

Examining Test-Time Adaptation for Personalized Child Speech Recognition

Zhonghao Shi^{1,2}, Xuan Shi¹, Anfeng Xu¹, Tiantian Feng¹, Harshvardhan Srivastava^{2,3}, Shrikanth Narayanan¹, Maja Mataric¹

¹Viterbi School of Engineering, University of Southern California, United States

²Sara Technology Inc., United States

³Fu Foundation School of Engineering and Applied Science, Columbia University, United States

zhonghas@usc.edu, xuanshi@usc.edu, anfengxu@usc.edu, tiantiaf@usc.edu,
hs3447@columbia.edu, shri@usc.edu, mataric@usc.edu

Abstract

Automatic speech recognition (ASR) models often experience performance degradation due to data domain shifts introduced at test time, a challenge that is further amplified for child speakers. Test-time adaptation (TTA) methods have shown great potential in bridging this domain gap. However, the use of TTA to adapt ASR models to the individual differences in each child’s speech has not yet been systematically studied. In this work, we investigate the effectiveness of two widely used TTA methods—SUTA, SGEM—in adapting off-the-shelf ASR models and their fine-tuned versions for child speech recognition, with the goal of enabling continuous, unsupervised adaptation at test time. Our findings show that TTA significantly improves the performance of both off-the-shelf and fine-tuned ASR models, both on average and across individual child speakers, compared to unadapted baselines. However, while TTA helps adapt to individual variability, it may still be limited with non-linguistic child speech.

Index Terms: children speech recognition, human-computer interaction, test time adaptation

1. Introduction

Child-AI interaction enabled by AI software agents [1] and socially assistive robots [2] has shown great potential for many application domains, for example education [3,4]. Conversational capabilities for these AI agents can support natural interaction with the child in achieving task goals [5]. To enable seamless human-like interaction, these AI agents and robots require accurate recognition of child speech [6]. Despite tremendous progress in machine learning methods for automatic speech recognition (ASR), a large body of recent work has shown that off-the-shelf ASR models do not generalize well to children’s speech data, due to the high amount of acoustic and linguistic variability [7], resulting in data domain shifts between the adult data used for pre-training and for testing [8,9].

Recent work on child ASR [10–18] has experimented with various supervised and unsupervised methods to adapt off-the-shelf ASR models at training time. These studies have proposed applying methods such as transfer learning [10–13,19], continued pre-training [14], adapters [15,16], low-rank adaptation [17] and unsupervised domain adaptation [20,21] to fine-tune and adapt the off-the-shelf ASR models with annotated children’s speech data at training time. The results from these prior work have shown that both supervised and unsupervised adaptation at training time can substantially improve models for recognizing children’s speech.

Despite recent progress, training-time adaptation is not always feasible for real-world model deployments due to several limitations: (1) individual variability: each new child

Table 1: *Definitions of different model adaptation settings for child speech recognition, adapted from [22]. Unlike other adaptation settings, TTA adapts off-the-shelf ASR models with child^c data (x^c) only at test time, without the need of annotations (y^c) or adult^a pre-training dataset (x^a , y^a) at training time.*

Setting	Train Loss	Test Loss
Supervised Pre-training	$L(x^a, y^a)$	-
Supervised Fine-Tuning	$L(x^c, y^c)$	-
Unsupervised Domain Adaptation	$L(x^a, y^a, x^c)$	-
Unsupervised Test-Time Adaptation (Ours)	-	$L(x^c)$

speaker introduces domain shifts at test time, and one-size-fits-all training-time adaptation fails to capture these individual differences; (2) annotation constraints: supervised adaptation requires labeled data, but obtaining annotations for each user poses significant financial, human labor, and logistical challenges; and (3) privacy and computational constraints: Users may prefer to keep their data private on local devices with limited computing and storage capacity, making both supervised and unsupervised training-time adaptation impractical.

