

# The Sound of Syntax: Finetuning and Comprehensive Evaluation of Language Models for Speech Pathology

Fagun Patel<sup>1\*</sup>, Duc Q. Nguyen<sup>2\*</sup>, Sang T. Truong<sup>1\*</sup>,  
Jody Vaynshtok<sup>3</sup>, Sanmi Koyejo<sup>1</sup>, Nick Haber<sup>1</sup>

<sup>1</sup>Stanford University <sup>2</sup>National University of Singapore <sup>3</sup>Sound Speech and Hearing Clinic  
\* Co-first Authors

## Abstract

According to the U.S. National Institutes of Health, more than 3.4 million children experience speech disorders that require clinical intervention. The number of speech-language pathologists (SLPs) is roughly 20 times fewer than the number of affected children, highlighting a significant gap in children’s care and a pressing need for technological support that improves the productivity of SLPs. State-of-the-art multimodal language models (MLMs) show promise for supporting SLPs, but their use remains underexplored largely due to a limited understanding of their performance in high-stakes clinical settings. To address this gap, we collaborate with domain experts to develop a taxonomy of real-world use cases of MLMs in speech-language pathologies. Building on this taxonomy, we introduce the first comprehensive benchmark for evaluating MLM across five core use cases, each containing 1,000 manually annotated data points. This benchmark includes robustness and sensitivity tests under various settings, including background noise, speaker gender, and accent. Our evaluation of 15 state-of-the-art MLMs reveals that no single model consistently outperforms others across all tasks. Notably, we find systematic disparities, with models performing better on male speakers, and observe that chain-of-thought prompting can degrade performance on classification tasks with large label spaces and narrow decision boundaries. Furthermore, we study fine-tuning MLMs on domain-specific data, achieving improvements of over 30% compared to base models. These findings highlight both the potential and limitations of current MLMs for speech-language pathology applications, underscoring the need for further research and targeted development<sup>1</sup>.

<sup>1</sup>To support continued progress, we publicly release our datasets, fine-tuned models, and benchmarking framework. Code: [https://github.com/sangttruong/slp\\_benchmark](https://github.com/sangttruong/slp_benchmark), Dataset: <https://huggingface.co/datasets/SAA-Lab/SLPHelmUltraSuitePlus>.

## 1 Introduction

Speech and language pathologies (SLP) in children can significantly impact communication, academic development, and long-term social outcomes (Hitchcock et al., 2015; Foster et al., 2023). Early detection and intervention by speech-language pathologists are critical to mitigating these adverse effects (Gibbard et al., 2004; Centers for Disease Control and Prevention, 2024). Unfortunately, the availability of qualified clinicians is characterized by an uneven distribution across geographic and socioeconomic contexts, with only one expert for every 20 affected children, resulting in significant disparities in access to care and leading to “missing intervention” for many children who could benefit from timely support (U.S. National Institute on Deafness and Other Communication Disorders, 2025; Tucker and McKinnon, 2020). This gap highlights an urgent need for scalable and supportive technological solutions to augment clinicians’ capabilities and expand the reach of vital interventions.

The shortage of qualified clinicians has led to significant gaps in diagnostic capacity, particularly in domains that require specialized expertise. Recent advancements in multimodal large language models (MLMs) present a promising opportunity to partially automate or augment diagnostic workflows (Lammert et al., 2025; Bhattacharya et al., 2024; Nagpal et al., 2025; Maqsood et al., 2024). Multimodal LLMs, such as GPT-4 and Gemini, exhibit state-of-the-art capabilities in speech processing and contextual reasoning. Effective integration of LLMs into clinical SLP workflow requires rigorous evaluation to establish their clinical validity (Cordella et al., 2025). This process relies on comprehensive and high-quality datasets that capture the variability of pediatric speech, especially disordered forms, and are annotated with clinical features. Currently, the evaluation is hindered by

two key challenges: the scarcity of well-curated pediatric speech corpora and the lack of comprehensive evaluation frameworks for analyzing children's speech (Suh et al., 2024).

In this study, we present a comprehensive evaluation of MLMs in SLP. We develop a systematic procedure for annotating symptoms and disorders in child speech, creating resources suitable for SLP-focused model evaluation and fine-tuning. The benchmark assesses models across five clinical scenarios covering a spectrum of tasks from foundational disorder detection to more granular symptoms, including Disorder Diagnosis, Transcription-Based Diagnosis, Transcription, Disorder Type Classification, and Symptom Classification. We systematically evaluate the capabilities and limitations of existing MLMs in SLP-relevant contexts and explore approaches for domain-specific adaptation. Our analysis reveals substantial performance gaps, with macro-F1 scores frequently falling below clinically acceptable thresholds, particularly on fine-grained tasks such as Disorder Type and Symptom Classification. To mitigate these shortcomings, we develop and assess fine-tuned MLMs, demonstrating significant improvements and advancing the state of the art on these specialized clinical tasks. Our contributions are summarized as follows.

- We release four curated pediatric speech datasets comprising approximately 30,000 speech samples across English and French, encompassing both typical and disordered speech. These datasets provide a publicly available, high-quality resource to support reproducible benchmarking in SLP.
- We introduce the first comprehensive evaluation framework for SLP, encompassing five essential clinical tasks. This framework enables consistent, task-aligned evaluation and facilitates direct comparison of speech LLM performance under a standardized protocol.
- We introduce fine-tuned speech LLMs that achieve state-of-the-art performance across all evaluated SLP tasks, illustrating the efficacy of domain-specific adaptation in enhancing diagnostic and transcriptional capabilities.
- We conduct extensive fine-grained analyses on model performance across various conditions, including demographic factors (gender, age), languages, and reasoning paradigms

(e.g., Chain-of-Thought), offering deeper insights into model robustness and potential biases.

## 2 Related Works

**Comprehensive AI Benchmarking** Comprehensive benchmarking has been instrumental in advancing speech-health research. The ADRess Challenge (Luz et al., 2020) established a balanced benchmark for Alzheimer's detection from spontaneous speech, standardizing evaluation via F1 and MMSE-regression metrics. Similarly, the Children's ASR Benchmark (Fan et al., 2024) introduced standardized splits and Whisper/Wav2Vec baselines for speech recognition in children aged 6-14, highlighting age-specific acoustic challenges. Nonetheless, systematic benchmarking of MLMs in clinical contexts remains limited.

**Finetuning Audio Models** Recent approaches have shown finetuning led to improvement in LLM audio understanding performance, especially for low-resource languages (Pillai et al., 2024). Models finetuned via approaches like Continued Pre-training (Ke et al., 2023) and Reasoning Preference Optimization (RPO) (Pang et al., 2024) have showcased improved reliability, factual accuracy, and cross-language generalization in clinical use cases (Kawakami et al., 2025). Additionally, finetuned LLMs have been shown to outperform classical DNN approaches across clinical use cases such as murmur detection (Florea et al., 2025). Furthermore, research has shown that while fine-tuning on specialized medical data enhances domain-specific knowledge, it can negatively impact a model's long-context understanding (Yang et al., 2024), highlighting the need for a balanced data composition during the fine-tuning process.

## 3 Method

### 3.1 Datasets, Annotation, and Core Scenarios

We use four datasets including Ultrasuite (Eshky et al., 2018), ENNI (Schneider et al., 2006), LeNormand (Le Normand, 1997), and Percept-GFTA (Benway et al., 2022). These datasets are publicly available, having been collected with informed consent and subsequently anonymized by their creators to protect participant privacy. The datasets encompass a range of child speech samples, both typical and disordered, and serve as the foundation for evaluating model performance

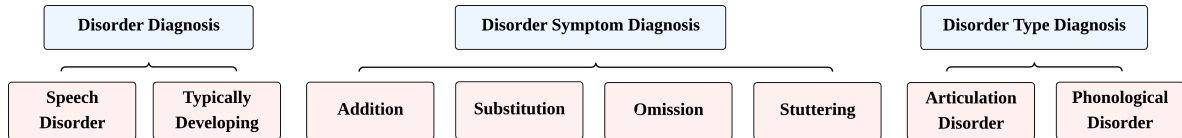


Figure 1: Taxonomy of classification tasks in SLPHelm. The benchmark includes three core diagnostic tasks: (i) disorder diagnosis, (ii) disorder symptom classification, and (iii) disorder type classification.

across diagnostic tasks. Their summary statistics are in Table 1. We randomly sample up to 1000 instances from each dataset for evaluation.

While foundational, existing datasets for children’s speech-language pathology research are largely confined to transcription and binary classification of speech as either disordered or typical (Benway et al., 2022; Eshky et al., 2018; Le Normand, 1997; Schneider et al., 2006). This scope is critically insufficient for developing and validating models for more nuanced clinical applications, such as identifying specific disorder types or their underlying symptoms/categories for which no large-scale, publicly available datasets currently exist. This lack of fine-grained annotation represents a significant bottleneck, impeding the development of automated tools that can support SLPs in differential diagnosis and therapy planning. To address this foundational gap, we collaborated closely with a certified SLP professional to develop a detailed annotation schema that captures both disorder types and their characteristic symptoms. For each speech sample, we assign the most prominent disorder type and symptom, prioritizing the most salient diagnostic features when multiple conditions may co-occur. Speech samples exhibiting no observable signs of speech disorder are annotated as typical. Our annotation protocol and chosen taxonomy are informed by clinical guidelines from the U.S. National Institutes of Health (Simon and Rosenbaum, 2016) and SLP best practices (American Speech-Language-Hearing Association, 2016). After initial manual labeling by our team, we conducted a verification phase in which the speech-language pathologist reviewed annotations to ensure consistency and clinical validity.

We evaluate five core tasks that collectively capture the essential stages of pediatric SLP, from initial screening to detailed diagnostic analysis: **(1) Disorder Diagnosis**, which assesses a model’s ability to distinguish between typical and disordered speech a critical early triage step for prioritizing clinical resources; **(2) Transcript-based Diagno-**

**sis**, which serves as a baseline for diagnostic accuracy by testing the assumption that speech from children with disorders deviates from expected utterances. This approach operates by matching model-generated transcripts to required spoken text, which offers a minimal, interpretation-free method that could be readily deployed in clinical settings. By benchmarking against this heuristic, we quantify the value added by more sophisticated multimodal LLMs; **(3) Transcription**, which measures the fidelity of automatic speech recognition (ASR) systems on child with disordered speech, a prerequisite for downstream diagnostic and documentation tasks; **(4) Disorder Type Classification**, which probes whether models can differentiate between *articulation disorders* motor-based speech errors such as lisps and *phonological disorders*, which involve rule-based sound pattern errors such as omission of a part of a word shape (e.g., *ca/cat*) or consistent shifts in the place of production (e.g., */k/ → /t/*). For example, a child who consistently replaces final consonants, such as saying “gape” for “gate,” is exhibiting a pattern-based phonological process. A successful model must identify this underlying rule, distinguishing it from an articulation disorder where a child might struggle to physically produce the */t/* sound in any context; **(5) Disorder Symptom Classification**, a more granular task, requires models to pinpoint the specific clinical symptoms that constitute these error patterns. This includes identifying substitutions, such as the final consonant changing in the production of “gape” for “gate”, or omissions, where a child might say “gore” instead of “gorge” by dropping the final sound. The task also includes identifying additions (the insertion of extra sounds) and stuttering (disruptions in speech fluency). Accurate identification of these individual symptoms is a critical prerequisite for downstream diagnostic analysis and informs the design of targeted therapy plans. Figure 1 illustrates an overview of classification tasks in our pipeline. Details prompts are presented in Appendix F.

