



# Large Language Models based ASR Error Correction for Child Conversations

Anfeng Xu<sup>\*1</sup>, Tiantian Feng<sup>\*1</sup>, So Hyun Kim<sup>2</sup>, Somer Bishop<sup>3</sup>, Catherine Lord<sup>4</sup>, Shrikanth Narayanan<sup>1</sup>

<sup>1</sup>Viterbi School of Engineering, University of Southern California, USA

<sup>2</sup>School of Psychology, Korea University, South Korea

<sup>3</sup>Weill Institute for Neurosciences, University of California, San Francisco, USA

<sup>4</sup>Semel Institute of Neuroscience and Human Behavior, University of California, Los Angeles, USA

anfengxu@usc.edu

## Abstract

Automatic Speech Recognition (ASR) has recently shown remarkable progress, but accurately transcribing children’s speech remains a significant challenge. Recent developments in Large Language Models (LLMs) have shown promise in improving ASR transcriptions. However, their applications in child speech including conversational scenarios are under-explored. In this study, we explore the use of LLMs in correcting ASR errors for conversational child speech. We demonstrate the promises and challenges of LLMs through experiments on two children’s conversational speech datasets with both zero-shot and fine-tuned ASR outputs. We find that while LLMs are helpful in correcting zero-shot ASR outputs and fine-tuned CTC-based ASR outputs, it remains challenging for LLMs to improve ASR performance when incorporating contextual information or when using fine-tuned autoregressive ASR (e.g., Whisper) outputs.

**Index Terms:** Automatic speech recognition, large language model, child speech.

## 1. Introduction

Automatic Speech Recognition (ASR) has made substantial advances in recent years, driven by Speech Foundation Models (SFM) [1], trained with advanced deep learning architectures, such as transformers [2] and conformers [3], while leveraging extensive training data. SFMs can be categorized into two categories. The first is end-to-end supervised models that leverage massive labeled datasets to jointly align acoustic and language information. Whisper [4] and Parakeet [5] are widely used models in this category. The second is models trained with self-supervised learning (SSL) regime, such as Wav2vec 2.0 [6], HuBERT [7], and WavLM [8], which learn speech representations from unlabeled audio data. When fine-tuned, the SSL-based models demonstrate competitive performance on ASR tasks.

However, technological improvements have primarily focused on improving adult speech recognition, while ASR for children’s speech remains a persistent challenge: recent evaluations have shown ASR error rates for child speech are 10 to 19 times worse than for adults with general models, and 6 times worse despite adaptation on children speech [9]. Accurately transcribing conversations involving children is crucial for various applications, including educational technology, clinical assessments, and developmental research [10]. Yet, current SFMs underperform on this task, as children’s speech patterns differ significantly from those of adults in terms of acoustic-phonetic characteristics [11, 12], vocabulary usage, prosodic features, and conversational dynamics [13]. These challenges are further

\*These authors contributed equally to this work

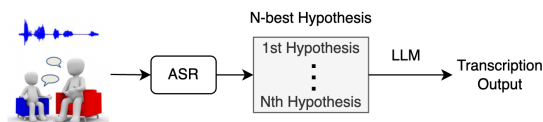


Figure 1: Overall pipeline for ASR with LLM error correction.

compounded by the relative scarcity of large-scale and naturalistic children’s speech datasets, resulting in ASR systems that fall short of generalizing to child-inclusive applications [14, 9].

Large Language Models (LLMs) have gained substantial attention in natural language processing through their advanced capabilities in processing large volumes of input data and generating sophisticated inferences and responses [15]. While initially developed for text-related tasks, these models have shown promising applications in ASR systems. For example, several studies [16, 17, 18] have shown that LLMs can perform ASR tasks when integrated with audio or speech encoders. Ogawa et al. [19] used LLMs to improve ASR transcriptions by rescoring ASR hypotheses using conversational contexts.