Recent advancements in unsupervised test-time adaptation (TTA) methods offer a computationally efficient alternative [22–24], enabling on-the-fly model adaptation to test data domains without requiring human annotations. However, no prior work has systematically evaluated the effectiveness of TTA methods in addressing domain adaptation challenges in child speech recognition. In this work, we investigate this gap by conducting experiments to answer three research questions, and report the following results and findings:

- *RQ1: Why is TTA needed for ASR models in child speech recognition?* We found that both off-the-shelf and fine-tuned ASR models do not generalize robustly across different child speakers. Substantial domain shifts occur both across and within individual speakers, highlighting the need for adaptation at test time.
- *RQ2: Can unsupervised TTA methods effectively adapt both off-the-shelf and fine-tuned ASR models for child speech recognition?* As shown in Table 2, we evaluated two popular TTA methods—SUTA and SGEM—across various off-the-shelf and fine-tuned model settings, using the MyST dataset. Our results indicate that SUTA is more robust than SGEM across different model settings. Notably, SUTA significantly improves both off-the-shelf models trained on out-

of-domain adult speech and models fine-tuned on in-domain child speech. These findings further validate the need for test-time adaptation to account for individual variability.

- *RQ3: When do TTA methods perform well, and when do they fail?* We conducted extensive analyses to probe the relationship between TTA performance and acoustic as well as linguistic characteristics of the input samples. We found that TTA may help adaptation to the individual background noise level, while non-linguistic speech may be particularly challenging and misleading for TTA’s adaptation process.

2. Methods

In this work, we examine two widely used test-time adaptation (TTA) methods: 1) Single-Utterance Test-Time Adaptation (SUTA) ; and 2) Sequential-Level Generalized Entropy Minimization (SGEM). Our goal is to evaluate their effectiveness in adapting various off-the-shelf and fine-tuned ASR models to each child speaker in a continuous and unsupervised fashion.

2.1. Problem Formulation

A canonical ASR model can be denoted as $z = f(x, \theta)$, where x is the input speech waveform, and θ refers to the ASR model parameters. Off-the-shelf ASR models are typically pre-trained on a primarily *adult*^a dataset $D_{adult} = \{(x_i^a, y_i^a)\}_I$ in a supervised or self-supervised manner, to estimate θ^a . The models are then used to transcribe \hat{y}_j^t from $x_j^t \in D_{test}$, which assumes identical distribution with D_{adult} . However, when recognizing child speech D_{child} , there exists a significant data domain shift between $D_{adult} = \{(x_i^a, y_i^a)\}_I$ and $D_{child} = \{(x_i^c, y_i^c)\}_I$ due to the wide acoustic and linguistic variability in child speech and lack of child-specific ASR dataset [8, 9].

To address this challenge, we use test-time adaptation (TTA) methods to adjust the parameters of the off-the-shelf ASR model $\theta^a \rightarrow \theta^c$, so that models can adapt to the domain shifts at test time for *child*^c speech. The goal of TTA is to design optimization objectives (\mathcal{L}) based on the output logits $z \in \mathbb{R}^{L \times C}$ to adapt ASR models to the current test child’s speech by continuously updating a small portion of model’s parameters, where $z \in \mathbb{R}^{L \times C}$ is the predicted context logits. L is the total number of timestamp, and C is the number of word class.

2.2. Test-Time Adaptation Methods

This work experimented with two widely used TTA methods: 1) SUTA [23]; and 2) SGEM [24] to adapt the feature extractor layers of Wav2vec2.0. The unsupervised optimization objective of SUTA [23] consists of two parts: 1) Shannon entropy minimization loss (\mathcal{L}_{em}) and 2) negative sampling loss (\mathcal{L}_{mcc}). With the weighting hyper-parameter α , the overall loss function is denoted as follows:

$$\mathcal{L}_{SUTA} = \alpha \mathcal{L}_{em} + (1 - \alpha) \mathcal{L}_{mcc} \quad (1)$$

On the other hand, the unsupervised optimization objective of SGEM [24] consists of two parts: 1) generalized Rényi entropy minimization (\mathcal{L}_{GEM}) and 2) negative sampling loss (\mathcal{L}_{NS}). With a weighting hyper-parameter λ , the overall loss function is denoted as follows:

$$\mathcal{L}_{SGEM} = \mathcal{L}_{GEM} + \lambda \mathcal{L}_{NS}, \quad (2)$$

2.3. Off-the-shelf and Fine-Tuned ASR Models

We examined the effectiveness of TTA on two sets of models: 1) off-the-shelf models: off-the-shelf ASR models trained on LibriSpeech, which primarily include out-of-domain non-child speech data; 2) fine-tuned models: off-the-shelf ASR models fine-tuned on the training set of MyST with in-domain child speech data. Specifically, we used two off-the-shelf ASR models (Wav2vec2-base-960h¹ and Wav2vec2-large-960h²) that are trained on 960 hours of data from Librispeech [25]. To further investigate the need for continuous adaptation to each child speaker’s characteristics after fine-tuning, we fine-tuned both Wav2vec2-base and Wav2vec2-large on the MyST dataset’s training set, which consists of child speech. We used a learning rate of $1e - 5$ with the Adam optimizer for 50,000 steps, and a batch size of 32. The best model was selected based on the word error rate (WER) on the validation set, which was evaluated every 200 steps.