Evaluation metrics for classification tasks include Macro F1, Micro F1, and Exact Match Accuracy, while transcription performance is assessed using Word Error Rate (WER), Match Error Rate (MER), and Word Information Preserved (WIP). The Macro F1 and Micro F1 scores were specifically chosen to provide a comprehensive assessment of classification performance. The Macro F1 score evaluates the model’s average performance on each class equally, ensuring that performance on rare but clinically significant disorder categories is not overlooked. In contrast, the Micro F1 score aggregates performance across all individual samples, offering a measure of overall classification correctness.

We implement two distinct model inference pipelines. The first, referred to as the audio-to-LLM prompting pipeline, is designed for models with native multimodal capabilities (e.g., GPT-4o-Audio, Gemini 2.0 Flash). Here, raw audio inputs are passed directly to the model alongside a task-specific prompt. The second pipeline, termed transcription-based prompting, targets language-only models (denoted with the -transcribe suffix). Here, audio inputs are first transcribed using a base automatic speech recognition (ASR) model (e.g., Whisper or GPT-4’s internal ASR), and the resulting text is embedded into a structured prompt for downstream reasoning. This two-pronged architecture compares models with native audio understanding to those using ASR-to-LLM pipelines, highlighting trade-offs between direct and transcription-mediated processing.

### 3.2 Models and Finetuning

We evaluate 15 speech models. Among the closed-source LLMs, we study GPT-4 family (4o-audio, 4o-mini-audio, 4o-transcribe, and 4o-mini-transcribe), Whisper, and the Gemini 2.0 family (2.0-flash, 2.0-flash-lite). For open-source models, we study Qwen families (2.5-omni-7b, 2.5-omni-3b, 2-audio-7b, audio-chat), the Phi-4, and IBM Granite series (3.3-8b, 3.3-3b, 3.2-8b). These models were chosen to cover a range of model sizes and families across both closed and open sources.

To investigate the impact of fine-tuning on model performance across multiple tasks, we explore two fine-tuning strategies. Our first strategy involves fine-tuning the model on a speech recognition task (Scenario 3, as described above), relying on the models intrinsic ability to transfer knowledge to improve performance on related tasks. In this setup,

both typical and disordered speech samples are labeled with the same expected transcriptions. However, assigning identical transcriptions to acoustically distinct inputs may introduce ambiguity and limit the model’s ability to learn disorder-specific patterns. To mitigate this, our second strategy - Finetuned w/ Markers - modifies the labeling of disordered speech by appending an asterisk to each word in its transcription. This lightweight labeling scheme serves to differentiate disordered speech from typical speech, thereby guiding the model to better recognize and transcribe disordered speech patterns without altering the overall task formulation. Details of fine-tuning prompts and hyperparameters are presented in Appendix D.

We hypothesize that fine-tuning on a general task (e.g., speech recognition) alone is insufficient to yield improvements on specialized clinical tasks unless the fine-tuning data contains explicit information relevant to those tasks. This stems from the theoretical premise that general-purpose models primarily optimize for surface-level acoustic-linguistic alignment, which may not encode the deeper, disorder-specific features, such as atypical phonological patterns or motor-based distortions, necessary for clinical inference (Shor et al., 2019; Dorfner et al., 2024).

## 4 Experiment

### 4.1 Per Scenario Results

Our findings suggest that both existing proprietary and open-source models currently fail to meet clinically acceptable performance thresholds. This limitation is likely attributable to the underrepresentation of disordered speech in training corpora, as such data is significantly less prevalent than typical speech samples available online. Existing FDA-approved diagnostic systems typically achieve F1 scores in the range of 0.80 to 0.85 (Fanni et al., 2023; Abramoff et al., 2018), which serves as a practical standard for clinical viability. Furthermore, model performance varies across different task scenarios, highlighting the absence of a universally robust model that can consistently address the diverse requirements of pediatric SLP applications. Figure 2 presents an overview of the performance of all models.

**Scenario 1: Disorder Diagnosis** In the disorder diagnosis task, performance remains limited, with no model exceeding a micro F1 score of 0.56. The



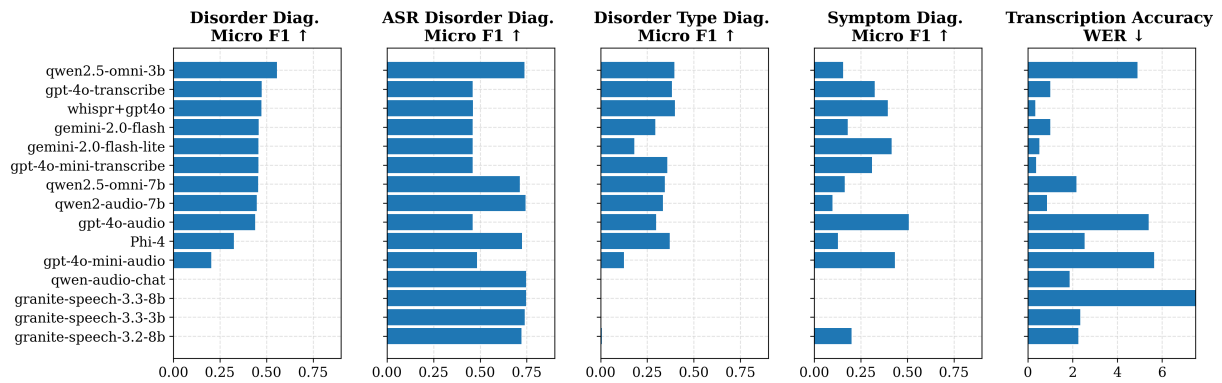


Figure 2: Metrics across all scenarios

best result is achieved by Qwen 2.5-Omni-3B, outperforming GPT-4o-Mini-Transcribe. The fact that these models use different pipelines audio-grounded vs. ASR+text suggests no clear advantage of one approach over the other. Smaller variants within each family exhibit similar performance. Audio-grounded Granite models perform poorly ( $F1 < 0.1$ ), likely due to their pretraining focus on speech-to-text and translation tasks (IBM Granite Team, 2025).

**Scenario 2: Transcription-based Diagnosis** In this scenario, we evaluate a naive baseline that performs diagnosis on transcribed text rather than directly from audio. Counter-intuitively, this two-stage “transcribe-and-compare” approach does not degrade performance compared to end-to-end acoustic reasoning. In fact, for a majority of the models, it proves to be a superior strategy, yielding higher F1 scores. This surprising result suggests that potential error propagation from the ASR system is less of a hindrance than the inherent difficulty models face in performing complex diagnostic reasoning directly on acoustic features. Consequently, this simple transcription-based method, far from being a conservative lower bound, establishes itself as a highly competitive and often preferable baseline for diagnosing speech disorders.

**Scenario 3: Transcription** WER varies widely from 0.33 to 6+, with the Gemini and OpenAI family of models performing the best. Importantly, transcription fidelity shows limited correlation with diagnostic accuracy: qwen2.5 omni 3b records a poor WER of 4.9 yet is consistently one of the strongest in classification tasks. These findings indicate that high-quality transcripts are neither necessary nor sufficient for dependable clinical reasoning.

**Scenario 4: Disorder Type Classification** Here, Whisper-GPT4o outperforms all models regardless of access or scale, hinting at architectural or pre-training advantages. The transcript-based inference consistently outperforms its audio-grounded counterparts. While for the larger GPT-4o model, audio-grounded inference performs close to its transcript-based counterpart, this trend reverses dramatically for the smaller 4o-mini model, where the ASR+LLM pipeline significantly outperforms the audio-grounded inference. This suggests that while large models can leverage the rich signals in raw audio, smaller models benefit substantially from the structured and potentially less noisy input of a transcript.

**Scenario 5: Disorder Symptom Classification** Accurate identification of these symptoms directly informs treatment goals and therapy design in speech-language pathology. GPT-4o leads in performance but still falls well short of clinically actionable accuracy. Moreover, transcription-based models underperform across all metrics, underscoring that critical acoustic cues needed for symptom detection are often lost or degraded during transcription.

## 4.2 Robustness Analysis

**The Effect of Finetuning** Fine-tuning large models can significantly enhance their performance on downstream tasks. Fine-tuning solely on automatic speech recognition (ASR) data, regardless of whether disordered speech is explicitly marked, leads to noticeable improvements in ASR-based tasks (Scenarios 2 and 3). However, not differentiating between typical and disordered speech introduces ambiguity in the input-label mapping, which degrades performance in other scenarios. Incorporating a simple asterisk mitigates this issue,

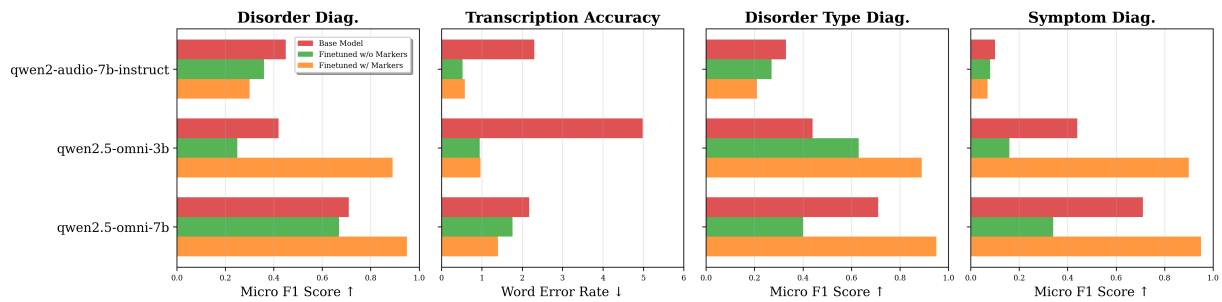


Figure 3: Model performance after finetuning

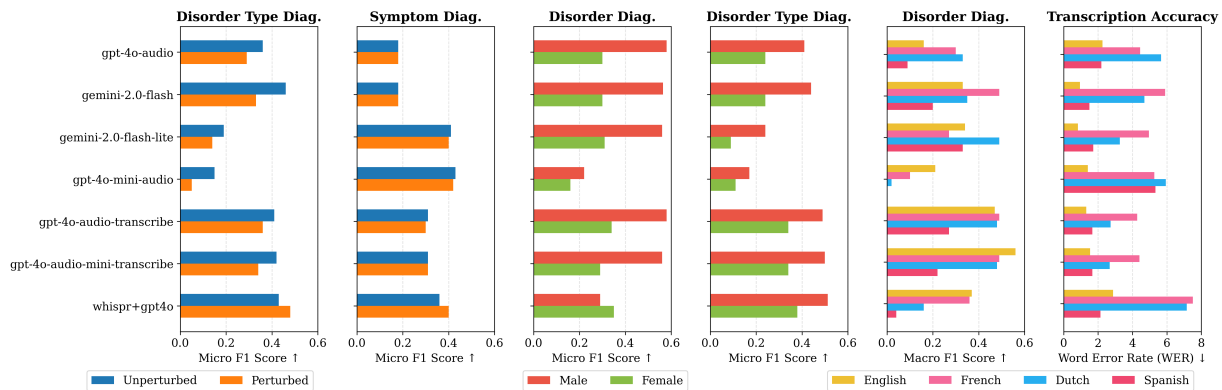


Figure 4: Model performance across robustness under noisy conditions, across gender, and languages

yielding more stable performance.

