

Contextualizing Parental Behaviors in Bilingual Datasets from In-Person and Telehealth Language Assessment

Anonymous ACL submission

Abstract

The increasing adoption of telehealth technologies presents both opportunities and challenges, offering greater convenience for patients while increasing clinicians workload, particularly in managing remotely collected data. Bilingual speech-language pathologists (SLPs) spent substantial effort in evaluating parent behaviors when conducting family-centered language assessments. In this study we collaborate with SLPs to examine how large language models (LLMs) can support clinical workflows and address real-world challenges in telehealth. We collected a detailed bilingual dataset of 59 Mandarin-English child language assessment sessions (16 in-person and 43 via telehealth) and benchmarked three open-source LLMs and one closed-source LLM on this task. All four LLMs are still inferior to human experts despite notable accuracy, and additional error analysis revealed that domain complexity, cultural context, and multimodal cues pose significant challenges for both LLMs and human annotators. This work highlights the need for domain-specific NLP advancement and evaluation methods that extend beyond standard benchmarks to include clinical utility, workflow integration, and cultural appropriateness in bilingual telehealth assessment.

1 Introduction

Language assessments play a critical role in the early detection and intervention of children’s communication disorders, particularly in bilingual contexts where assessments must account for linguistic proficiency in multiple languages (Wang et al., 2020; Gorman et al., 2015; Wang et al., 2024). With the rapid expansion of telehealth technologies recently, especially after COVID-19, more bilingual families can access essential medical resources remotely, such as interacting with web-based child language assessment tools at home (Pratt et al., 2022; Dam and Pham, 2023).

However, bilingual speech-language pathologists (SLPs) are already in severe shortage in the United States because of the challenging but necessary expertise to evaluate children’s language abilities across different linguistic and cultural contexts (Du et al., 2020; Pratt et al., 2022; Dam and Pham, 2023). The adoption of telehealth technologies further poses a significant workload for SLPs because they need to spend a considerable amount of time reviewing large-scale data collected remotely.

During remote assessments, parents are required to supervise and facilitate children’s interaction with the telehealth tool. For instance, parents could help provide technical assistance with digital platforms and offer behavioral support (Pozniak et al., 2024; Edwards-Gaither et al., 2023). However, due to the lack of specialized training, parents can unintentionally exhibit interference behaviors, such as repeating assessment instructions and “leaking” the correct answer, which can compromise the validity of assessment results (Du et al., 2020; Tomlinson et al., 2018). Identifying these behaviors typically requires SLPs to conduct manual transcription and meticulous behavioral coding from video-recorded sessions, which is an extremely time-intensive and laborious process (Sun et al., 2024; Cao et al., 2019; Lønfeldt et al., 2023).

Recent advances in Natural Language Processing, particularly large language models (LLMs), have shown promise in automating complex behavioral coding tasks in domain-specific contexts like motivational interviewing and mental health counseling (Cao et al., 2019; Tavabi et al., 2020; Mayer et al., 2024; Pellemans et al., 2024). Nevertheless, these applications are mainly limited to monolingual, adult, and face-to-face interactions, and the application scenarios, unlike medical assessments, often do not require highly professional expertise and strict step-by-step requirements. Little research has explored the potential of

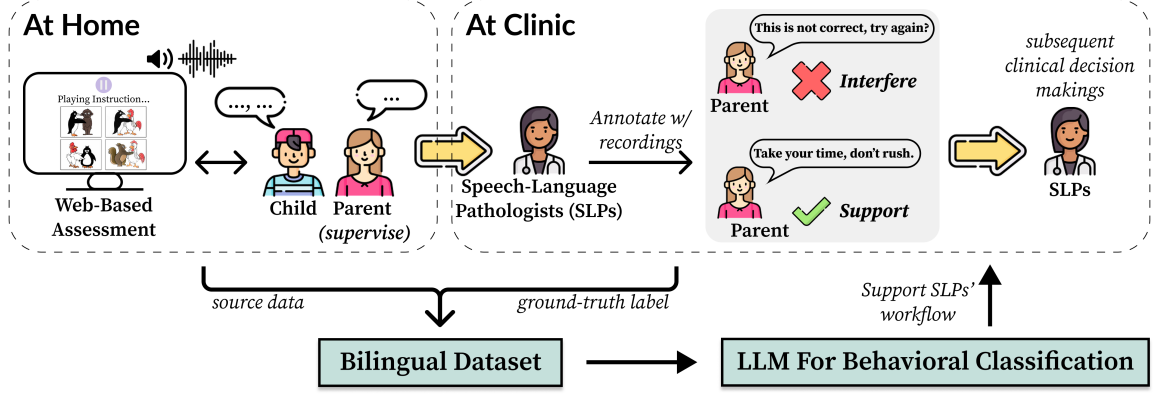


Figure 1: The clinical workflow of speech-language pathologists (SLPs) for remote patient assessment.