LLMs’ capability to process language structure, context, and semantic relationships enables them to effectively correct ASR errors by considering both narrow linguistic patterns and broader semantic context [15]. Recent works [20, 21] have demonstrated the effectiveness of leveraging LLMs for ASR error correction by selecting and refining the most probable transcription using N-best hypotheses from ASR systems. However, to our knowledge, limited work examines LLMs for ASR error correction in conversational settings, particularly in child-inclusive contexts. Additionally, these works have mainly focused on improving zero-shot ASR outputs, while LLMs’ capabilities to improve fine-tuned ASR outputs remain unclear.

In this work, we investigate approaches for improving ASR accuracy with child-adult conversations by error correction using LLMs. Our method builds upon the HyParadise [20] benchmark for LLM-based ASR error correction by adapting it to handle conversational child speech. The main contributions of this work are summarized below.

- We incorporate LLMs for ASR error correction for children’s conversational speech. To the best of our knowledge, this work is one of the earliest attempts in this domain.
- We investigate the effectiveness of LLMs in improving ASR transcription across multiple scenarios: applying LLMs to both zero-shot and fine-tuned outputs, from supervised and self-supervised ASR models.
- We investigate utilizing conversational context to improve LLM-based error correction by leveraging previous utterances in the conversation history as additional input.

Error Correction Prompt without Context
<p><b>[System Prompt]</b> You're a helpful assistant that help to correct transcriptions between a child and a clinician.</p> <p><b>[User Prompt]</b> Below is the best-hypotheses transcribed from speech recognition system between interactions between a child and a clinician, and the speaker of this sentence is the {speaker}. Please revise it using the words which are only included into other-hypothesis, and only write the response for the true transcription. ### Best-hypothesis: {best} ### Other-hypothesis: {others}</p>

Figure 2: LLM prompt without context.

Error Correction Prompt with Context
<p><b>[System Prompt]</b> You're a helpful assistant that help to correct transcriptions between a child and a clinician.</p> <p><b>[User Prompt]</b> Here is the previous {num_context} utterances. {prev_sentences}. Below is the best and other hypotheses transcribed from a speech recognition system for the current utterance by {speaker}. Please revise it using the words which are only included into other-hypothesis, and only write the response for the true transcription. ### Best-hypothesis: {best} ### Other-hypotheses: {others}</p>

Figure 3: LLM prompt with context.

## 2. Methods

### 2.1. LLM Prompt Design

Our approach to LLM-based error correction for child-adult conversations builds upon the benchmark framework established by HyPoradise [20], which uses N-best hypotheses from ASR for LLMs as illustrated in Figure 1. While HyPoradise focused on general ASR error correction, we specifically adapted their methodology for conversational speech between children and adults. We use 5-best hypotheses from ASR outcomes and train an LLM for error correction. We use LLaMA3 [22] models for our experiments.

#### 2.1.1. ASR Error Correction without Context

We first examine if LLMs can help correct ASR errors without previous conversational context. The prompt we use is shown in Figure 2. {*speaker*} is replaced with either “Child” or “Adult”, while {*best*} and {*others*} are replaced with the top-1 ASR hypothesis and remaining top ASR hypotheses, respectively.

#### 2.1.2. ASR Error Correction with Context

Our previous work has shown the utility of dialog context modeling to improve child ASR [23]. We hence investigate whether incorporating prior conversational context can also guide LLMs to correct ASR errors in child-adult interactions. We also note that young children frequently echo adults’ speech in conversational interaction, providing natural repetition that could serve as meaningful contextual information for error detection. We experiment by including either 1 or 3 previous utterances. Figure 3 shows the prompt we used. {*speaker*}, {*best*}, and {*others*} are replaced similarly as for the case without context. {*num\_context*} is replaced with the number of previous utterances, and {*prev\_sentences*} is replaced by the previous 1 or 3 utterances with speaker tags (e.g., “Adult: how are you?”). The ground-truth previous utterances are used during the training, while inferred previous utterances are used for testing. When there is no previous utterance, it is replaced by “None.”

### 2.2. ASR Models

To evaluate the effectiveness of LLM-based ASR error correction, we conduct experiments using two distinct ASR architectures to assess our approach across different ASR paradigms.