**The Effect of Gender** Figure 4 showcases that across two diagnostic tasks, the models exhibit a systematic gender performance gap that favors male speech. For both tasks, we run model evaluation on 1000 utterances for each gender on the UltraSuite dataset since it makes demographic identifiers available through its metadata. The results showcase a remarkably consistent pattern: almost every model posts a positive male-female differential. Notably, the gap is not confined to a particular modeling strategy; it appears in fully audio-grounded systems as well as in transcript-conditioned variants. Persistent speech-sound disorders are more commonly associated with boys than girls (Wren et al., 2016), which can lead to an imbalance in training data and downstream model performance. The magnitude of the divergence suggests practical consequences for clinical deployment, as female speech receives both lower sensitivity and lower precision across disorder categories. Taken together, the results underscore the need for targeted auditing and, potentially, gender-balanced fine-tuning to ensure equitable diagnostic performance across child speakers.

**The Effect of Language** Figure 4 shows a diverging pattern: classification accuracy shows no clear trend between languages, yet WER in French and Dutch is markedly worse. A plausible explanation lies in the way these systems were pre-trained. Their ASR components are heavily optimized on English text-to-speech pairs, so lexical recognition degrades when confronted with French/Dutch phonotactics, inflating WER. By contrast, the diagnostic classifiers operate on higher-level acoustic embeddings learned during large-scale audio pre-training that is largely language-agnostic (Klempí and Krupika, 2024). Those embeddings could still capture phonological and articulatory cues relevant to speech-disorder detection, so classification accuracy can be divergent from word-level transcription accuracy. This lack of correlation between the diagnostic capabilities of a model and its performance under transcription tasks highlights that while the classification performance of some models generalizes across languages, transcription is highly language-dependent.

In a preliminary analysis to evaluate model performance on tonal languages, we assessed disorder diagnosis capabilities on datasets of Taiwanese (Tsay, 2007) and Cantonese (Edwards and Beckman, 2008) speech. The results indicate a sig-

nificant failure of current models to generalize to these linguistic contexts. As shown in Figure 8, performance is markedly degraded, with the best-performing model achieving a Micro F1 score below 0.10 in Taiwanese, while performance in Cantonese approached zero across all models.

Given that these evaluation datasets consisted entirely of speech from typically developing children, the low F1 scores reveal that the models systematically misdiagnosed the vast majority of samples as disordered. This suggests that the acoustic features these models rely on for disorder detection may be conflating the inherent tonal variations of these languages with pathological speech patterns. These findings underscore a critical limitation, indicating that current architectures are unable to accurately comprehend low-resource tonal languages within the SLP context.

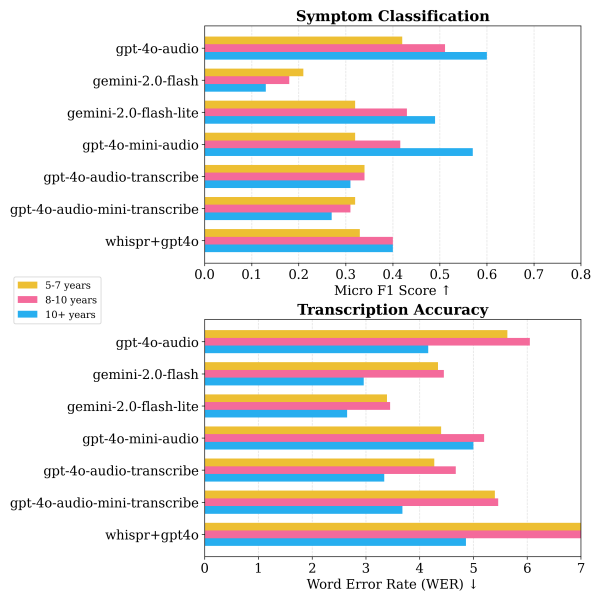


Figure 5: Model performance across age groups

**The Effect of Age** We analyzed model performance across three age ranges: early-elementary (5-7 years old), mid-elementary (8-10 years old), and post-elementary (10+ years old). The Ultra-Suite dataset provided structured age data for the children in its dataset and was used to run this analysis. The results reveal a distinct trend in symptom classification. Audio-native models demonstrate a significant drop in performance for younger children, as indicated by a lower Micro F1 Score for the 5-7 year-old group. In contrast, ASR+LLM pipelines maintain more consistent classification accuracy across all age ranges. For transcription accuracy, nearly all models performed best on the

speech of the oldest children (10+). Performance gradually decreased for the younger age brackets.

The results showcase that variations in younger children’s speech pose a significant challenge for audio-native models in complex classification tasks. The models are likely optimized for adult speech, and performance improves as children’s speech patterns approach this norm. While transcription accuracy is also best for older children, the two-stage ASR+LLM approach appears more resilient to age-related speech variations for downstream analytical tasks like symptom classification, given that it relies on the model’s text understanding capabilities to make its classification rather than audio understanding.

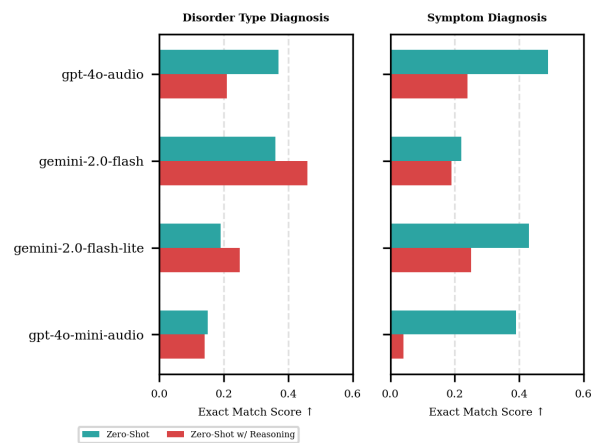


Figure 6: Compares model’s classification performance with/without CoT prompting

**The Effect of Reasoning** We analysed the effects of introducing an explicit chain-of-thought (CoT) prompt across the Gemini 2.0 and GPT 4o model families. We saw that CoT systematically depressed F1 scores on the symptom diagnosis task and produced a mixed picture on the disorder type diagnosis, as can be seen in Figure 6. The pattern aligns with recent evidence that CoT can hamper tasks where the optimal decision boundary is compact or where answer formatting is unforgiving, because the additional reasoning tokens introduce distraction or bleed into the predicted label (Liu et al., 2024). However, even when CoT degrades final accuracy, the explicit rationale provides crucial insight into a model’s failure modes. By inspecting the reasoning trace, we can diagnose why a model arrived at its conclusion, rather than simply observing the error itself. We analysed 200 reasoning traces, each for the 4 models for the Disorder-Type scenario. The CoTs are analysed along the follow-

ing axes to reveal a clear profile of the model’s reasoning process (Gandhi et al., 2025). A subset of the analysed traces can be found in Appendix G.

- **Subgoal Setting**, where a complex problem is broken down into manageable steps (e.g., “To solve this, we first need to...”).
- **Rule following** or the application of taxonomical rules during classification (e.g. repeated errors in the last consonant indicate a phonological disorder rather than articulation).
- **Error Detection** or whether the model can recognize individual errors in speech, on which it will determine its classification.

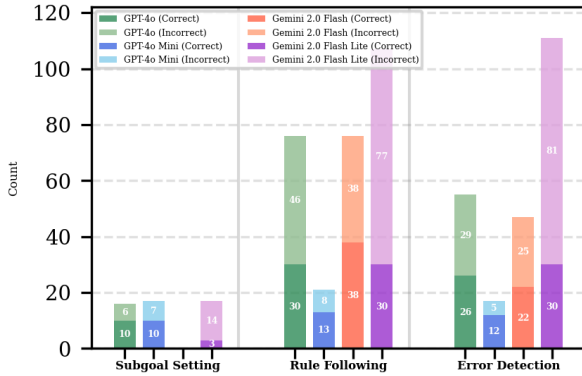


Figure 7: CoT Behaviour Analysis - Breakdown of behaviors showcased by the model. The bar for each model indicates the number of times the given behaviour is observed. Each bar is broken down by the times the given behaviour resulted in a correct/incorrect answer.

As shown in Figure 7, our analysis of CoT reasoning reveals key differences in the models’ problem-solving approaches. Rule Following and Error Detection were the most prominent strategies employed across all models. However, the effective application of these strategies varied significantly. For GPT-4o, a majority of reasoning traces using Rule Following and Error Detection led to an incorrect final answer. This high rate of flawed application likely explains its reduced classification performance with CoT. Conversely, Gemini 2.0 Flash’s improved performance correlates with more reliable reasoning; its application of Rule Following, for example, was correct in 50% of instances, compared to 39.5% for GPT-4o. This suggests that the fidelity of a model’s reasoning directly impacts its accuracy. Interestingly, Gemini 2.0 Flash Lite exhibits the highest frequency of Rule Following and Error Detection

behaviors but with the lowest accuracy (28% and 27%, respectively). Despite this, its classification performance improves, highlighting that the frequency of these model behaviours also contributes to overall model performance. This analysis indicates that for SLP tasks, CoT prompting improves performance only when the label space is limited with well-defined boundaries and the model can reliably apply rule-following and error-detection strategies in its reasoning traces.

**Effect of Ensemble** We implemented an ensemble strategy that combines predictions from multiple models. For classification tasks, the ensemble uses a majority voting system where the final prediction is determined by selecting the most frequently predicted label across all models, effectively leveraging the collective wisdom of multiple classifiers to reduce individual model biases and improve robustness. We created ensembles of three models to avoid tie scenarios where possible. If encountered, a tie is broken by arbitrarily selecting one of the predicted labels.

The ensemble analysis in Table 2 reveals a nuanced performance landscape. For the primary task of disorder diagnosis, a performance plateau is observed, with three distinct ensembles - the “Qwen Pure Ensemble” and both “Google + OpenAI” combinations achieving similar F1 scores. However, significant differences emerge in the more granular classification tasks. The mixed-vendor Google + OpenAI ensembles demonstrate substantially better performance in symptom identification, achieving a Micro F1 score of 0.389, more than double that of the Qwen-based ensembles. Conversely, the Qwen ensembles show a slight advantage in disorder type classification. These findings suggest that a simple ensemble strategy does not guarantee superior performance, and that single model family ensembles are not guaranteed to be dominant. The effectiveness of an ensemble appears to be highly task-dependent and may benefit from combining models with diverse architectures, as seen with the Google and OpenAI ensembles’ success in symptom classification. This highlights that the optimal ensemble configuration is not merely a combination of the strongest individual models but requires empirical testing to find complementary strengths for specific tasks.

