimization

GRPO algorithm is adopted to train the model. For each input, a group of generated outputs are sampled and ranked within the group based on task-specific rewards. This approach effectively combines positive rewards that reinforce correct outputs and negative rewards that penalize uncertain or incorrect predictions.

The reward r for GRPO training is the combination of WER reward r_{WER} and CoT format reward r_{format} :

$$r = r_{\text{WER}} + r_{\text{format}}. \quad (1)$$

The WER reward evaluates the accuracy of the transcription between <answer> and </answer>. Sequences with a lower WER are assigned higher rewards, offering a direct guidance for speech recognition, defined as:

$$r_{\text{WER}} = 1 - \frac{\text{Sub} + \text{Del} + \text{Ins}}{N}, \quad (2)$$

where *Sub*, *Del*, *Ins* are substitution, deletion and insertion errors, respectively. And N is the number of words in ground-truth. The format reward is:

$$r_{\text{format}} = \begin{cases} 1 & \text{if } y \in F, \\ 0 & \text{otherwise,} \end{cases} \quad (3)$$

where F is the format: < think > ... < /think > < answer > ... < /answer > and y is the output string. The format reward encourages the model to perform the full reasoning process and ensures correct reasoning structure.

3. EXPERIMENT

3.1. Dataset

Based on the approach in LibriMix [24], this study generates mixed speech datasets of two speakers (Libri2Mix) and three speakers (Libri3Mix) using the 960-hour LibriSpeech [22] single-speaker corpus¹. Following the data scale in [13], we create mixed samples for each source utterance, alternately selecting each speaker in the mixture as the target. Reference speech segments are randomly selected from the single utterances of the same speaker. The duration of the Libri2Mix and Libri3Mix datasets is 1920 and 2880 hours, respectively. When combined with the original 960 hours of single-speaker data, the total training set reaches about 5760 hours. The test set is generated from LibriSpeech test-clean, with Libri2Mix and Libri3Mix constructed using the mixture and reference speech provided by the open-source SpeakerBeam toolkit².

3.2. Experiment settings

The proposed model uses Qwen2.5-0.5B-Instruct [25] as the backbone LLM. It adopts a pre-trained Data2Vec2 model as the audio encoder, which is trained on one million hours of unsupervised speech data. The number of parameters of Data2Vec2 is 0.3 billion. During the SFT stage, the AdamW optimizer [26] is applied with a learning rate schedule combining linear warmup and cosine annealing. The learning rate is gradually increased from 1e-7 to 7e-5 over the first 3000 steps and decayed to a minimum of 1e-6. The model is trained for 240000 steps, equivalent to one full epoch, with a global batch size of 112. To prevent input truncation, the maximum text length is set to 512. For RL training, the model is fine-tuned on 20000 selected samples for 2000 steps with a learning rate of 1e-6 and a

¹<https://github.com/JorisCos/LibriMix>

²<https://github.com/BUTSpeechFIT/speakerbeam/tree/main/egs>

temperature of 1. For each input, 8 responses are generated, with a maximum length of 512 tokens.

3.3. Evaluation metrics

In this experiment, the evaluation metric is WER, which is calculated by extracting the text between `<answer>` and `</answer>` tags from the decoded output of the LibriSpeech, Libri2Mix, and Libri3Mix test sets. For the small number of outputs with formatting errors or incomplete tags, the result is treated as empty for evaluation. To reflect overall performance, a weighted average is computed with the sample sizes as weights.

4. RESULT AND DISCUSSION

4.1. Comparison to baselines

We compare our proposed model with both traditional TS-ASR methods and LLM-based approaches. As shown in Table 1, our model achieves significant improvements, with WERs of 4.79% on Libri2Mix and 12.23% on Libri3Mix.

Compared with MT-LLM, the current best-performing LLM-based TS-ASR model, our TS-ASR *BASE* model without CoT data has a lower performance. However, training with CoT data achieves a significant improvement, with the average WER relatively reduced by 26.18%. It demonstrates that logic reasoning helps the model more accurately transcribe the target speaker. Based on the model trained on CoT data, RL fine-tuning is applied on selected data, further reducing the average WER to 8.33%. It indicates that RL training can help the model further improve the reasoning capabilities. Our model outperforms all baseline models in average WER, achieving new SOTA performance.

Table 1. TS-ASR performance on LibriSpeech, Libri2Mix and Libri3Mix test-clean dataset, as well as their weighted average. Evaluated by WER(%).