2.4. Datasets

This work used both the training and test set of My Science Tutor (MyST) dataset [26], currently one of the largest publicly available datasets for child speech recognition³. After removing recordings and utterances with missing annotations in the MyST dataset, we included data from 2622 children in the training set, and data from 438 children in the validation set, to fine-tune off-the-shelf ASR models. We used 91 childrens’ data in the test set of MyST to evaluate TTA methods in all model settings, as listed in Table 2. We only used utterances of less than 30s for the fine-tuning, following the setting in [18], resulting in 24.1 (SD=13.9) and 22.7 (SD=10.2) utterances for the training and validation dataset, respectively. For the test set, there were on average 140 (SD=99) utterances for each child speaker in the dataset. The duration of the utterances varied from less than 1 second to 111.4 seconds across all speakers. Overall, 129.3, 19.2, and 28.0 hours of data were included in the training, validation, and test set, respectively.

2.5. Evaluation Experimental Setup

We developed our codebase based on the open-sourced implementation from [24]⁴. All experimental settings such as learning rate and optimizer were kept consistent between settings. For both TTA methods, the adaptation step was set to $N = 10$. The weighting hyper-parameter was set to the following default values: $\alpha = 0.3$ for SUTA and $\lambda = 0.3$ for SGEM. We evaluated ASR models using word error rate (WER). Due to the data distribution imbalance – a small group of speakers had many utterances – we report the unweighted average WER based on the speaker rather than the utterance, so the results accurately represent the ASR system’s performance for each child.

3. Experiments, Results, and Discussion

3.1. Why is TTA needed for ASR models in child speech recognition?

Prior work has primarily reported average WER as a measure of model performance for child speech recognition [10–18], but overall average WER tends to overlook individual WER differ-

¹<https://huggingface.co/facebook/wav2vec2-base-960h>

²<https://huggingface.co/facebook/wav2vec2-large-960h>

³<https://catalog.ldc.upenn.edu/LDC2021S05>

⁴<https://github.com/drump/SGEM>

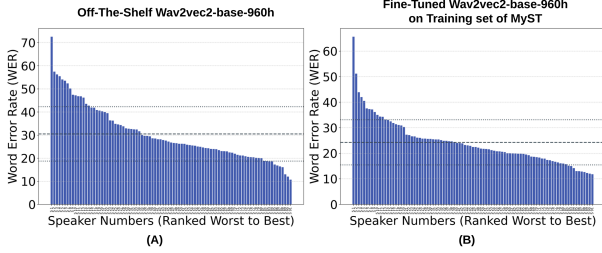


Figure 1: *Unadapted baseline performance for both LibriSpeech-off-the-shelf and MyST-fine-tuned Wav2vec2-base models in word error rate (WER).*

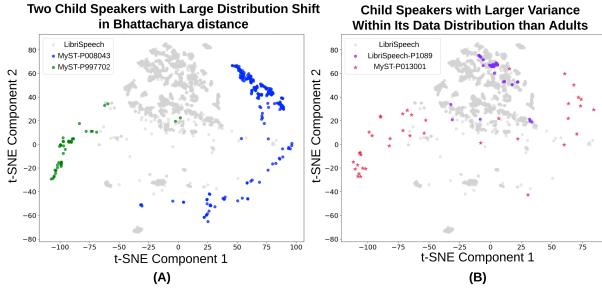


Figure 2: *Domain shift across and within child speakers. (A) Significant domain shift between child speakers; (B) Significant variance and domain shift within child speakers.*

ences between child speakers. As shown in Figure 1, before we applied test-time adaptation (TTA) methods, we first analyzed the performance of off-the-shelf ASR models for each child speaker individually. As shown in Figure 1, for both unadapted off-the-shelf and fine-tuned Wav2vec2-base, we discovered substantial differences in WER across speakers. By fine-tuning the model with child speech (training set of MyST), we observed noticeable decreases in both WER average (from 30.5% to 24.3%) and standard deviation (from 11.8% to 8.8%) across individual speakers, indicating that fine-tuning is effective to improve ASR performance across child speakers. However, as visualized in Figure 1(A) and (B), the differences quantified by the standard deviation still suggest that neither off-the-shelf nor fine-tuned ASR models generalize robustly across different child speakers, highlighting the need for more individual-level test-time adaptation to ensure model robustness for all speakers.