**The Effect of Background Noise** We analyze model robustness by evaluating their performance with three added artificial perturbations, including



road noise, classroom noise, and office noise. We aggregate the results to assess model performance. Each type of background noise was added at 20 dB. As shown in Figure 4, model performance degrades on average by 10 basis points across all models for disorder type diagnosis, while performance on symptom diagnosis remains virtually unchanged. These observations suggest that noise resilience is not strictly determined by model and inference architecture. It might instead be influenced by model-specific factors such as design, scale, and training data.

## 5 Conclusion & Future Work

Our empirical findings reveal that even the best-performing models fall short of clinical-grade reliability, revealing considerable room for improvement. This performance gap presents a significant clinical risk, as the deployment of an unreliable model could lead to patient misdiagnosis or delayed treatment. Our fine-tuning experiments with the Qwen2.5 family demonstrate that this performance can be substantially improved, highlighting the effectiveness of task-specific adaptation and the potential for developing specialized SLP models that generalize well across tasks.

Our robustness analysis identifies a consistent performance disparity favoring male speakers, contradicting the clinical principle of equitable care and highlighting the need for bias mitigation through strategies such as gender-balanced fine-tuning and targeted data augmentation. Cross-linguistic evaluations further demonstrate that audio-grounded models retain competitive diagnostic performance even when transcription quality degrades, suggesting that higher-order acoustic features support language-agnostic reasoning.

Future work could expand coverage to low-resource languages and neurodiverse populations, and evaluate model explanations for clinical faithfulness. We also plan to investigate privacy-preserving fine-tuning paradigms to facilitate deployment in sensitive pediatric settings. Collectively, these directions aim to bridge the gap between promising laboratory advances and the development of clinically robust, ethically sound AI systems for SLP.

## Limitations and Ethical Considerations

Despite the comprehensive scope of our benchmark, several limitations and core ethical consid-

erations warrant explicit discussion prior to any consideration of clinical deployment.

### 1. Privacy and Data Sensitivity

The ethical considerations surrounding privacy and consent are central to work in sensitive clinical domains. While this study was conducted on public datasets collected with informed consent and subsequently anonymized to protect participant privacy, our current setup does not incorporate a formal privacy-preserving learning or evaluation framework. Addressing this is essential to safeguard patient data and build trust among clinicians and families. As such, future work must investigate privacy-preserving fine-tuning paradigms to facilitate responsible deployment in pediatric settings.

### 2. Bias and Equity

A primary finding of this work is the identification of significant performance disparities, which present a direct challenge to the principle of equitable care.

- **Gender Bias:** Our evaluation revealed a consistent performance disparity that favors male speakers across multiple models. This indicates an urgent need for bias mitigation strategies, such as gender-balanced fine-tuning and targeted data augmentation, which were not implemented in this study but are critical next steps.
- **Data Representation:** The datasets employed are drawn primarily from English and French speakers, leading to an underrepresentation of other languages and dialects. This limits the generalizability of our findings and risks creating tools that are not effective for more linguistically and culturally diverse populations.
- **Annotation Bias:** The annotation process, though guided by an SLP professional, required prioritizing the most prominent disorder when multiple conditions co-occurred, introducing a degree of subjectivity.

### 3. Clinical Reliability and Validation

A core ethical requirement for deployment is robust clinical validation.

- **Performance Gap:** This study establishes a benchmark for clinical viability based on existing FDA-approved systems (F1 scores of

0.80 to 0.85). Our findings demonstrate that even the best-performing models currently fall short of this standard.

- **Risk of Misdiagnosis:** This performance gap represents a significant clinical risk, as the deployment of an unreliable model could lead to patient misdiagnosis or delayed treatment. The results underscore that current MLMs are not yet clinically robust and require substantial further development and validation before they can be considered for supportive roles in real-world SLP workflows.

## Acknowledgments

SK acknowledges support by NSF 2046795 and 2205329, IES R305C240046, ARPA-H, the MacArthur Foundation, Schmidt Sciences, Stanford HAI, and Google. NH acknowledges the National AI Institute for Exceptional Education (Institute of Education Sciences, U.S. Department of Education, through Grant 22298673 (NSF)).

## References

- Michael D. Abràmoff, Patrick T. Lavin, Mary Birch, Nicholas Shah, and John C. Folk. 2018. [Pivotal trial of an autonomous AI-based diagnostic system for detection of diabetic retinopathy in primary care offices](#). *NPJ Digital Medicine*, 1:39.
- American Speech-Language-Hearing Association. 2016. Speech sound disorders: Articulation and phonology. <https://www.asha.org/practice-portal/clinical-topics/articulation-and-phonology/>. [Practice Portal].
- Nicholas R. Benway, Jonathan L. Preston, Elaine R. Hitchcock, Adam Salekin, Harsh Sharma, and Tara McAllister. 2022. PERCEPT-R: An Open-Access American English Child/Clinical Speech Corpus Specialized for the Audio Classification of R. In *Proceedings of Interspeech 2022*, pages 2408–2412, Incheon, Republic of Korea. International Speech Communication Association (ISCA).
- Anish Bhattacharya and 1 others. 2024. [Clinical efficacy of pre-trained large language models through the lens of aphasia](#). *Scientific Reports*, 14(1):15573.
- Centers for Disease Control and Prevention. 2024. Why act early if you're concerned about development? <https://www.cdc.gov/ncbddd/actearly/whyActEarly.html>. Accessed May 10, 2025.
- Claire Cordella, Manuel J. Marte, Hantian Liu, and Swathi Kiran. 2025. [An introduction to machine learning for speech-language pathologists: Concepts, terminology, and emerging applications](#). *Perspectives of the ASHA Special Interest Groups*, 10(2):432–450.
- Felix J Dorfner, Amin Dada, Felix Busch, Marcus R Makowski, Tianyu Han, Daniel Truhn, Jens Kleesiek, Madhumita Sushil, Jacqueline Lammert, Lisa C Adams, and 1 others. 2024. Biomedical large language models seem not to be superior to generalist models on unseen medical data. *arXiv preprint arXiv:2408.13833*.
- Jan Edwards and Mary E. Beckman. 2008. Methodological questions in studying consonant acquisition. *Clinical Linguistics and Phonetics*, 22(12):937–956.
- Paul E Engelhardt and 1 others. 2021. [A review of automated speech and language features for assessment of cognitive and thought disorders](#). *Frontiers in Psychology*.
- Ahmed Eshky, Mário S. Ribeiro, Jane Cleland, Korin Richmond, Zoe Roxburgh, James Scobbie, and Alan Wrench. 2018. Ultrasuite: A repository of ultrasound and acoustic data from child speech therapy sessions. In *Proceedings of Interspeech*, pages 2342–2346, Hyderabad, India. International Speech Communication Association (ISCA).
- Zhaoxi Fan and 1 others. 2024. A benchmark for automatic speech recognition on child speech. In *Proc. INTERSPEECH*.
- S. C. Fanni, A. Marcucci, F. Volpi, S. Valentino, E. Neri, and C. Romei. 2023. [Artificial intelligence-based software with ce mark for chest x-ray interpretation: Opportunities and challenges](#). *Diagnostics*, 13(12):2020.
- Adrian Florea, Xilin Jiang, Nima Mesgarani, and Xiaofan Jiang. 2025. [Exploring finetuned audio-llm on heart murmur features](#). *Preprint*, arXiv:2501.13884.
- M. E. Foster, A. L. Choo, and S. A. Smith. 2023. [Speechlanguage disorder severity, academic success, and socioemotional functioning among multilingual and english children](#). *Frontiers in Psychology*, 14:1096145.
- Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh, Nathan Lile, and Noah D. Goodman. 2025. [Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars](#). *Preprint*, arXiv:2503.01307.
- D. Gibbard, L. Coglan, and J. MacDonald. 2004. [Costeffectiveness analysis of a preschool speech and language therapy service](#). *International Journal of Language & Communication Disorders*, 39(1):1–11.
- John Heilmann and 1 others. 2023. [Automation of language sample analysis](#). *Journal of Speech, Language, and Hearing Research*.

- E. R. Hitchcock, D. Harel, and T. M. Byun. 2015. [Social, emotional, and academic impact of residual speech errors in school-age children](#). *Seminars in Speech and Language*, 36(4):283–294.
- IBM Granite Team. 2025. Granite speech model documentation. <https://www.ibm.com/granite/docs/models/speech/>. Version last updated 14 May 2025; accessed 15 May 2025.
- Alon Jacovi, Avi Caciularu, Jonathan Herzig, Roei Aharoni, Bernd Bohnet, and Mor Geva. 2023. [A comprehensive evaluation of tool-assisted generation strategies](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13856–13878, Singapore. Association for Computational Linguistics.
- Wataru Kawakami, Keita Suzuki, and Junichiro Iwasawa. 2025. [Stabilizing reasoning in medical llms with continued pretraining and reasoning preference optimization](#). *Preprint*, arXiv:2504.18080.
- Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. 2023. [Continual pre-training of language models](#). *Preprint*, arXiv:2302.03241.
- Ondej Klempí and Radim Krupika. 2024. [Analyzing wav2vec 1.0 embeddings for cross-database parkinsons disease detection and speech features extraction](#). *Sensors*, 24(17):5520. Shows wav2vec embeddings generalize across Italian and English PD datasets.
- Jessica M. Lammert, Angela C. Roberts, Ken McRae, Laura J. Batterink, and Blake E. Butler. 2025. [Early identification of language disorders using natural language processing and machine learning: Challenges and emerging approaches](#). *Journal of Speech, Language, and Hearing Research*, 68(2):705–718.
- Marie-Thérèse Le Normand. 1997. Early morphological development in french children. In A. S. Olofsson and Sven Strömquist, editors, *Cross-Linguistic Studies of Dyslexia and Early Language Development*, pages 59–79. Office for Official Publications of the European Communities, Luxembourg.
- Ryan Liu, Jiayi Geng, Addison J. Wu, Iliia Sucholutsky, Tania Lombrozo, and Thomas L. Griffiths. 2024. [Mind your step \(by step\): Chain-of-thought can reduce performance on tasks where thinking makes humans worse](#). *arXiv preprint arXiv:2410.21333*.
- Saturnino Luz, Fasih Haider, Sofia de la Fuente, Davida Fromm, and Brian MacWhinney. 2020. [Alzheimers Dementia Recognition through Spontaneous Speech: The ADReSS Challenge](#). In *Proc. INTERSPEECH*.
- Umer Maqsood and 1 others. 2024. [Large language models for dysfluency detection in stuttered speech](#). *arXiv preprint arXiv:2406.11025*.
- Chirag Nagpal, Subhashini Venugopalan, Jimmy Tobin, Marilyn Ladewig, Katherine Heller, and Katriin Tomanek. 2025. [Speech recognition with llms adapted to disordered speech using reinforcement learning](#). In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.
- Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar, and Jason Weston. 2024. [Iterative reasoning preference optimization](#). *Preprint*, arXiv:2404.19733.
- Leena G Pillai, Kavya Manohar, Basil K Raju, and Elizabeth Sherly. 2024. [Multistage fine-tuning strategies for automatic speech recognition in low-resource languages](#). *Preprint*, arXiv:2411.04573.
- Giacomo Sanguedolce and 1 others. 2024. [Universal speech disorder recognition: Towards a foundation model for cross-pathology generalisation](#). <https://aphasia.talkbank.org/publications/2024/Sanguedolce24b.pdf>. Accessed May 2025.
- Phyllis Schneider, Denyse Hayward, and Rebecca V. Dubé. 2006. [Storytelling from pictures using the edmonton narrative norms instrument](#). *Journal of Speech-Language Pathology and Audiology*, 30:224–238.
- Alyssa Schwartz and 1 others. 2024. [Perceptions of artificial intelligence and chatgpt by speech-language pathologists and students](#). *American Journal of Speech-Language Pathology*.
- Joel Shor, Dotan Emanuel, Oran Lang, Omry Tuval, Michael Brenner, Julie Cattiau, Fernando Vieira, Maeve McNally, Taylor Charbonneau, Melissa Nollstadt, Avinatan Hassidim, and Yossi Matias. 2019. [Personalizing asr for dysarthric and accented speech with limited data](#). In *Interspeech 2019*, pages 784–788.
- Patti Simon and Sara Rosenbaum. 2016. *Speech and language disorders in children: Implications for the social security administration’s supplemental security income program*. National Academies Press.
- Hyewon Suh, Aayushi Dangol, Hedda Meadan, Carol A. Miller, and Julie A. Kientz. 2024. [Opportunities and challenges for ai-based support for speech-language pathologists](#). In *Proceedings of the 3rd Annual Meeting of the Symposium on Human-Computer Interaction for Work, CHIWORK ’24*, New York, NY, USA. Association for Computing Machinery.
- Jane Tsay. 2007. [Construction and automatization of an Minnan child speech corpus with some research findings](#). *Computational Linguistics and Chinese Language Processing*, 12(4):411–442.
- Dawn A. Tucker and Stella A. McKinnon. 2020. [Addressing the shortage of speech-language pathologists in rural areas: Barriers and solutions](#). *Journal of Rural Health*, 36(4):620–628.
- U.S. National Institute on Deafness and Other Communication Disorders. 2025. [Quick statistics about voice, speech, language](#).