LLMs to support clinicians’ workflow within the unique context of remote language assessments for bilingual children, where the interplay of linguistic, cultural, and technological factors significantly complicates analysis (Zhang et al., 2023a,b; Lin et al., 2022; Karacan et al., 2024).

In this paper, we address this critical gap by systematically examining how state-of-the-art (SOTA) LLMs, including open-sourced and closed-sourced ones, can support bilingual SLPs’ clinical workflows by automating the coding of parental behaviors during Mandarin-English language assessments. We introduce a novel and comprehensive dataset consisting of transcripts and meticulously annotated behavioral descriptions from 59 bilingual parent-child dyads, including both in-person (16) and telehealth (43) sessions. Our dataset encompasses 1,304 annotated parental behaviors categorized into eight clinically validated labels that are collaboratively developed with bilingual SLPs.

We benchmark the performance of four SOTA LLMs: GPT-4 (Achiam et al., 2023), Llama 3 (Grattafiori et al., 2024), Qwen2 (Yang et al., 2024a), and DeepSeek-V3 (Liu et al., 2024) on this bilingual dataset using zero-shot and few-shot prompting techniques. Our analysis reveals that while some models (GPT-4 and DeepSeek-V3) achieve commendable accuracy, all models are consistently outperformed by expert human annotators, particularly when faced with Mandarin-English utterances. An in-depth error analysis further elucidates critical challenges posed by cultural nuances, contextual complexities, and intricate clinical procedures, which highlight substantial space for future exploration in multilingual and culturally sensitive NLP methodologies.

To our knowledge, this study introduces the

first publicly available bilingual Mandarin-English dataset specifically to encode parental behaviors in bilingual language assessment contexts. Our findings provide a rigorous benchmark and an ecologically valid challenge for NLP research in clinical scenarios, emphasizing the urgent need to enhance the multilingual and domain-specific capabilities of LLMs. By advancing in this direction, NLP technologies hold the potential to significantly alleviate clinician workload, enhance the accuracy and efficiency of clinical assessments, and ultimately improve patient outcomes.

2 Related Work

2.1 Multilingual LLMs for Real-World Tasks

Recent studies have explored LLMs’ capabilities in real-world scenarios that require domain expertise, such as education (Chen et al., 2023; Nayeem and Rafiei, 2024; Valentini et al., 2023; Ghanizadeh and Dousti, 2025) and health care (Rina et al., 2024; Labrak et al., 2024). LLMs like GPT-4 (OpenAI, 2023), Llama (Touvron et al., 2023a,b), Qwen2 (Yang et al., 2024b), and DeepSeek (Bi et al., 2024; Peng et al., 2025; Neha and Bhati) have been trained on multilingual data and demonstrated impressive performance, in tasks like question answering and logical reasoning (Wei et al., 2021; Sanh et al., 2021; Chung et al., 2022), although performances may differ due to language features (e.g., morphosyntax) (Hlavnova and Ruder, 2023; Weissweiler et al., 2023) or use of in-context learning (ICL) (Brown et al., 2020; Zhang et al., 2022; Rubin et al., 2022; Li et al., 2023) or multimodal prompting (Yang et al., 2024c). Prior NLP research has focused on various domain-specific tasks, such as assessment (Wang et al., 2020; Gorman et al., 2015; Laverghetta Jr and Licato, 2023), behavioral anal-

ysis (Van Aken et al., 2021; Sun et al., 2024; Cao et al., 2019; Yang et al., 2023), and narrative tasks (Prud’hommeaux and Roark, 2015; Chen et al., 2023).