**Whisper** [4] uses an attention-based encoder-decoder architecture, trained jointly for ASR and related tasks (e.g., translation) with weak supervision, using 680k hours of multi-language speech content gathered from online sources. It has shown competitive results on ASR benchmarks and works ro-

bustly across different recording scenarios. Specifically, we apply the Whisper-small (**WSP-S**), Whisper-large-v3 (**WSP-L**), and Whisper-large-v3-turbo (**WSP-L-T**) models. We use all these models to generate zero-shot ASR outputs and report the fine-tuned ASR results with the WSP-L-T. During inference, we apply beam search decoding with a beam size of 60.

**WavLM** [8] is an SSL-based model that builds upon Wav2vec 2.0 and HuBERT. It learns universal speech representations through masked speech prediction in pre-training. It shows strong transfer learning performance across diverse speech processing tasks, thus being the current state-of-the-art for the SUPERB benchmark [24]. We only use fine-tuned WavLM outputs for ASR experiments rather than zero-shot inference. In our experiments, we fine-tune WavLM-large (**WavLM-L**) with CTC loss for character predictions. For WavLM inference, we employ beam search decoding with a smaller beam size of 10.

### 2.3. ASR Fine-tuning

In addition to LLM-based error correction for zero-shot ASR outputs, we investigate whether they are effective at correcting fine-tuned ASR outputs. To this end, we generate fine-tuned ASR outputs for both training and test sets. We conduct 2-fold cross-validation on the training set and use the validation ASR outputs as the training set for LLM error correction. For the test sets, we use the fine-tuned ASR models trained on the full training set. We use WSP-L-T and WavLM-L as the ASR models.

## 3. Experiments

### 3.1. Evaluation

We report mean Word Error Rate (WER) across all utterances. Before calculating WER for each utterance, we pass the ground truth transcript and ASR outputs to the Whisper normalizer.

### 3.2. Dataset

We consider two child conversational datasets: My Science Tutor (MyST) Children’s speech corpus [25] and ADOS-Mod3 corpus of Autism diagnostic administration [26]. Our research complies with all Institutional Review Board (IRB) protocols and follows the Data Use Agreements (DUAs) established by the original data providers.

**MyST** [27] includes transcribed conversations between children and virtual tutors. The children were recruited from grade 3 to grade 5, which corresponds to around 8 to 12 years of age. The topics include 8 areas of science, such as biology, physics,

and others. Similar to [10], we filter out samples longer than 30s (maximum length for Whisper). However, unlike benchmark results reported from [10], we did not filter out the samples based on the number of words or WER from zero-shot Whisper-large, allowing us to get a more comprehensive assessment of the ASR performance. In addition, we use this dataset only for the experiments for ASR Error Correction *without* Context, as the corpus does not include speech or transcript data from the virtual tutors. We use the official training and test split for the LLM ASR error correction experiments.

**ADOS-Mod3** dataset [26] contains 352 sessions collected from 180 children during two specific sections of the ADOS-2 autism diagnostic protocol: “Social Difficulties” and “Annoyance and Emotional” tasks. The children ranged in age from 2 to 13 years, with 45 being female. Approximately half received an autism spectrum disorder diagnosis, while the remaining children were diagnosed with various other conditions, including ADHD and mental or language disorders. 96 children and 84 children were recorded at two different medical centers in Chicago (CHIC) and Michigan (MICH). On average, each session contains 25.9 child utterances and 30.0 adult utterances, with mean durations of 2.58s and 2.06s, respectively. For LLM-based ASR error correction experiments, we use data collected from CHIC as the training set and MICH as the test set. Since this dataset contains both child and adult speech utterances, we also report the individual WERs.

### 3.3. Experimental Setup: ASR Fine-tuning

For Whisper finetuning, we choose WSP-L-T because of its performance and relatively smaller size compared to WSP-L. We train for 2000 steps with a learning rate of  $1e - 6$ . For the WavLM, we train for 30000 steps with a learning rate of  $3e - 4$ . Adam optimizer is used with a batch size of 32 for both models. The same configurations are used for both datasets. We choose the best model based on the validation WER.