- Yvette Wren, Linda L. Miller, Tim J. Peters, Ashley Emond, and Susan Roulstone. 2016. [Prevalence and predictors of persistent speech sound disorder at eight years old: Findings from a population cohort study](#). *J Speech Lang Hear Res*, 59(4):647–73.
- Qimin Yang, Rongsheng Wang, Jiexin Chen, Runqi Su, and Tao Tan. 2024. [Fine-tuning medical language models for enhanced long-contextual understanding and domain expertise](#). *Preprint*, arXiv:2407.11536.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. [Llamafactory: Unified efficient fine-tuning of 100+ language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.



## A Dataset Statistic

Table 1: Dataset statistics

Dataset	# Children	# Samples	Age Range
Ultrasuite	66	8338	5–13
ENNI	377	16546	4–9
LeNormand (French)	17	329	3–8
PERCEPT-GFTA	350	3664	6–17

## B Additional Related Works

**AI in Speech Language Pathology Assessment** The use of LLMs in clinical speech-language assessment has gained increasing attention in recent years. Several recent studies have demonstrated the utility of LLMs in detecting and characterizing speech and language disorders. For instance, [Bhattacharya et al. \(2024\)](#) showed that pre-trained LLMs could effectively identify both the presence and type of aphasia, suggesting that these models can serve as viable tools for clinical screening and diagnosis of language disorders. Beyond perception studies, a growing body of technical literature examines the use of speech and language features for automated assessment. [Engelhardt et al. \(2021\)](#) reviewed computational features used to assess cognitive and thought disorders, highlighting the relevance of acoustic and linguistic cues in differential diagnosis. Similarly, [Heilmann et al. \(2023\)](#) demonstrated that automatic language sample analysis tools can support clinical workflows, providing reliable linguistic metrics with reduced human effort.

**LLMs for Disordered Speech Analysis** A recent survey of SLPs and graduate students revealed a combination of cautious optimism and skepticism regarding the integration of LLMs such as ChatGPT into diagnostic and therapeutic workflows ([Schwartz et al., 2024](#)). These practitioner attitudes highlight critical socio-technical barriers to the clinical adoption of AI-driven systems in speech-language pathology. Recent research has explored the adaptation of LLMs for disordered speech processing. [Sanguedolce et al. \(2024\)](#) proposed a more generalized framework by fine-tuning Whisper on a dataset of stroke patients, resulting in a universal disordered-speech detection model. Their approach exhibited strong generalization across multiple neurological conditions, underscoring the potential of foundation models for broad-spectrum clinical speech applications.

## C Additional Analysis

**The Effect of Fewshot Examples** The results of the GPT-4 family across the first three scenarios under few-shot prompting indicate that few-shot examples do not consistently enhance the model’s intrinsic capabilities; the benefits of prompting are not uniformly evident. For instance, while few-shot prompting significantly improves the performance of GPT-4o-Mini-Transcribe and GPT-4o-Transcribe in Scenario 1, it leads to a reduction in accuracy in Scenario 2. Given that the examples provided to the model describe the expected transcription, few-shot prompting may bias the model to transcribe what the child was expected to say rather than accurately transcribe what the child is saying, which would explain the reduced diagnostic accuracy in Scenario 2. Our observations are consistent with prior findings on text-only LLMs ([Jacovi et al., 2023](#)).

## D Fine-tuning details

### D.1 Single Finetuning Round

We perform supervised fine-tuning on three models, including Qwen2-Audio 7B, Qwen2.5-Omni 3B, and Qwen2.5-Omni 7B using LLaMA-Factory framework ([Zheng et al., 2024](#)). We reuse the same UltraSuite dataset, but filtered out children whose speech in the evaluation set to avoid data contamination. We set the same fine-tuning hyperparameters across all models. Specifically, we used a LoRA rank of 32 with a LoRA alpha of 64, applying LoRA to all linear layers. The maximum token length was 4096, and training

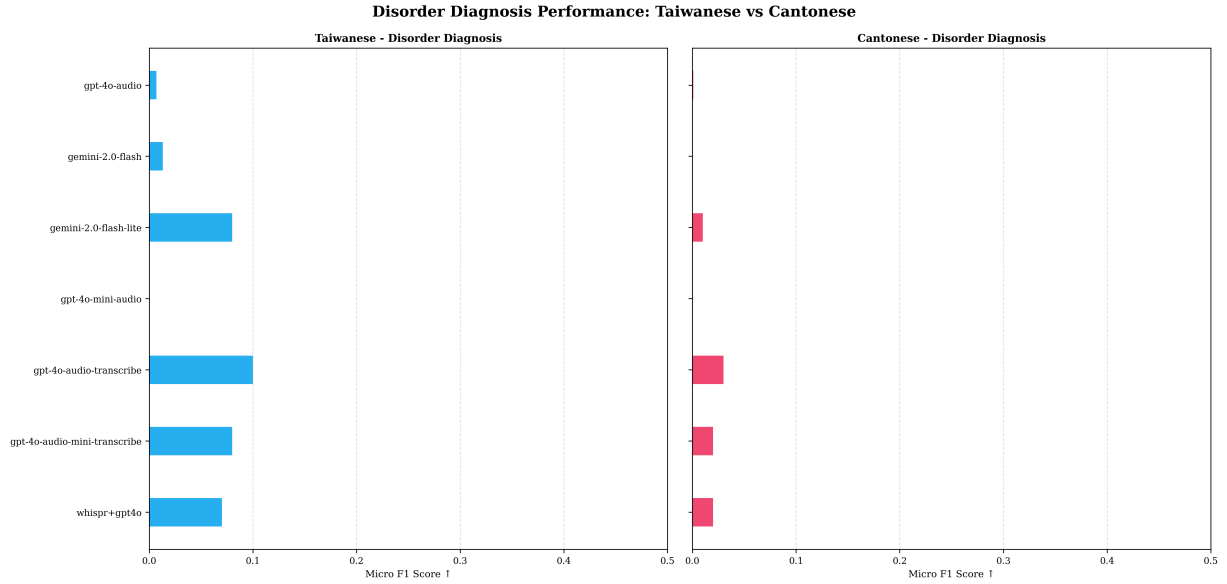


Figure 8: Disorder Diagnosis performance on tonal languages

was performed with a batch size of 32 for three epochs. We adopted a learning rate of 0.0001 with cosine scheduling and a warm-up ratio of 0.1. Regarding the prompts used for the three ablation settings for fine-tuning models, we present them as follows.

1. ASR-only without asterisk: `<audio>Transcribe this sound into text.`
2. ASR-only with asterisk: `<audio>Transcribe this sound into text. If the speech is disordered, please mark the words with an asterisk.`

## D.2 Multiple Finetuning Rounds

To investigate the models capacity for self-evolution, we conducted an iterative finetuning experiment in which the training data were generated by the model itself. The procedure is summarized as follows. In each iteration, we prompted the current model to generate up to 2,000 samples comprising both normal and disordered textspeech pairs. The details of the data generation process are described below. The generated data were then used to finetune the same model, after which the updated model was evaluated using our proposed framework. This procedure was repeated for five iterations. We use the Qwen 2.5 Omni 3B model for this experiment, and the finetuning hyperparameters were kept identical to those used in the single-round finetuning experiment described above.

**Data generation process.** We employed the ViLLM framework to deploy the current model and prompt it to produce both disordered and normal text and corresponding speech. We use simple sentences from the Tatoeba dataset<sup>2</sup> as the input text. The prompt used to generate text and speech is as below.

- **Normal text and speech:** Repeat exactly the following sentence and do not alter it. Return only the sentence and nothing else. The sentence is: `{transcript}`
- **Disordered text and speech:** Assuming you are a child with speech disordered. `{disorder_description}` Generate a corresponding transcript with the `'{disorder_type} disorder'` when the child say: `"{transcript}"`. Return only the transcript text and nothing else.

<sup>2</sup>[https://huggingface.co/datasets/Helsinki-NLP/tatoeba\\_mt](https://huggingface.co/datasets/Helsinki-NLP/tatoeba_mt)

If the model failed to generate disordered text, we manually introduced disorder markers into the text (including addition, substitution, omission, and stuttering). We then instructed the model to regenerate the disordered text and corresponding speech using the prompt as above. All generated audio was further processed through a *child-en* pipeline, which applied a series of signal processing operations (including pitch shifting, low-frequency reduction, high-frequency boosting, and slight speed-up) to approximate the characteristics of childrens speech. Subsequently, each sample was verified for detecting disorder using the same current model. If verification failed, generation was repeated up to five times; if it still failed, the sample was discarded, and the process moved on to the next sample. The prompt we used for verification is presented as follows.

```
You are a highly experienced Speech-Language Pathologist (SLP). An audio recording of a child's voice will be
→ provided. The child is attempting to say: "{transcript}". Based on your professional expertise,
→ analyze the child's speech production in the recording to determine whether it indicates a potential
→ speech-language disorder (addition, substitution, omission, stuttering). Respond with one of the
→ following labels only: 'speech_disorder', 'typically_developing'. Do not include any explanation
→ or additional text in your response.
```

We present the results of the self-evolving experiment in Table 2. As shown, the transcription task exhibits a clear reduction in word error rate (WER), which can be attributed to the availability of the ground-truth transcript for each generated speech sample. In contrast, the performance on disorder diagnosis remains largely unchanged across iterations. This is expected, as the same model is used to verify whether the speech is disordered, thereby providing no additional learning signal. Consequently, the model's performance on this task does not surpass that of the original model. A promising direction for future work is to employ a stronger verification model and/or incorporate supervision on the disorder type, which may yield richer signals for iterative finetuning.