We also analyzed and visualized the MFCC features extracted using Librosa [27] using the T-distributed stochastic neighbor embedding (t-SNE). We calculated the pair-wise Bhattacharya distance across all pairs of child speakers in MyST data, and visualized the pair with one of the largest distances in Figure 2(A). The results validate that there may be a significant distribution shift between child speakers. As shown in Figure 2(B), we also calculated the variance within each child speaker’s embedding from MyST (Mean=1545.5, SD=915.3) and found that it is significantly larger than the variance within each adult speaker from Librispeech (Mean=259.1, SD=275.0). Our findings suggest that child speakers may also have larger distribution shifts both across and within their data distribution due to more expressive and personal speech characteristics, unique non-linguistic and unintelligible speech, and environment noise, which need to be addressed at test time. These findings further motivate the need for applying test-time adaptation for child speech recognition to ensure robust model gen-

Table 2: *WER comparisons of TTA methods (SUTA and SGEM) with the unadapted baseline on test set of MyST dataset in different off-the-shelf and fine-tuned model settings.*

Model Setting	TTA Method	WER	Δ	Stat. Sig.
Off-The-Shelf Wav2vec2-base	Unadapted	30.5%	—	—
	SUTA	27.5%	-3%	$p < .001$
	SGEM	27.2%	-3.3%	$p < .001$
Off-The-Shelf Wav2vec2-large	Unadapted	26.6%	—	—
	SUTA	25.1%	-1.5%	$p < .001$
	SGEM	25.3%	-1.3%	$p < .001$
Fine-Tuned Wav2vec2-base	Unadapted	24.3%	—	—
	SUTA	22.8%	-1.5%	$p < .001$
	SGEM	23.8%	-0.5%	$p < .001$
Fine-Tuned Wav2vec2-large	Unadapted	23.2%	—	—
	SUTA	22.5%	-0.7%	$p < .001$
	SGEM	23.4%	0.2%	$p < .001$

eralization to individual variability in real-world applications.

3.2. Can unsupervised TTA methods effectively adapt both off-the-shelf and fine-tuned ASR models for child speech recognition?

Off-the-shelf models. As shown in Table 2, on average, both SUTA and SGEM effectively reduced WER from the unadapted baseline. SUTA outperformed the unadapted baseline by 3% in WER when using the off-the-shelf Wav2vec2-base, and SUTA outperformed the unadapted baseline by 1.5% with the off-the-shelf Wav2vec2-large. In addition to reporting average model performance, we conducted individual two-sided Wilcoxon signed-rank tests between the TTA condition (N=91) and the unadapted condition (N=91), to further validate whether the improvements extend across all child speakers and are not driven by gains from a small group of speakers. We found that both SUTA and SGEM significantly improved from the unadapted baseline across different speakers for both Wav2vec2-base and Wav2vec2-large. These results suggest that TTA can help adapt ASR models from adult (LibriSpeech) to child (MyST) from the off-the-shelf models.

On an individual speaker level, as shown in Figure 3 (A) and (C), we visualized the performance gain enabled by SUTA over the unadapted baseline for each child speaker. A darker blue indicates a larger performance gain from TTA, and the speakers are numbered using the ranking of unadapted WER presented in Figure 1 (A). For off-the-shelf models, the majority of speakers benefited from TTA. In Figure 3(A), we found that S7, who benefited from SUTA the most when using off-the-shelf Wav2vec2-base, had a 10.3% gain with SUTA. In Figure 3(C), we found that S20, who benefited from SUTA the most when using off-the-shelf Wav2vec2-large, had a 7.1% gain with SUTA. However, we also observed inconsistent cases such as S3 (WER drop of 3.1%) in Figure 3(A) and S3 (WER drop of 12.5%), in Figure 3(C). These results suggest that although TTA effectively improved WER for most speakers, the current TTA method may still not be sufficiently capable of providing performance gains for every child speaker.