Model	LibriSpeech			
	Single	Libri2Mix	Libri3Mix	Avg
WavLM + TSE [9]	-	7.6	-	-
TS-VAD(460h) [12]	-	6.61	14.81	-
MT-LLM [13]	2.3	6.7	16.2	11.09
TS-ASR Base	8.1	7.4	17.24	12.45
TS-ASR + CoT	4.08	5.21	13.48	9.19
TS-ASR + CoT + RL	3.47	4.81	12.23	8.33

4.2. Ablation studies for CoT

We explore the impact of different CoT data and training strategies on model performance to identify an effective configuration. As shown in Table 2, the model trained on the data using continuous similarity scores is compared with that with discrete similarity levels. The results show that the use of discrete similarity levels outperforms using continuous similarity scores, which indicates that replacing continuous floating values with discrete levels in CoT data leads to more stable and effective model convergence. The benefit is particularly evident in multi-speaker scenarios, indicating that similarity information indeed helps the model identify the target speaker from overlapping speech, which is consistent with our expectations. In addition, the proposed random reasoning strategy mentioned in 2.2.2 is implemented. Experimental results show that this approach obtains a marginal improvement compared with using full CoT data, especially on Libri2Mix and Libri3Mix.

Table 2. Ablation Studies for CoT, evaluated by WER(%). FC indicate the use of Full CoT data and RR denotes the random reasoning.

Model	LibriSpeech			
	Single	Libri2Mix	Libri3Mix	Avg
TS-ASR + CoT				
FC + similarity score	3.8	5.39	14.09	9.52
FC + similarity level	3.88	5.29	13.57	9.23
RR + similarity level	4.08	5.21	13.48	9.19

4.3. Ablation studies for RL

Table 3 presents the impact of different training data selection strategies for RL on model performance. Specifically, "Random Sampling" denotes randomly selecting a subset from the full training dataset for RL training. "Balanced Correct&Error" selects correct and error samples in a 1:5 ratio from the model trained on CoT data for RL training, aiming to focus on difficult cases while preserving general performance. The "Error-only Sampling" strategy uses only incorrectly predicted samples for subsequent training.

Experimental results show that the "Random Sampling" strategy performs worst, with an average WER of 8.73%, while the "Error-only Sampling" strategy achieves the best result, reducing WER to 8.33%. This demonstrates that in reinforcement learning for TS-ASR, challenging error samples contribute more to performance improvement than correct samples. Therefore, training with a higher proportion of error samples is more effective for helping the model correct mistakes and improve overall accuracy.

Table 3. Ablation Studies for RL, evaluated by WER(%). RR indicate the random reasoning strategy.

Model	LibriSpeech			
	Single	Libri2Mix	Libri3Mix	Avg
TS-ASR + CoT (with RR + similarity level) + RL				
+ Random Sampling	4.08	5.01	12.70	8.73
+ Balanced Correct&Error	3.91	4.8	12.39	8.47
+ Error-only Sampling	3.47	4.81	12.23	8.33

5. CONCLUSION

In this work, we proposed a novel framework that incorporates CoT and RL training into the TS-ASR task to enhance reasoning capabilities in cocktail party scenarios. The approach began with training a TS-ASR *BASE* model using the LALM architecture. Subsequently, a novel methodology was proposed to generate CoT training samples by extracting salient features from overlapped speech, such as number of speakers, speech duration, gender information, and similarity to the reference speech. These features were systematically organized into structured reasoning data to fine-tune the TS-ASR *BASE* model. Furthermore, a subset of mispredicted samples from the CoT training phase was randomly selected for additional RL-based refinement. Experimental results showed a significant improvement of TS-ASR performance with CoT and RL training, which not only demonstrated the effectiveness of our constructed CoT dataset and the proposed training framework but also established a SOTA performance for TS-ASR on comparable benchmark datasets. The successful integration of CoT and RL into TS-ASR demonstrated the efficacy of logical reasoning in cocktail party scenarios. As LALMs continue to advance with CoT and RL techniques, we believe this framework has significant potential for further solving the cocktail party problem.