2.2 Annotating Clinical Assessment Data

Computational researchers have attempted to annotate clinical data (Leeson et al., 2019) using a variety of speech processing techniques (Pérez-Rosas et al., 2021; Narayanan and Georgiou, 2013) and NLP approaches such as topic modeling, multimodal models (Tavabi et al., 2020; Leeson et al., 2019). Differing from other clinical annotation tasks, assessment tasks may directly impact the diagnostic accuracy of patient care, the annotation accuracy requires psychometric standards to ensure assessment validity and reliability (Abbasi et al., 2021; Wang et al., 2020; Gorman et al., 2015; Laverghetta Jr and Licato, 2023). Applying NLP techniques to bilingual telehealth datasets remained limited due to difficulties in accessing patient data and the high cost of human annotations (Chen et al., 2022). Therefore, to improve clinician workflow and accuracy, novel approaches need to be developed with clinically informed guidelines to support clinicians.

2.3 Behavioral Coding in Clinical NLP

Behavioral coding is a common data analysis methodology in social science (Wang et al., 2022; Black et al., 2013) and has been widely adopted in health and clinical research (Tavabi et al., 2020; Cao et al., 2019; Sun et al., 2024; Mayer et al., 2024; Pellemans et al., 2024). Prior work have utilized computational behavioral recognition for coding video and audio data from parent-child interactions (Lønfeldt et al., 2023). However, conducting a child language assessment requires more fine-grained coding for linguistic features in addition to assessment behaviors (Wang et al., 2020; Gorman et al., 2015), and bilingual data may introduce additional complexity for analysis due to issues such as code-switching (Du et al., 2020).

3 Bilingual Dataset

We collaborated with bilingual (Mandarin-English) speech language pathologists (SLP) to collect and annotate a text-based dataset of bilingual child language assessment sessions. The dataset comprises 59 parent-child dyads using the Mandarin-English Receptive Language Screener

(MERLS) tool. The dataset consists of 16 in-person sessions and 43 telehealth sessions, where the telehealth data was collected through Zoom recordings during COVID-19. An illustration of the clinical workflow, our dataset collection, and model development process is shown in Figure 1.

MERLS is a multi-modal web-based platform designed for assessing language comprehension skills for Mandarin-English-speaking children (Sheng et al., 2021; Du et al., 2020). The assessment consists of a Mandarin module with 44 test items and an English module with 36 test items. For each test item, MERLS plays audio instructions in one language and asks bilingual children to select a picture on the interface that matches the instruction (see Figure 2). Parents are expected to supervise and provide technical support to children during assessments.

3.1 Data Annotation Process

Raw session video recordings were transcribed for speaker utterances verbatim and documented for parents’ verbal and non-verbal behaviors by two research assistants. Next, two bilingual SLPs independently annotate parental behaviors using an established video analysis codebook (Du et al., 2020) developed via Clinical Discourse Analysis (Damico, 1985) to ensure the classification of behaviors was informed based on the existing clinical framework of family-centered assessment practices (Crais et al., 2006). Disagreement between annotators was resolved via member checking method through discussions and codebook refinement until consensus (Birt et al., 2016). Inter-observer agreement reached 97% (in-person split) and 86.1% (remote split) as shown in Table 3.

Parental behavior classification conducted by SLPs is an eight-class classification task, with each behavior assigned one correct label. Table 1 presents the two primary classes and four corresponding sub-categories. “*Interference*” behaviors represent when parents negatively impacted the assessment including “*Repeating Questions (RQ)*”, “*Answering Questions (AQ)*”, “*Analyzing Items (AI)*”, and “*Judging of Correctness (JC)*”; “*Support*” behaviors represent incidents when parents positively facilitated the assessment including “*Encouragement (E)*”, “*Technical Support (TS)*”, “*Broadcasting (B)*”, and “*Miscellaneous (M)*”. For NLP models, each input includes the current test item, a description of the child’s actions, and the parent’s behavior.

Top-Level Category	Sub-Level Category	Definition
Interfere	Repeating Questions (RQ)	Repeating the <Voiceover> audio before and/or during the process of a child selecting the picture on the web.
	Answering Questions (AQ)	Using verbal or gestural cues to suggest or select a correct answer for the child.
	Analyzing Items (AI)	Elaborating on the critical linguistic components by labeling objects and actions, making emphasis via prosodic cues, or breaking down complex sentences from <Voiceover>.
	Judging of Correctness (JC)	Verbally evaluating the child’s response as correct or incorrect.
Support	Encouragement (E)	Showing verbal and/physical affirmation for the child to continue, saying good job/excellent to reinforce the child’s selection, expressing empathy (e.g., it’s okay) on struggled items.
	Technical Support (TS)	Offering verbal (labeled as “Technical Support Verbal”) and/or physical assistance (labeled as “Technical Support Physical”) to the child related to interacting with the website and the computer.
	Broadcasting (B)	After the child makes a selection, describing the selection via a word, a phrase, or a sentence.
	Miscellaneous (M)	Initiating and/or responding to events that redirected a child’s attention, sharing personal opinions about test procedures and stimuli, or other verbal and nonverbal behaviors that were out of the child’s view.