Fine-tuned models. A body of prior work has shown that fine-tuning off-the-shelf models on in-domain child speech data can effectively adapt for the adult-child domain gap. However, based on our findings from Section 3.1, we hypothesized that

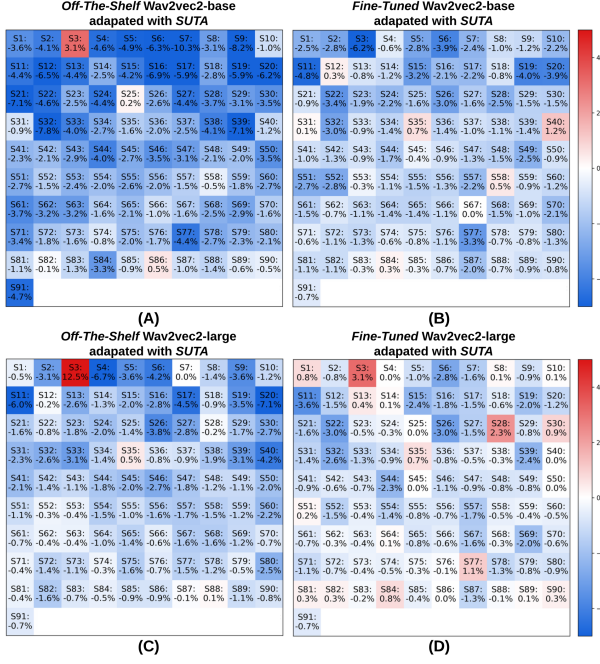


Figure 3: Heatplots of WER (%) performance gains for all 91 child speakers in four model settings adapted by SUTA. Note: speakers are numbered (1-91) using the ranking listed in Figure 1 (A).

the individual variability between each child speaker still needs to be adapted at test time. As shown in Table 2, our results validated our hypothesis. Even after in-domain fine-tuning, on average, SUTA outperformed the unadapted baseline by 1.5% in WER using the fine-tuned Wav2vec2-base, and by 0.7% in the most challenging case of using the fine-tuned Wav2vec2-large. We noticed that SGEM performed worse than SUTA for the fine-tuned ASR models, suggesting that SUTA may be the more robust TTA method across different model settings.

At the individual speaker level, as shown in Figure 3(B) and (D), the majority of the speakers also benefited from TTA for fine-tuned models. However, the performance gains were less than the corresponding off-the-shelf models as expected. In Figure 3(B), we found that S3, which benefited from SUTA the most when using the fine-tuned Wav2vec2-base, had a 6.2% gain with SUTA after fine-tuning. In Figure 3(D), we found that S11, which benefited from SUTA the most when using the off-the-shelf Wav2vec2-large, had a 3.6% gain with SUTA after fine-tuning. However, we also observed performance drops, such as S40 (WER drop of 1.2%) in Figure 3(A) and S3 (WER drop of 3.1%) in Figure 3(D). When using the fine-tuned Wav2vec2-large, as shown in Figure 3(D), 11 speakers experienced performance declines after adaptation with SUTA. While this result is expected given that fine-tuned Wav2vec2-large is the most accurate unadapted model setting, it further suggests that the current TTA methods may be less capable of selectively adapting to individual variability while preventing performance degradation for highly accurate models.

3.3. When do TTA methods perform well and when do they fail?

Background noise level can also be a source of individual variability that needs to be adapted at test time. To validate this, we

Table 3: Acoustic and linguistic feature analysis. We conducted a Spearman’s rank-order correlation test to examine the relationship between TTA’s performance gain and three acoustic and linguistic features: (1) effective mean squared (EMS) energy, and (2) word duration.