For the MyST dataset, we use the official training and validation sets for fine-tuning, and we report the WERs on the test set. However, to prepare training data for the LLM instruction tuning for fine-tuned ASR models, we split the training data in half for fine-tuning ASR models. We then apply the fine-tuned ASR models to the other half of the training dataset. We conduct the fine-tuning for each of the 2-splits.

For the ADOS-MOD3 dataset, we randomly select 80% from the CINC as the training set and the rest of 20% as the validation set, with the data from MICH as the test set. Similar to the MyST dataset setup, we conduct a 2-split fine-tuning to generate the fine-tuned ASR outputs for the train set.

### 3.4. Experimental Setup: LLM Instruction tuning

We experiment with the LLaMa 3.1-8B and LLaMa 3.2-1B models for ASR correction. We train the LLMs for 5 and 10 epochs on the MyST and the ADOS datasets, respectively. We apply a learning rate of  $5e - 4$  in all LLM fine-tuning experiments. Our system prompt for fine-tuning the ADOS dataset is shown in Figure 2, while the system prompt for fine-tuning the MyST dataset is “You are a helpful assistant that helps to correct transcriptions from a child in a tutoring session.” We empirically tested that ASR correction results remain robust across varying temperature values, with different temperatures yielding similar ASR correction outputs. Therefore, we set the temperature to 0.2 during the inference phase in all experiments. Even though LLMs generally produce reasonable outputs, we identify instances where they could generate repeated or hallu-

Table 1: WER comparison with LLM for zero-shot ASR error correction using ADOS-Mod3 and MyST dataset.

ASR	LLaMA3	ADOS			MyST
		Overall	Child	Adult	Child
WSP-S	Unused	46.67	63.73	32.23	22.33
	1B	47.19	64.64	32.41	22.20
	8B	<b>43.96</b>	<b>62.71</b>	<b>28.10</b>	<b>20.60</b>
WSP-L-T	Unused	40.77	55.84	28.07	20.01
	1B	39.11	54.29	26.30	19.66
	8B	<b>37.09</b>	<b>53.87</b>	<b>22.94</b>	<b>18.35</b>
WSP-L	Unused	40.26	55.19	27.65	19.58
	1B	39.55	54.48	26.93	19.50
	8B	<b>36.70</b>	<b>52.63</b>	<b>23.24</b>	<b>18.41</b>

Table 2: WER comparison with LLM for fine-tuned ASR output error correction using ADOS-Mod3 and MyST dataset.

ASR	LLaMA3	ADOS			MyST
		Overall	Child	Adult	Child
WSP-L-T	Unused	<b>32.11</b>	<b>46.99</b>	<b>19.47</b>	<b>14.31</b>
	1B	33.25	47.93	20.77	14.55
	8B	32.92	47.47	20.56	14.40
WavLM-L	Unused	66.33	88.05	47.87	27.54
	1B	56.83	78.34	38.54	19.93
	8B	<b>50.58</b>	<b>72.24</b>	<b>16.45</b>	<b>16.45</b>

inated lengthy content. Thus, we set the ASR output as the best hypothesis whenever the generated output exceeds the best ASR hypothesis by more than three words.

## 4. Results and Discussion

### 4.1. Can LLMs Improve *zero-shot* Child ASR Results?

Table 1 shows the LLM error correction results for the zero-shot Whisper ASR outputs. We see consistent reductions in WERs across all three ASR models for each dataset when LLaMA 3.1-8B model is used. The improvements are less substantial when the LLaMA 3.2-1B model is applied. Interestingly, the 3.2-1B model slightly increases the WER for WSP-S with the ADOS-Mod3 dataset. Thus, the parameter size of the LLM used is critical for the ASR error correction in this domain, and larger LLMs tend to perform better in ASR corrections. In summary, these results show that when the audio resources or the ASR models are unavailable for training, LLMs with larger parameter sizes can help refine the transcriptions.