6. REFERENCES

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al., “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [2] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al., “Llama: Open and efficient foundation language models,” *arXiv preprint arXiv:2302.13971*, 2023.
- [3] Rongjie Huang, Mingze Li, Dongchao Yang, et al., “AudioGPT: Understanding and generating speech, music, sound, and talking head,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 21, pp. 23802–23804, Mar. 2024.
- [4] Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu, “Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities,” in *EMNLP (Findings)*, 2023.
- [5] Wenyi Yu, Changli Tang, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and Chao Zhang, “Connecting speech encoder and large language model for asr,” in *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2024, pp. 12637–12641.
- [6] Ye Bai, Jingping Chen, Jitong Chen, Wei Chen, Zhuo Chen, Chuang Ding, Linhao Dong, Qianqian Dong, Yujiao Du, Kepan Gao, et al., “Seed-asr: Understanding diverse speech and contexts with llm-based speech recognition,” *arXiv preprint arXiv:2407.04675*, 2024.
- [7] Naoyuki Kanda, Shota Horiguchi, Ryoichi Takashima, Yusuke Fujita, Kenji Nagamatsu, and Shinji Watanabe, “Auxiliary interference speaker loss for target-speaker speech recognition,” *Interspeech 2019*, pp. 236–240, 2019.
- [8] Yang Zhang, Krishna C. Puvvada, Vitaly Lavrukhin, and Boris Ginsburg, “Conformer-based target-speaker automatic speech recognition for single-channel audio,” in *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.
- [9] Zili Huang, Desh Raj, Paola García, and Sanjeev Khudanpur, “Adapting self-supervised models to multi-talker speech recognition using speaker embeddings,” in *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.
- [10] Pengcheng Guo, Xuankai Chang, Hang Lv, Shinji Watanabe, and Lei Xie, “Sq-whisper: Speaker-querying based whisper model for target-speaker asr,” *IEEE Transactions on Audio, Speech and Language Processing*, vol. 33, pp. 175–185, 2025.
- [11] Hao Ma, Zhiyuan Peng, Mingjie Shao, Jing Li, and Ju Liu, “Extending whisper with prompt tuning to target-speaker asr,” in *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2024, pp. 12516–12520.
- [12] Chikara Maeda, Muhammad Shakeel, and Yui Sudo, “Joint target-speaker asr and activity detection,” in *Proc. Interspeech 2025*, 2025, pp. 1683–1687.
- [13] Lingwei Meng, Shujie Hu, Jiawen Kang, Zhaoqing Li, Yuejiao Wang, Wenxuan Wu, Xixin Wu, Xunying Liu, and Helen Meng, “Large language model can transcribe speech in multi-talker scenarios with versatile instructions,” in *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2025, pp. 1–5.
- [14] Daya Guo, Dejian Yang, Haowei Zhang, et al., “Deepseek-r1 incentivizes reasoning in llms through reinforcement learning,” *Nature*, vol. 645, pp. 633–638, 2025.
- [15] Andrew Rouditchenko, Saurabhchand Bhati, Edson Araujo, Samuel Thomas, Hilde Kuehne, Rogerio Feris, and James Glass, “Omni-r1: Do you really need audio to fine-tune your audio llm?,” *arXiv preprint arXiv:2505.09439*, 2025.
- [16] Ziyang Ma, Zhuo Chen, Yuping Wang, Eng Siong Chng, and Xie Chen, “Audio-cot: Exploring chain-of-thought reasoning in large audio language model,” *arXiv preprint arXiv:2501.07246*, 2025.
- [17] Prashanth Gurunath Shivakumar, Yile Gu, Ankur Gandhe, and Ivan Buluko, “Group relative policy optimization for speech recognition,” *arXiv preprint arXiv:2509.01939*, 2025.
- [18] Alexei Baevski, Arun Babu, Wei-Ning Hsu, and Michael Auli, “Efficient self-supervised learning with contextualized target representations for vision, speech and language,” in *Proceedings of the 40th International Conference on Machine Learning*, 23–29 Jul 2023, vol. 202 of *Proceedings of Machine Learning Research*, pp. 1416–1429, PMLR.
- [19] Jian Cheng and Sam Nguyen, “Speech few-shot learning for language learners’ speech recognition,” in *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2025, pp. 1–5.
- [20] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al., “Wavlm: Large-scale self-supervised pre-training for full stack speech processing,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1505–1518, 2022.
- [21] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, et al., “Deepseekmath: Pushing the limits of mathematical reasoning in open language models,” *arXiv preprint arXiv:2402.03300*, 2024.
- [22] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, “Librispeech: An asr corpus based on public domain audio books,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 5206–5210.
- [23] Hui Wang, Siqi Zheng, Yafeng Chen, Luyao Cheng, and Qian Chen, “Cam++: A fast and efficient network for speaker verification using context-aware masking,” in *Interspeech 2023*, 2023, pp. 5301–5305.
- [24] Joris Cosentino, Manuel Pariente, Samuele Cornell, Antoine Deleforge, and Emmanuel Vincent, “Librimix: An open-source dataset for generalizable speech separation,” 2020.
- [25] Qwen Team, “Qwen2.5: A party of foundation models,” September 2024.
- [26] Ilya Loshchilov and Frank Hutter, “Decoupled weight decay regularization,” in *International Conference on Learning Representations*, 2019.