Table 1: The classification categories and corresponding definitions of parent behaviors based on established clinical guidelines. Two high-level categories (“interfere” and “support”) consist of four individual sub-categories.

	In-Person (n=16)		Virtual (n=43)			
	English	Mandarin	All	English	Mandarin	All
RQ	41	101	142	8	30	38
AQ	0	12	12	0	0	0
AI	18	68	86	3	3	6
JC	20	33	53	1	12	13
E	39	68	107	50	83	133
TS	14	78	92	164	206	370
B	14	42	56	4	10	14
M	12	18	30	77	75	152
Interference	79	214	293	12	45	57
Support	79	206	285	295	374	669

Table 2: MERLS dataset (in-person n=16 and virtual n=43) statistics. The top row shows the label distribution across different datasets and test languages.

3.2 Dataset Description & Statistics

The dataset is structured to include the following components: (1) **Time stamps**: Precise time stamps for each assessment item and corresponding parent-child behavior; (2) **GUI descriptions**: Textual descriptions of the graphical user interface (GUI) elements displayed on the MERLS platform; (3) **Audio transcriptions**: Transcriptions of the audio recordings, with annotations that identify different speakers in each voiceover; (4) **Behavior descriptions**: Textual descriptions of parents’ verbal and non-verbal behaviors. An example of such data is illustrated in Figure 6.

Table 2 presents overall statistics for the In-person and Virtual sessions. The two datasets exhibit imbalances in their label distributions: the Virtual split contains fewer interference behaviors and more technical support behaviors. It may be due to (1) an instructional video in the MERLS system aimed to reduce interference behaviors, or

(2) the use of the MERLS system during telehealth assessment which increased the need for technical support activities conducted by parents.

4 Behavior Classification with LLMs

Our experiments focus on zero-shot (ZS) and few-shot (FS) in-context learning (ICL) prompting strategies for LLMs to investigate whether LLMs can reliably classify parental behaviors during child language assessments.

Prompts. Our zero-shot prompt in Figure 3 provides instructions, explains the input format, and defines each of the eight labels. From the test example itself, the model is shown (1) the text of the current question, (2) a description of the child’s behavior, and (3) the description of the parent’s behavior. The few-shot prompt is similar but includes one demonstrative example (by a clinical expert to ensure validity) under each label definition. Figures 4 and 5 in the Appendix show the few-shot prompts, split over multiple pages. Both prompts include (1) **Voiceover**, the text of the current question; (2) **Child behavior**, a description of the child’s behavior, and (3) **Parent behavior**, the utterance and/or a description of the action performed by the parent.

Models. To evaluate NLP performance on our tasks, we experiment with four LLMs: the open-weight models Llama-3-8B-Instruct, DeepSeek-V3, and Qwen2-7B, and the closed-source model GPT-4 (Turbo-2024-04-09). While Llama-3 is primarily English-based, its pre-training data includes data from 30 other languages.¹

¹<https://ai.meta.com/blog/meta-llama-3/>

Micro F1 / Macro F1 (%) on:	In-Person (n=16)			Virtual (n=43)		
	English	Mandarin	Overall	English	Mandarin	Overall
Llama3 ZS	48.7/37.4	43.6/33.7	45.0/34.6	29.6/21.9	31.5/20.8	30.7/21.0
Llama3 FS	45.6/31.8	39.8/26.2	41.3/27.6	22.1/9.3	23.6/10.1	23.0/9.9
GPT-4 ZS	65.8/58.5	60.5/55.7	61.9/ 57.2	45.3/22.3	52.7/38.2	49.6/33.8
GPT-4 FS	66.5/ 61.1	55.0/49.1	58.1/52.4	50.2/ 28.8	53.9/38.2	52.3/ 36.2
DeepSeek-V3 ZS	61.7/53.5	58.5/51.2	59.4/52.4	47.6/23.6	53.2/ 38.9	50.8/34.3
DeepSeek-V3 FS	64.6/54.9	63.5/ 56.2	63.8/56.3	52.4/25.2	53.7/36.2	53.2/33.3
Qwen2 ZS	31.4/19.6	28.3/19.6	29.2/20.0	38.8/13.9	40.7/20.3	39.9/17.8
Qwen2 FS	20.1/15.2	25.6/21.5	24.1/19.7	41.6/17.0	42.1/24.2	41.9/22.1
Human Experts	96.84	96.43	97.0	86.93	81.82	86.1