### 4.2. Can LLMs Improve *Fine-tuned* Child ASR Results?

Table 2 shows WERs from the fine-tuned ASR outputs before and after applying LLM error corrections. Since the previous benchmark applied additional dataset filtering (e.g. removing utterances with WER > 50%), the WER for MyST in this study is higher than those reported in [10]. The results show that the WSP-L-T model substantially outperforms WavLM-L for both datasets, similar to [10]. Based on our manual error inspections, one plausible reason is that the CTC-based character prediction approach produces more spelling errors, especially for children with underdeveloped pronunciation capabilities. In contrast, the Whisper model is less likely to make spelling mistakes since it generates the transcriptions by predicting byte-level Byte-Pair Encoding (BPE) tokens autoregressively during inference.

For both datasets, LLMs show substantial improvements for WavLM ASR outputs, while they do not show improve-

Table 3: WER comparison with LLaMA 3.1-8B for ASR error correction using context. ADOS-Mod3 dataset is used. (ft) indicates whether the ASR model is fine-tuned or not.

ASR (ft)	# Context	Overall	Child	Adult
WSP-L-T (No)	1	38.06	55.67	23.21
	3	37.79	54.65	23.56
WSP-L-T (Yes)	1	33.02	47.03	21.1
	3	37.87	55.58	22.81
WavLM-L (Yes)	1	52.99	74.46	34.73
	3	54.98	78.47	35.02

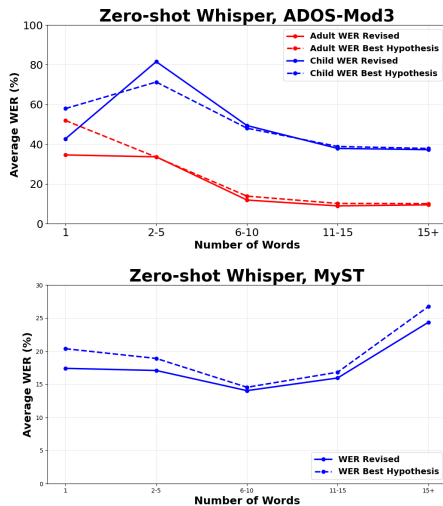


Figure 4: WERs by utterance lengths with zero-shot Whisper ASR (WSP-L-T). Results from both datasets.

ments for Whisper ASR outputs. We observe that LLMs help correct spelling errors that the fine-tuned WavLM produces. However, we reason that LLMs show limited advantages for Whisper outputs because both systems use similar autoregressive decoding, where each token depends on previous predictions. In addition, unlike the Whisper decoder, LLM decoders can not access the speech features through cross-attention.

#### 4.3. Does Context Improve LLM Error Correction?

Table 3 presents the WERs obtained using LLM-based ASR error correction with contextual information from the previous 1 or 3 utterances. Contrary to our initial hypothesis, incorporating previously predicted utterances leads to increased WERs compared to LLM-based ASR error correction without context, across all our experimental conditions. Furthermore, using the context of 3 utterances yields higher error rates than using the context of a single utterance. **One plausible reason behind this performance degradation is error propagation**, as the previously predicted utterances already include recognition errors.

#### 4.4. Analysis on utterance length

To have a more comprehensive understanding of when LLMs can help correct ASR errors, we examine how the WER changes when varying the number of words in the utterances. Figure 4 shows the results for zero-shot ASR with WSP-L-T using both ADOS-Mod3 and MyST datasets. The zero-shot ASR results show that LLM correction is most effective for single-word utterances. This is likely because Whisper models of-

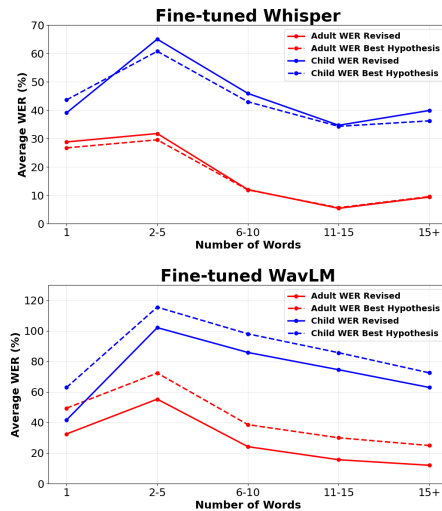


Figure 5: WERs by utterance lengths with fine-tuned ASR models (WSP-L-T, WavLM-L), using the ADOS-Mod3 dataset.