Model Setting	EMS Energy	Word Duration
Off-the-shelf Wav2vec2-base	$r = 0.22$ $p = .04$	$r = -0.30$ $p < .01$
Off-the-shelf Wav2vec2-large	$r = 0.02$ $p > .05$	$r = -0.56$ $p < .001$
Fine-Tuned Wav2vec2-base	$r = -0.09$ $p = 0.37$	$r = 0.41$ $p < .001$
Fine-Tuned Wav2vec2-large	$r = 0.23$ $p = 0.03$	$r = -0.29$ $p < .01$

conducted a Spearman’s rank-order correlation test with Holm-Bonferroni (HB) correction [28] between speakers’ effective mean squared (EMS) energy level during non-speech regions and TTA’s performance gain. We obtained the non-speech regions using Silero voice activity detector [29] and calculated EMS energy using Librosa [27]. We found moderate correlations in both the off-the-shelf Wav2vec2-base and fine-tuned Wav2vec2-large model settings, suggesting TTA may help to adapt for background noise level at test time.

In addition, we also examined the speech data of speakers who did not benefit from TTA, such as speaker S3 in Figure 3(A), (C) and (D). Comparing them qualitatively with speakers who experienced greater improvements from TTA, we observed that those with performance degradation tended to produce a high proportion of prolonged and non-linguistic speech. Since these speech segments were not reflected in the ground-truth transcription, they may have misled the optimization process during TTA. Quantitatively analyzing these speech patterns would require human annotations to accurately identify non-linguistic speech. However, we noticed that such non-linguistic speech often occurred before or after linguistic speech, effectively increasing the overall duration of the audio clip beyond what was necessary. To approximate this effect, we calculated the duration of the word as:

$$\text{Word Duration} = \frac{\text{Total Speech Duration (seconds)}}{\text{Number of Words in Transcription}}$$

Despite the limitations of this approximation, we found a significant correlation between speaking rate and TTA’s performance gain in all four modeling settings adapted by SUTA ($p < .01$) after HB correction, suggesting that the presence of non-linguistic speech may contribute to performance degradation in certain speakers.

4. Conclusion

In this work, we systematically examined the needs and performance of two widely used test-time adaptation (TTA) methods—SUTA, SGEM—to enable continuous adaptation of child speech data at test time. Our findings show that TTA significantly improved the performance of both off-the-shelf and their fine-tuned versions for child speech recognition, both on average and across individual child speakers, compared to unadapted baselines. Our acoustic and linguistic analysis also discovered that TTA may help adaptation to the individual background noise levels, while finding non-linguistic speech may be particularly challenging and misleading for test-time adaptation.

5. Acknowledgment

This work was supported by National Science Foundation (IIS-1925083), Simons Foundation (SFI-AR-HUMAN-00004115-03, 655054), and Sara Technology Inc. The authors alone are responsible for the content and conclusions.