## E Detailed Results

2 presents our evaluation results. Values are reported as mean  $\pm$  bootstrap standard deviation, providing both the point estimate and a measure of uncertainty in the metric estimates.

## F Prompts

### Scenario 1: Binary Classification

```
You are a highly experienced Speech-Language Pathologist (SLP). An audio recording will be provided,
→ typically consisting of a speech prompt from a pathologist followed by a child's repetition. The
→ prompt the child is trying to repeat is as follows: {words}. Based on your professional expertise:
→ 1. Assess the child's speech in the recording for signs of typical development or potential speech-
→ language disorder. 2. Conclude your analysis with one of the following labels only: typically
→ developing or speech disorder. 3. Provide your response as a single letter without any additional
→ explanation, commentary, or unnecessary text.
```

### Scenario 2: ASR-Based Classification

```
You are a highly experienced Speech-Language Pathologist (SLP). An audio recording is provided to you,
→ typically consisting of a speech prompt from a pathologist followed by a child's repetition. Based
→ on your expertise transcribe the child's speech into text. Do not make any assumptions about the
→ words the child is expected to say. Only transcribe based on the words that the child actually says.
→ Only respond with the text transcription, no other text or commentary.
```

### Scenario 3: Transcription Accuracy

```
You are a highly experienced Speech-Language Pathologist (SLP). An audio recording will be provided,
→ typically consisting of a speech prompt from a pathologist followed by a child's repetition. Based
→ on your expertise transcribe the child's speech into text. Try to understand what the child is
→ expected to say. And only respond with the transcription of the child's speech. Not the pathologist's
→ prompt or any other commentary. Only respond with the text transcription, no other text,
→ commentary or punctuations.
```

### Scenario 4: Disorder Type Classification

```
You are a highly experienced Speech-Language Pathologist (SLP). An audio recording will be provided,
→ typically consisting of a speech prompt from a pathologist followed by a child's repetition. The
→ prompt text the child is trying to repeat is as follows: {words}. Based on your professional
→ expertise: 1. Assess the child's speech in the recording for signs of typical development or
→ potential speech-language disorder. 2. Conclude your analysis with one of the following labels only:
→ A - 'typically developing' (child's speech patterns and development are within normal age-
→ appropriate ranges), B - 'articulation' (difficulty producing specific speech sounds correctly, such
→ as substituting, omitting, or distorting sounds), C - 'phonological' (difficulty understanding and
→ using the sound system of language, affecting sounds of a particular type). 3. Provide your response
→ as a single letter without any additional explanation, commentary, or unnecessary text
```