ten generate phonetically similar but contextually inappropriate words (including words from other languages) when processing short or unclear utterances. The language understanding of the LLMs helps filter out these mismatched transcriptions and convert them to more probable conversational utterances. We have also found that the improvements are more observable for adult speech in ADOS-Mod3 and child speech in MyST than for child speech in ADOS-Mod3. This is likely due to the challenges in correcting ASR outputs of children with less developed language skills prevalent in ADOS-Mod3 data.

Moreover, Figure 5 shows the results for fine-tuned ASR with WSP-L-T and WavLM-L, using only the ADOS-Mod3 dataset. The results demonstrate that LLMs do not show improvements for the Whisper model other than for single utterances from children. For WavLM, the improvements are consistent across utterances with different utterance lengths, where we found most of the improvements are from correcting spelling errors as discussed in Section 4.1.

## 5. Conclusion

This paper has investigated the use of LLMs for ASR error correction in child conversations, making several key findings. First, larger LLMs consistently improve zero-shot ASR performance across different Whisper models, while smaller LLMs show limited benefits. Second, for fine-tuned ASR systems, LLMs substantially improve CTC-based self-supervised ASR outputs by correcting spelling errors but show minimal improvements for the outputs from Whisper, a supervised ASR model with the attention-based encoder-decoder architecture. Third, contrary to initial hypotheses, incorporating conversational context degrades error correction performance, likely due to error propagation from previous utterances. These findings advance our understanding of LLM capabilities in child speech recognition while highlighting some of its limitations. Future work may include investigating larger LLMs and designing more effective strategies to incorporate conversational context.

## 6. Acknowledgment

This work was supported by SIMONS FOUNDATION (SFI-AR-HUMAN-00004115-03, 655054).