6. References

- [1] W. Huang, K. F. Hew, and L. K. Fryer, “Chatbots for language learning—are they really useful? a systematic review of chatbot-supported language learning,” *Journal of Computer Assisted Learning*, vol. 38, no. 1, pp. 237–257, 2022.
- [2] T. Belpaeme, P. Baxter, J. De Greeff, J. Kennedy, R. Read, R. Looije, M. Neerinx, I. Baroni, and M. C. Zelati, “Child-robot interaction: Perspectives and challenges,” in *Social Robotics: 5th International Conference, ICSR 2013, Bristol, UK, October 27-29, 2013, Proceedings 5*. Springer, 2013, pp. 452–459.
- [3] W. Holmes and I. Tuomi, “State of the art and practice in ai in education,” *European Journal of Education*, vol. 57, no. 4, pp. 542–570, 2022.
- [4] M. Li, S. Zha, W. Gong, and Y. Jia, “A survey of human-computer interaction technology in intelligent home for children’s learning,” *Journal of Computer-Aided Design & Computer Graphics*, vol. 35, no. 2, pp. 248–261, 2023.
- [5] S. S. Narayanan and A. Potamianos, “Creating conversational interfaces for children,” *IEEE Transactions on Speech and Audio Processing*, vol. 10, no. 2, pp. 65–78 [IEEE Signal Processing Society Best Paper Award Winner, 2005], feb 2002.
- [6] A. Potamianos and S. S. Narayanan, “Robust recognition of children’s speech,” *IEEE Transactions on Speech and Audio Processing*, vol. 11, no. 6, pp. 603–616, nov 2003.
- [7] 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.
- [8] V. Bhardwaj, M. T. Ben Othman, V. Kukreja, Y. Belkhier, M. Bajaj, B. S. Goud, A. U. Rehman, M. Shafiq, and H. Hamam, “Automatic speech recognition (asr) systems for children: A systematic literature review,” *Applied Sciences*, vol. 12, no. 9, p. 4419, 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] S. Shraddha, S. Kumar *et al.*, “Child speech recognition on end-to-end neural asr models,” in *2022 2nd International Conference on Intelligent Technologies (CONIT)*. IEEE, 2022, pp. 1–6.
- [11] R. Jain, A. Barcovich, M. Yiwere, P. Corcoran, and H. Cucu, “Adaptation of whisper models to child speech recognition,” in *INTERSPEECH*, 2023.
- [12] J. Thienpondt and K. Demuynck, “Transfer learning for robust low-resource children’s speech asr with transformers and source-filter warping,” in *INTERSPEECH*, 2022.
- [13] T. Rolland, A. Abad, C. Cucchiari, and H. Strik, “Multilingual transfer learning for children automatic speech recognition,” in *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, N. Calzolari, F. Béchet, P. Blache, K. Choukri, C. Cieri, T. Declerck, S. Goggi, H. Isahara, B. Maegaard, J. Mariani, H. Mazo, J. Odijk, and S. Piperidis, Eds. Marseille, France: European Language Resources Association, Jun. 2022, pp. 7314–7320. [Online]. Available: <https://aclanthology.org/2022.lrec-1.795/>
- [14] A. A. Attia, D. Demszky, T. Ogunremi, J. Liu, and C. Espy-Wilson, “Continued pretraining for domain adaptation of wav2vec2.0 in automatic speech recognition for elementary math classroom settings,” *arXiv preprint arXiv:2405.13018*, 2024.
- [15] T. Rolland and A. Abad, “Exploring adapters with conformers for children’s automatic speech recognition,” in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 12 747–12 751.
- [16] R. Fan and A. Alwan, “Draft: A novel framework to reduce domain shifting in self-supervised learning and its application to children’s asr,” in *INTERSPEECH*, 2022.
- [17] W. Liu, Y. Qin, Z. Peng, and T. Lee, “Sparsely shared lora on whisper for child speech recognition,” in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 11 751–11 755.
- [18] R. Fan, N. B. Shankar, and A. Alwan, “Benchmarking children’s asr with supervised and self-supervised speech foundation models,” in *INTERSPEECH*, 2022.
- [19] A. A. Attia, J. Liu, W. Ai, D. Demszky, and C. Espy-Wilson, “Kid-whisper: Towards bridging the performance gap in automatic speech recognition for children vs. adults,” in *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, vol. 7, 2024, pp. 74–80.
- [20] R. Duan and N. F. Chen, “Senone-aware adversarial multi-task training for unsupervised child to adult speech adaptation,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 7758–7762.
- [21] R. Duan and N. Chen, “Unsupervised feature adaptation using adversarial multi-task training for automatic evaluation of children’s speech,” in *INTERSPEECH*, 2020.
- [22] D. Wang, E. Shelhamer, S. Liu, B. Olshausen, and T. Darrell, “Tent: Fully test-time adaptation by entropy minimization,” in *International Conference on Learning Representations*, 2021. [Online]. Available: <https://openreview.net/forum?id=uXl3bZLkr3c>
- [23] G.-T. Lin, S.-W. Li, and H.-y. Lee, “Listen, adapt, better wer: Source-free single-utterance test-time adaptation for automatic speech recognition,” in *INTERSPEECH*, 2022.
- [24] C. Kim, J. Park, H. Shim, and E. Yang, “Sgem: Test-time adaptation for automatic speech recognition via sequential-level generalized entropy minimization,” in *INTERSPEECH*, 2023.
- [25] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: an asr corpus based on public domain audio books,” in *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.
- [26] W. Ward, R. Cole, and S. Pradhan, “My science tutor and the myst corpus,” *Boulder Learning Inc*, 2019.
- [27] B. McFee, C. Raffel, D. Liang, D. P. Ellis, M. McVicar, E. Battenberg, and O. Nieto, “librosa: Audio and music signal analysis in python,” in *SciPy*, 2015, pp. 18–24.
- [28] S. Holm, “A simple sequentially rejective multiple test procedure,” *Scandinavian journal of statistics*, pp. 65–70, 1979.
- [29] S. Team, “Silero vad: pre-trained enterprise-grade voice activity detector (vad), number detector and language classifier,” <https://github.com/snakers4/silero-vad>, 2024.