Table 2: Detailed Model Performance Results

Model	Micro F1 ↑				Macro F1 ↑				Exact Match ↑				WER ↓			
	Disorder Diag.	ASR Diag.	Disorder Type	Symptom Diag.	Disorder Diag.	ASR Diag.	Disorder Type	Symptom Diag.	Disorder Diag.	ASR Diag.	Disorder Type	Symptom Diag.	WER ↓	MER ↓	WIP ↓	
gemini-2.0-flash-lite	0.457 ±0.015	0.459 ±0.011	0.180 ±0.012	0.415 ±0.016	0.341 ±0.010	0.316 ±0.012	0.174 ±0.011	0.209 ±0.024	0.457 ±0.015	0.459 ±0.011	0.180 ±0.012	0.415 ±0.016	0.51 ±0.016	0.51 ±0.016	0.49 ±0.016	
gemini-2.0-flash	0.458 ±0.014	0.460 ±0.011	0.204 ±0.015	0.180 ±0.015	0.330 ±0.009	0.315 ±0.011	0.202 ±0.011	0.072 ±0.005	0.458 ±0.014	0.460 ±0.011	0.204 ±0.015	0.180 ±0.015	1.00 ±0.000	1.00 ±0.000	1.00 ±0.000	
gpt-4o-mini-audio	0.204 ±0.011	0.438 ±0.009	0.125 ±0.011	0.433 ±0.017	0.208 ±0.009	0.371 ±0.008	0.114 ±0.008	0.112 ±0.005	0.204 ±0.011	0.438 ±0.009	0.125 ±0.011	0.433 ±0.017	5.64 ±0.006	0.96 ±0.006	0.94 ±0.006	
gpt-4o-audio	0.439 ±0.014	0.459 ±0.014	0.208 ±0.015	<b>0.507 ±0.016</b>	0.238 ±0.005	0.316 ±0.013	0.166 ±0.007	0.189 ±0.009	0.439 ±0.014	0.459 ±0.014	0.208 ±0.015	<b>0.507 ±0.016</b>	5.40 ±0.177	0.68 ±0.015	0.32 ±0.015	
gpt-4o-mini-transcribe	0.456 ±0.018	0.460 ±0.012	0.358 ±0.016	0.310 ±0.016	0.360 ±0.012	0.310 ±0.014	0.274 ±0.015	0.209 ±0.014	0.456 ±0.018	0.460 ±0.012	0.358 ±0.016	0.310 ±0.016	0.37 ±0.015	0.37 ±0.015	0.63 ±0.015	
gpt-4o-transcribe	0.474 ±0.017	0.460 ±0.011	0.382 ±0.016	0.325 ±0.015	0.373 ±0.014	0.320 ±0.010	<b>0.520 ±0.017</b>	0.181 ±0.019	0.474 ±0.017	0.460 ±0.011	0.382 ±0.016	0.325 ±0.015	1.00 ±0.000	1.00 ±0.000	1.00 ±0.000	
whisper-gpt4o	0.473 ±0.015	0.461 ±0.015	<b>0.399 ±0.016</b>	0.395 ±0.016	0.369 ±0.013	0.310 ±0.014	0.245 ±0.013	<b>0.563 ±0.015</b>	0.473 ±0.015	0.461 ±0.015	<b>0.399 ±0.016</b>	0.395 ±0.016	<b>0.33 ±0.018</b>	<b>0.29 ±0.015</b>	<b>0.71 ±0.015</b>	
gpt-4o-mini-audio	0.455 ±0.006	0.713 ±0.008	0.345 ±0.005	0.163 ±0.005	<b>0.425 ±0.007</b>	0.4675 ±0.007	0.202 ±0.007	0.105 ±0.007	0.455 ±0.006	0.713 ±0.008	0.345 ±0.005	0.163 ±0.005	2.17 ±0.009	0.74 ±0.009	0.22 ±0.009	
qwen2.5-omni-7b	<b>0.556 ±0.005</b>	0.738 ±0.007	0.306 ±0.016	0.155 ±0.016	0.401 ±0.008	0.4283 ±0.006	0.273 ±0.015	0.097 ±0.016	<b>0.556 ±0.005</b>	0.738 ±0.007	0.306 ±0.016	0.155 ±0.016	4.90 ±0.150	0.79 ±0.014	0.19 ±0.023	
qwen2.5-omni-7b-ft	0.449 ±0.009	<b>0.747 ±0.004</b>	0.300 ±0.009	0.090 ±0.009	0.300 ±0.009	0.4276 ±0.008	0.090 ±0.009	0.090 ±0.009	0.449 ±0.009	<b>0.747 ±0.004</b>	0.300 ±0.009	0.090 ±0.009	0.57 ±0.068	0.58 ±0.062	0.40 ±0.062	
qwen2.5-omni-3b	0.447 ±0.005	0.743 ±0.006	0.335 ±0.005	0.097 ±0.003	0.206 ±0.002	<b>0.4803 ±0.005</b>	0.203 ±0.003	0.083 ±0.008	0.447 ±0.005	0.743 ±0.006	0.335 ±0.005	0.097 ±0.003	0.86 ±0.007	0.46 ±0.001	0.52 ±0.002	
phi-multimodal	0.325 ±0.005	0.746 ±0.003	0.000 ±0.000	0.000 ±0.000	0.000 ±0.000	0.4526 ±0.004	0.000 ±0.000	0.000 ±0.000	0.325 ±0.005	0.746 ±0.003	0.000 ±0.000	0.000 ±0.000	9.11 ±0.117	0.58 ±0.015	0.41 ±0.005	
granite-speech-3.3-8b	0.000 ±0.000	0.739 ±0.005	0.000 ±0.000	0.000 ±0.000	0.000 ±0.000	0.4564 ±0.003	0.000 ±0.000	0.000 ±0.000	0.000 ±0.000	0.739 ±0.005	0.000 ±0.000	0.000 ±0.000	2.34 ±0.018	0.71 ±0.002	0.25 ±0.002	
gemini-3.0-flash-32b	0.000 ±0.000	0.721 ±0.007	0.007 ±0.001	0.200 ±0.005	0.000 ±0.000	0.4797 ±0.006	0.007 ±0.001	0.096 ±0.002	0.000 ±0.000	0.721 ±0.007	0.007 ±0.001	0.200 ±0.005	2.25 ±0.017	0.46 ±0.004	0.48 ±0.003	
Finetuned Models with Asterisk																
qwen2.5-omni-7b-ft	<b>0.971 ±0.003</b>	0.611 ±0.020	<b>0.971 ±0.003</b>	<b>0.971 ±0.003</b>	<b>0.951 ±0.003</b>	0.511 ±0.010	<b>0.951 ±0.002</b>	<b>0.951 ±0.002</b>	<b>0.951 ±0.003</b>	0.631 ±0.020	<b>0.951 ±0.002</b>	<b>0.951 ±0.002</b>	1.40 ±0.011	0.52 ±0.006	0.41 ±0.009	
qwen2.5-omni-3b-ft	0.912 ±0.005	0.626 ±0.015	0.911 ±0.005	0.912 ±0.005	0.890 ±0.003	0.497 ±0.010	0.890 ±0.003	0.890 ±0.003	0.912 ±0.005	0.626 ±0.015	0.890 ±0.003	0.890 ±0.003	0.97 ±0.014	0.53 ±0.005	0.39 ±0.007	
qwen2.5-audio-instruct-ft	0.314 ±0.006	<b>0.609 ±0.018</b>	0.132 ±0.002	0.065 ±0.002	0.015 ±0.014	<b>0.296 ±0.005</b>	0.214 ±0.004	0.073 ±0.003	0.314 ±0.006	<b>0.609 ±0.018</b>	0.132 ±0.004	0.073 ±0.003	<b>0.58 ±0.017</b>	<b>0.43 ±0.004</b>	<b>0.50 ±0.008</b>	
Performance Models without Asterisk																
qwen2.5-omni-7b-ft	<b>0.67 ± 0.008</b>	<b>0.60 ± 0.012</b>	<b>0.40 ± 0.011</b>	<b>0.34 ± 0.007</b>	<b>0.81 ± 0.013</b>	<b>0.49 ± 0.015</b>	<b>0.28 ± 0.016</b>	<b>0.16 ± 0.016</b>	<b>0.67 ± 0.013</b>	<b>0.60 ± 0.012</b>	<b>0.40 ± 0.016</b>	<b>0.34 ± 0.016</b>	1.76 ± 0.019	0.46 ± 0.005	0.47 ± 0.008	
qwen2.5-omni-3b-ft	0.54 ± 0.009	0.59 ± 0.011	0.36 ± 0.010	0.16 ± 0.008	0.35 ± 0.013	0.46 ± 0.014	0.24 ± 0.016	0.08 ± 0.013	0.54 ± 0.013	0.59 ± 0.011	0.36 ± 0.016	0.16 ± 0.013	0.95 ± 0.014	0.49 ± 0.006	0.43 ± 0.007	
qwen2.5-audio-instruct-ft	0.36 ± 0.005	0.55 ± 0.004	0.27 ± 0.008	0.08 ± 0.007	0.01 ± 0.002	<b>0.49 ± 0.006</b>	0.05 ± 0.015	0.01 ± 0.009	0.36 ± 0.005	0.55 ± 0.004	0.27 ± 0.015	0.08 ± 0.009	<b>0.52 ± 0.015</b>	<b>0.38 ± 0.004</b>	<b>0.56 ± 0.006</b>	
Freshnet Prompting																
gpt-4o-mini-audio	0.23	0.09	-	-	0.19	0.01	-	-	0.23	0.09	-	-	1.58	0.65	0.28	
gpt-4o-audio	0.71	<b>0.12</b>	-	-	0.21	<b>0.04</b>	-	-	0.71	<b>0.12</b>	-	-	1.73	<b>0.62</b>	<b>0.30</b>	
gpt-4o-mini-transcribe	<b>0.72</b>	0.01	-	-	0.44	0.01	-	-	<b>0.72</b>	0.01	-	-	1.08	0.80	0.12	
gpt-4o-transcribe	<b>0.72</b>	0.01	-	-	<b>0.45</b>	0.01	-	-	<b>0.72</b>	0.01	-	-	<b>1.01</b>	0.79	0.14	
whisper-gpt4o	<b>0.72</b>	0.01	-	-	0.24	0.00	-	-	<b>0.72</b>	0.01	-	-	1.89	0.77	0.16	
Performance across multiple finetuning iterations																
Iteration 0 (baseline)	<b>0.555 ± 0.027</b>	0.738 ± 0.002	0.390 ± 0.029	0.150 ± 0.033	0.401 ± 0.021	0.428 ± 0.002	0.264 ± 0.019	0.094 ± 0.021	<b>0.555 ± 0.026</b>	0.738 ± 0.003	0.390 ± 0.016	0.150 ± 0.028	4.90 ± 0.115	<b>0.787 ± 0.034</b>	0.185 ± 0.023	
Iteration 1	0.527 ± 0.032	0.740 ± 0.001	0.537 ± 0.031	0.123 ± 0.028	0.360 ± 0.020	0.425 ± 0.002	0.255 ± 0.020	0.106 ± 0.018	0.527 ± 0.017	0.740 ± 0.002	0.537 ± 0.034	0.123 ± 0.022	2.77 ± 0.23	0.865 ± 0.017	0.114 ± 0.025	
Iteration 2	0.529 ± 0.034	0.740 ± 0.003	0.505 ± 0.022	0.140 ± 0.021	<b>0.407 ± 0.033</b>	0.425 ± 0.003	0.265 ± 0.024	0.122 ± 0.030	0.529 ± 0.031	0.740 ± 0.001	0.505 ± 0.032	0.140 ± 0.016	2.63 ± 0.11	0.851 ± 0.028	0.126 ± 0.019	
Iteration 3	0.526 ± 0.024	0.740 ± 0.003	0.552 ± 0.018	0.125 ± 0.027	0.379 ± 0.026	0.425 ± 0.002	0.239 ± 0.028	0.080 ± 0.030	0.526 ± 0.025	0.740 ± 0.002	0.552 ± 0.033	0.125 ± 0.029	2.63 ± 0.17	0.854 ± 0.029	0.121 ± 0.017	
Iteration 4	0.528 ± 0.038	0.740 ± 0.003	0.553 ± 0.028	0.131 ± 0.024	0.378 ± 0.028	0.425 ± 0.002	0.254 ± 0.021	0.085 ± 0.019	0.528 ± 0.032	0.740 ± 0.001	0.553 ± 0.027	0.131 ± 0.021	2.63 ± 0.17	0.854 ± 0.029	0.121 ± 0.017	
Iteration 5	0.531 ± 0.026	0.740 ± 0.003	0.554 ± 0.017	0.140 ± 0.024	0.385 ± 0.033	0.425 ± 0.003	0.250 ± 0.018	0.097 ± 0.029	0.531 ± 0.019	0.740 ± 0.002	0.554 ± 0.027	0.140 ± 0.029	<b>2.14 ± 0.14</b>	0.846 ± 0.023	0.129 ± 0.023	
Model Ensemble Performance																
Google + OpenAI Ensemble	0.455	-	0.313	0.389	0.319	-	0.215	0.184	0.455	-	0.313	0.389	0.96	0.57	0.03	
Gemini 2.0 Flash, G. Gemini 2.0 Flash Lite, OpenAI GPT-4o	0.455	-	0.313	0.389	0.319	-	0.215	0.184	0.455	-	0.313	0.389	0.92	0.54	0.07	
Qwen + OpenAI Variant	0.455	-	0.313	0.389	0.319	-	0.215	0.184	0.455	-	0.313	0.389	0.92	0.54	0.07	
Qwen + OpenAI Variant	0.455	-	0.313	0.389	0.319	-	0.215	0.184	0.455	-	0.313	0.389	0.92	0.54	0.07	
Qwen + OpenAI Hybrid	0.403	-	0.338	0.167	0.300	-	0.189	0.109	0.403	-	0.338	0.167	1.88	1.75	0.3	
Qwen + OpenAI Hybrid	0.403	-	0.338	0.167	0.300	-	0.189	0.109	0.403	-	0.338	0.167	1.88	1.75	0.3	
Performance by Age																
Gemini 2.0 Flash	0.7629 ±0.0291	0.7700 ±0.0064	0.475 ±0.037	0.1329 ±0.012	0.4326 ±0.0094	0.4350 ±0.0347	0.239 ±0.026	0.0557 ±0.008	0.7629 ±0.0291	0.7700 ±0.0064	0.475 ±0.037	0.1329 ±0.012	4.346 ±0.115	1.000 ±0.000	0.000 ±0.000	
pre-elementary	0.4623 ±0.023	0.4500 ±0.0099	-	0.1836 ±0.014	0.303 ±0.013	0.3103 ±0.0184	0.156 ±0.009	0.0729 ±0.009	0.4623 ±0.023	0.4500 ±0.0099	-	0.1836 ±0.014	4.452 ±0.129	0.998 ±0.001	0.000 ±0.000	
elementary	0.280 ±0.026	0.2800 ±0.0281	0.200 ±0.024	0.1329 ±0.011	0.228 ±0.019	0.2188 ±0.0101	0.159 ±0.019	0.0557 ±0.007	0.280 ±0.026	0.2800 ±0.0281	0.200 ±0.024	0.1329 ±0.011	<b>2.960 ±0.008</b>	0.998 ±0.001	0.000 ±0.000	
Gemini 2.0 Flash Lite	0.7288 ±0.0302	0.7700 ±0.0067	0.306 ±0.028	0.4930 ±0.018	0.4214 ±0.0101	0.4350 ±0.0345	0.245 ±0.021	<b>0.2102 ±0.013</b>	0.7288 ±0.0302	0.7700 ±0.0067	0.306 ±0.028	0.4930 ±0.018	<b>4.290 ±0.091</b>	1.000 ±0.000	0.000 ±0.000	
pre-elementary	0.428 ±0.023	0.4450 ±0.0256	-	0.4361 ±0.019	0.322 ±0.015	0.3000 ±0.0427	-	<b>0.2419 ±0.015</b>	0.428 ±0.023	0.4450 ±0.0256	-	0.4361 ±0.019	<b>3.458 ±0.102</b>	0.990 ±0.002	0.000 ±0.000	
elementary	0.271 ±0.026	0.2800 ±0.0094	0.106 ±0.017	0.4930 ±0.017	0.222 ±0.018	0.2188 ±0.0118	0.119 ±0.017	<b>0.2102 ±0.012</b>	0.271 ±0.026	0.2800 ±0.0094	0.106 ±0.017	<b>0.4930 ±0.017</b>	2.650 ±0.076	0.992 ±0.002	0.000 ±0.000	
GPT-4o Audio	0.490 ±0.0294	0.7700 ±0.0154	0.561 ±0.034	<b>0.6049 ±0.021</b>	0.4281 ±0.0096	0.4350 ±0.0347	0.238 ±0.021	0.1648 ±0.011	0.490 ±0.0294	0.7700 ±0.0154	0.561 ±0.034	<b>0.6049 ±0.021</b>	3.650 ±0.155	1.000 ±0.000	0.000 ±0.000	
pre-elementary	0.414 ±0.023	0.4600 ±0.0223	-	<b>0.5115 ±0.020</b>	0.296 ±0.012	0.296 ±0.012	0.166 ±0.016	0.1206 ±0.013	0.414 ±0.023	0.4600 ±0.0223	-	<b>0.5115 ±0.020</b>	3.086 ±0.180	0.998 ±0.001	0.000 ±0.000	
elementary	0.271 ±0.026	0.2800 ±0.0260	0.261 ±0.025	<b>0.6049 ±0.019</b>	0.219 ±0.018	0.2188 ±0.0110	0.166 ±0.016	0.1648 ±0.010	0.271 ±0.026	0.2800 ±0.0229	0.261 ±0.025	0.6049 ±0.019	4.162 ±0.132	1.000 ±0.000	0.000 ±0.000	
GPT-4o Mini Audio	0.0984 ±0.0202	0.7900 ±0.0137	0.232 ±0.028	0.574 ±0.020	0.0868 ±0.0180	0.5200 ±0.0477	0.101 ±0.013	0.1226 ±0.010	0.0984 ±0.0202	0.7900 ±0.0137	0.232 ±0.028	0.574 ±0.020	4.162 ±0.132			



## Scenario 5: SymptomClassification

You are a highly experienced Speech-Language Pathologist (SLP). An audio recording will be provided, typically consisting of a speech prompt from a pathologist followed by a child's repetition. The target phrase the child is attempting to repeat is: \{words\}. Based on your professional expertise, assess the child's speech in the recording and identify any abnormal features. These features can be one of the following: A - 'substitution' (the child replaces one word, syllable, or sound with another), B - 'omission' (the child omits a word, syllable, or sound), C - 'addition' (the child adds an extra word, syllable, or sound), D - 'typically developing' (the child's speech is appropriate for their age), or E - 'stuttering' (the child exhibits repetition, prolongation, or difficulty initiating speech). Provide your response as a single letter (A-E) only, without any additional explanation or commentary.

## G Gemini 2.0 Flash Reasoning Traces

Below, we provide some examples of reasoning traces output by the Gemini 2.0 Flash model when tasked with disorder type diagnosis.