## 7. References

- [1] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill *et al.*, “On the opportunities and risks of foundation models,” *arXiv preprint arXiv:2108.07258*, 2021.
- [2] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, and A. N. Gomez, “L. u. kaiser, and i. polosukhin,” “attention is all you need,” *Advances in neural information processing systems*, vol. 30, pp. 5998–6008, 2017.
- [3] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, and R. Pang, “Conformer: Convolution-augmented transformer for speech recognition,” in *Interspeech 2020*, 2020, pp. 5036–5040.
- [4] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, “Robust speech recognition via large-scale weak supervision,” in *Proceedings of the 40th International Conference on Machine Learning*, 2023, pp. 28 492–28 518.
- [5] D. Rekish, N. R. Koluguri, S. Kriman, S. Majumdar, V. Noroozi, H. Huang, O. Hrinchuk, K. Puvvada, A. Kumar, J. Balam *et al.*, “Fast conformer with linearly scalable attention for efficient speech recognition,” in *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2023, pp. 1–8.
- [6] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *Advances in neural information processing systems*, vol. 33, pp. 12 449–12 460, 2020.
- [7] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhota, R. Salakhutdinov, and A. Mohamed, “Hubert: Self-supervised speech representation learning by masked prediction of hidden units,” *IEEE/ACM transactions on audio, speech, and language processing*, vol. 29, pp. 3451–3460, 2021.
- [8] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z. Chen, J. Li, N. Kanda, T. Yoshioka, X. 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.
- [9] P. Gurunath Shivakumar and S. Narayanan, “End-to-end neural systems for automatic children speech recognition: An empirical study,” *Computer Speech & Language*, vol. 72, p. 101289, 2022. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0885230821000905>
- [10] R. Fan, N. Balaji Shankar, and A. Alwan, “Benchmarking children’s asr with supervised and self-supervised speech foundation models,” in *Interspeech 2024*, 2024, pp. 5173–5177.
- [11] S. Lee, A. Potamianos, and S. S. Narayanan, “Acoustics of children’s speech: Developmental changes of temporal and spectral parameters,” *Journal of the Acoustical Society of America*, vol. 105, no. 3, pp. 1455–1468, mar 1999, selected Research Article.
- [12] S. Lee, A. Potamianos, and S. Narayanan, “Developmental acoustic study of american english diphthongs,” *The Journal of the Acoustical Society of America*, vol. 136, no. 4, pp. 1880–1894, 2014.
- [13] A. Potamianos and S. Narayanan, “Robust recognition of children’s speech,” *IEEE Transactions on Speech and Audio Processing*, vol. 11, no. 6, pp. 603–616, 2003.
- [14] R. Southwell, W. Ward, V. A. Trinh, C. Clevenger, C. Clevenger, E. Watts, J. Reitman, S. D’Mello, and J. Whitehill, “Automatic speech recognition tuned for child speech in the classroom,” in *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2024, pp. 12 291–12 295.
- [15] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, “A survey of large language models,” *arXiv preprint arXiv:2303.18223*, 2023.
- [16] Y. Fathullah, C. Wu, E. Lakomkin, J. Jia, Y. Shangguan, K. Li, J. Guo, W. Xiong, J. Mahadeokar, O. Kalinli *et al.*, “Prompting large language models with speech recognition abilities,” in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 13 351–13 355.
- [17] Z. Ma, G. Yang, Y. Yang, Z. Gao, J. Wang, Z. Du, F. Yu, Q. Chen, S. Zheng, S. Zhang *et al.*, “An embarrassingly simple approach for llm with strong asr capacity,” *arXiv preprint arXiv:2402.08846*, 2024.
- [18] W. Yu, C. Tang, G. Sun, X. Chen, T. Tan, W. Li, L. Lu, Z. Ma, and C. 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. 12 637–12 641.
- [19] A. Ogawa, N. Kamo, K. Matsuura, T. Ashihara, T. Moriya, T. Kano, N. Tawara, and M. Delcroix, “Applying llms for rescoring n-best asr hypotheses of casual conversations: Effects of domain adaptation and context carry-over,” *arXiv preprint arXiv:2406.18972*, 2024.
- [20] C. Chen, Y. Hu, C.-H. H. Yang, S. M. Siniscalchi, P.-Y. Chen, and E.-S. Chng, “Hyporadise: An open baseline for generative speech recognition with large language models,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [21] R. Ma, M. Qian, P. Manakul, M. Gales, and K. Knill, “Can generative large language models perform asr error correction?” *arXiv preprint arXiv:2307.04172*, 2023.
- [22] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan *et al.*, “The llama 3 herd of models,” *arXiv preprint arXiv:2407.21783*, 2024.
- [23] M. Kumar, S. H. Kim, C. Lord, T. Lyon, and S. Narayanan, “Leveraging linguistic context in dyadic interactions to improve automatic speech recognition for children,” *Computer, Speech and Language*, vol. 63, 2020.
- [24] S.-w. Yang, P.-H. Chi, Y.-S. Chuang, C.-I. J. Lai, K. Lakhota, Y. Y. Lin, A. T. Liu, J. Shi, X. Chang, G.-T. Lin *et al.*, “Superb: Speech processing universal performance benchmark,” *arXiv preprint arXiv:2105.01051*, 2021.
- [25] S. Pradhan, R. Cole, and W. Ward, “My science tutor (myst)—a large corpus of children’s conversational speech,” in *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, 2024, pp. 12 040–12 045.
- [26] R. Lahiri, M. Nasir, M. Kumar, S. H. Kim, S. Bishop, C. Lord, and S. Narayanan, “Interpersonal synchrony across vocal and lexical modalities in interactions involving children with autism spectrum disorder,” *JASA express letters*, vol. 2, no. 9, 2022.
- [27] W. Ward, R. Cole, D. Bolanos, C. Buchenroth-Martin, E. Svirsky, S. V. Vuuren, T. Weston, J. Zheng, and L. Becker, “My science tutor: A conversational multimedia virtual tutor for elementary school science,” *ACM Transactions on Speech and Language Processing (TSLP)*, vol. 7, no. 4, pp. 1–29, 2011.