Identify the Target Sounds: The target words are "gorge," "gate," "get," and "luggage." The key sounds within these words include /g/, /t/, /e/, /l/. Listen for Substitutions/Distortions: The child says "gore" for "gorge", "gape" for "gate", "get" appears correct, and "nugget" for "luggage". Analyze the Error Patterns: "Gore" for "gorge" seems to be a mispronunciation or substitution of the vowel and omitting the final sound /dʒ/. "Gape" for "gate" also a mispronunciation or substitution of the vowel sound /e/. "Nugget" for "luggage" involves multiple sound changes. The initial /l/ is replaced with /n/, the vowel /ʌ/ is present, and the final sound /t/ is replaced by /t/. Determine the Type of Speech Error: The sound changes observed are primarily substitutions and omissions, indicating difficulties with producing specific sounds accurately. Classify the Speech Pattern: Based on the analysis, the most appropriate classification is articulation, as the errors relate to the physical production of individual speech sounds.

Okay, let's analyze the child's speech based on the prompt "owrow awraw ooroo". The prompt presents a challenge, focusing on vowel sounds within a simple consonant-vowel-consonant (CVC) structure. The child's repetitions are "O wow", "A wa", and "Uu Uu". \* \*\*owrow\*\* becomes "O wow": This indicates a possible difficulty in producing the /r/ sound in the middle of the word. The vowel sounds are relatively close but the /r/ is replaced by a /w/. \* \*\*awraw\*\* becomes "A wa": Similar to the first attempt, the child simplifies or substitutes the /r/ sound with a /w/. The vowel sounds are again reasonably close. \* \*\*ooroo\*\* becomes "Uu Uu": This repetition is more successful, with the child accurately producing the vowel sound, but there is no attempt to use the /r/. This suggests that the child is either not hearing the /r/ sound, or they are hearing it but unable to produce it. The consistent difficulty with the /r/ sound suggests an articulation issue specifically related to that sound. The vowel sounds seem relatively well approximated, pointing away from significant phonological problems where a whole sound class (e.g., all back vowels) would be affected. Conclusion: B - 'articulation'

Okay, let's analyze the child's speech based on the recording. The prompt was "lie reek hair strewn." The child said: "lie reek hair shm." \* \*\*lie:\*\* Correct production. \* \*\*reek:\*\* Correct production. \* \*\*hair:\*\* Correct production. \* \*\*strewn:\*\* "shm". This indicates a cluster reduction, where the /str/ consonant cluster is reduced to /shm/. The error on "strewn" demonstrates difficulty with consonant clusters, specifically /str/. There's no evidence of distortions, substitutions, or omissions on single phonemes in the other words. The error is on a later-developing sound. Given that the child is having difficulty with consonant clusters it suggests that a phonological process is at play. Conclusion: C - 'phonological'

Okay, let's analyze this speech sample. The child attempts the phrase "coot packer toot tie". \* \*\*coot:\*\* The child says "koot" which sounds accurate. \* \*\*packer:\*\* The child says "packa". The "r" sound is missing or distorted. \* \*\*toot:\*\* The child says "toot" which sounds accurate. \* \*\*tie:\*\* The child says "tie" which sounds accurate. The primary error observed is with the "r" sound in "packer." It's possible that the "r" is simply omitted, or it may be a distorted production. The other sounds are articulated relatively clearly. This difficulty with the "r" sound suggests an issue with producing a specific sound, rather than a broader issue with the phonological system. Conclusion: B - 'articulation'

Okay, let's analyze the child's speech. The prompt is: "You get milk from cows." The child's response is: "You get milk from dows." The child substitutes the /k/ sound in "cows" with a /d/ sound, resulting in "dows." This is a substitution error affecting a specific sound. There's no evidence of errors impacting an entire sound class (e.g., all fricatives, all velar sounds). The error involves difficulty with the articulation of the /k/ sound. Therefore, the most appropriate label is 'articulation'. Reasoning: \* The child's speech contains a sound substitution. \* The substitution affects one specific sound: /k/ becoming /d/. \* The errors involve specific sounds, not an entire sound class (which would indicate a phonological issue). Answer: B

Okay, let's analyze the child's speech based on the provided audio.

The prompt was: "kai meek poke cow" The child responded with: "tai, meet, poat, dow"

Here's a breakdown:

**\*\*kai -> tai:\*\*** The child substitutes /k/ with /t/. **\*\*meek -> meet:\*\*** The vowel sounds are similar; no concerning error. **\*\*poke -> poat:\*\*** The child substitutes /k/ with /t/. **\*\*cow -> dow:\*\*** The child substitutes /k/ with /d/.

The child consistently substitutes /k/ with /t/ and /k/ with /d/. This pattern suggests a difficulty with the /k/ sound. Since the errors are consistent across multiple instances of the /k/ sound, it points to a potential issue with the phonological system rather than just a simple articulation error.

**\*\*Conclusion:\*\***

C - 'phonological'

The child's speech contains several sound errors. "Crab" is produced as "gap," suggesting difficulty with the /kr/ consonant cluster. "Biscuits" is produced as "bisits," which contains a sound distortion of /k/. "Helicopter" is produced as "hedocter," involving multiple sound substitutions and omissions. These errors point to difficulties with the articulation and phonological representation of sounds.

ANSWER: B

Okay, let's analyze the child's repetition of "DDK ppp ttt kkk ptk patticake."

The prompt is designed to assess the child's ability to produce rapid, alternating sounds (diadochokinesis or DDK) and to articulate specific phonemes in isolation and in a word.

**\*\*DDK (ppp ttt kkk):\*\*** The child attempts the DDK task and does show signs of struggling. The child seems to understand the prompt, and can repeat the p, t and k sounds. However the sounds appear jumbled together. **\*\*ptk:\*\*** The child attempts ptk. The child appears to jumble this prompt. **\*\*patticake:\*\*** The child repeats this prompt several times. The child shows great difficulty producing this word and does not enunciate the word well.

Based on these observations, the child exhibits difficulties with articulation and sequencing sounds. While single sounds can be produced, transitioning between them appears challenging. This suggests the primary difficulty is in the motor production of the sounds.

**\*\*Conclusion:\*\*** B - 'articulation'

Okay, I've listened to the audio. Here's my analysis:

**\*\*Reasoning:\*\***

The target words contain the // phoneme (as in "sh"). The child says "shore," then what sounds like "tosh," then the child attempts to repeat "puss." I hear what sounds like a substitution of // with /s/ in the word "puss," making it sound similar to the target word of "puss." This consistent substitution of // with /s/ (or perhaps a distortion close to /s/) across multiple words containing the // phoneme suggests a potential phonological pattern, specifically affecting fricative sounds.

**\*\*Conclusion:\*\***

C - 'phonological'

Okay, let's analyze the child's speech based on the recording. - **\*\*Sock:\*\*** The child's repetition of "sock" sounds accurate. - **\*\*Shoop:\*\*** The child's repetition of "shoop" sounds accurate. - **\*\*Shire:\*\*** The child's repetition of "shire" sounds like "show." This suggests a potential difficulty with the /t/ sound after the /sh/ sound. - **\*\*Shallow:\*\*** The child's repetition of "shallow" sounds like "show." This suggests a potential difficulty with the /l/ sound after the /sh/ sound.

Based on these observations, the child demonstrates difficulty with the production of the /t/ and /l/ sounds in the context of the /sh/ blend. There is evidence of a potential speech sound disorder. Given the errors are consistent with particular sound classes, it is most indicative of a phonological pattern. Final Answer: The final answer is C

Okay, let's analyze the child's speech. The prompt "owrow awraw ooroo" contains several instances of the "r" sound in different syllable positions. The SLP in the recording notes that there was barely an "r" in the first attempt and it sounded like "all to ohs." The subsequent attempt on "owrow" was considered much better. The final attempt on "ooroo" appears to be a distortion, possibly an "r" sound replaced with a vowel or gliding sound. Based on the limited data, there is an apparent difficulty with the production of the /r/ sound. This could indicate an articulation issue, where the child struggles with the motor movements required for the /r/ sound. Since the difficulty is centered on one specific sound, articulation is the most appropriate label.

**\*\*Answer: B\*\***

Okay, I will analyze the provided audio and provide a diagnosis based on my expertise as an SLP.

Here's my analysis:

The child's repetition of "A dozen green ghosts" sounds generally clear. However, a closer listening reveals a potential difficulty with the 'th' sound, as in "ghosts". The repetition of "ghosts" is not as clear compared to the phrase, "dozen green".

Reasoning: 1. The child produces "dozen green" relatively accurately. 2. There is a possible sound substitution affecting the /s/ sound. 3. Due to the specific nature of the sound, and a mild misarticulation it suggests a more localized articulation difficulty rather than a broader phonological issue.

Therefore, my diagnosis is:

B - 'articulation'

Okay, let's analyze the child's speech based on the provided audio.

The target words are "book boy". The child says something that sounds like "buh... book... boy."

Here's a breakdown of what I'm hearing: \* The child seems to produce the "b" sound correctly at the beginning of both words. \* For "book," they seem to be attempting the vowel sound and the "k" sound. It doesn't sound perfect, but it could be due to typical developmental imprecision, or it could indicate a slight articulation issue with the "k" sound. \* For "boy," the "b" sounds correct, and the vowel sound sounds within an acceptable range.

Given the information and the presence of the "b" at the beginning of both words, the rest of the sounds aren't clear enough to indicate a Phonological delay. However, there are some slight errors in production.

Therefore, my assessment is:

ANSWER: B - 'articulation'

Okay, I will listen to the audio and analyze the child's speech based on the prompt "sheen shake sham shock."

**\*\*Analysis:\*\***

The child appears to have difficulty with the initial // sound in all four words. While the vowel sounds seem relatively intact, the consistency of the // distortion suggests a pattern rather than an isolated error. It is difficult to determine the exact nature of the distortion without a more thorough assessment (e.g., is it a substitution, omission, or distortion?). The consistency across multiple words with the same sound indicates a possible phonological issue.

**\*\*Conclusion:\*\***

C - 'phonological'

Okay, let's analyze the child's speech sample.

The prompt words are: "die," "got," "hag," and "Hague".

Here's a breakdown of potential issues:

**\*\*\*"die"\*\*\*:** The child repeats "die" which sounds accurate. **\*\*\*"got"\*\*\*:** The child repeats "got" which sounds accurate. **\*\*\*"hag"\*\*\*:** The child repeats "hag" which sounds accurate. **\*\*\*"Hague"\*\*\*:** The child repeats "Hake" which involves changing the final sound.

The one slight discrepancy could be related to the place of articulation. Because this appears to be a single instance, I will not classify this as a phonological issue. I'm confident that this instance points to a lack of production of the "g" sound in "Hague."

**\*\*Conclusion:\*\***

B - 'articulation'

## H Computational Resources

To ensure reproducibility and transparency, we detail the computational resources required for running the SLPHelm framework as follows. Closed-source models were accessed via API, with an average inference time of approximately 30 minutes per model. In contrast, open-source models required significantly more time, averaging 3–4 hours per model on high-performance hardware. All open-source experiments were conducted using an NVIDIA H200 GPU with 140GB of VRAM, enabling efficient handling of large-scale model inference and evaluation.

## I AI Tool Usage Declaration

We used ChatGPT to format paragraphs and correct grammatical mistakes. We are responsible for the content and quality of the study.