Table 3: Micro/Macro F1 results on assessment language for In-Person and Virtual dataset. ZS = zero-shot, FS = few-shot. The best-performing values for each metric are highlighted.

Macro F1 on:	In-Person (n=16)		Virtual (n=43)	
	Interference	Support	Interference	Support
Llama3 ZS	75.2	70.5	18.0	76.4
Llama3 FS	73.5	64.6	20.3	69.3
GPT-4 FS	87.6	86.7	37.4	90.1
DeepSeek-V3 ZS	78.0	81.4	20.4	86.1
DeepSeek-V3 FS	85.0	85.9	29.2	90.6
Qwen2 ZS	54.7	66.0	25.5	87.0
Qwen2 FS	56.3	57.4	28.2	84.7

Table 4: Macro F1 on the binary classification version ("Interference" vs. "Support" behaviors) of In-Person and Virtual dataset. ZS = zero-shot, FS = few-shot.

Qwen2 (Yang et al., 2024b) achieved strong benchmarks across approximately 30 languages. GPT-4 has also demonstrated strong performance on Chinese language understanding benchmarks (Xu et al., 2023; Zhu et al., 2024). Deepseek-V3 is optimized for computational efficiency and excels in complex linguistic and reasoning tasks with minimal supervised data (Liu et al., 2024).

Evaluation metrics. To accurately evaluate parent behaviors across different test items in Mandarin and English, we compute three metrics: (1) Macro F1 score (MACRO): prediction performance addressing the effects of dataset imbalance; (2) Micro F1 score (MICRO): precision and recall computed over all prediction instances; (3) Item-level Accuracy (ITEMACC): the proportion of items with at least one behavior where all behaviors are predicted correctly. ITEMACC is calculated as the number of items with at least one behavior that is predicted entirely correctly, divided by the total number of items that contained at least one behavior in each item.

5 Experimental Results

5.1 Main Results

Table 3 presents the micro- and macro-F1 scores for each model in the assessment items of English

and Mandarin in both the in-person and virtual settings. Overall, GPT-4 and DeepSeek-V3 consistently outperform Llama3 and Qwen2 across all settings. Their stronger performance is in line with established scaling laws in large language models (LLMs), rather than indicating domain-specific adaptation (Wei et al., 2021)

The consistent gap between macro and micro F1 across settings reflects the severe class imbalance in our dataset. Since macro F1 gives equal weight to each class, even infrequent behaviors such as Analyzing Items, which appears only three times in the Virtual dataset (Table 2), can significantly impact the overall score. In the following interpretations, we emphasize macro F1 because it provides a more rigorous and fair evaluation across all behavior categories, especially in imbalanced datasets where some behaviors are rare but clinically meaningful.

On the In-Person dataset, GPT-4 achieves the highest macro F1 score on English assessment items (61.1%), while DeepSeek-V3 leads in Mandarin (56.3%). The strong Mandarin performance of DeepSeek-V3 likely reflects its exposure to high-quality Chinese-language sources (Guo et al., 2024). Llama3 shows moderate performance in both languages, consistently lagging behind GPT-4 and DeepSeek-V3 but outperforming Qwen2, which remains the weakest model, especially in English, even with few-shot support.

On the Virtual dataset, all models perform worse than in the In-Person setting. GPT-4 and DeepSeek-V3 remain the top performers, while Llama3 shows moderate performance, and Qwen2 ranks lowest. Interestingly, this performance drop occurs despite a higher number of annotated parental behaviors in Virtual sessions. This suggests that the challenge lies not in annotation spar-

sity but in the increased variability and complexity of behaviors exhibited during virtual interactions, which may be harder for models to learn and classify accurately. These findings highlight the unique difficulty of modeling parental behaviors in remote formats, where behavior patterns may be different from those in in-person settings.

In both in-person and virtual settings, model performance varied across clinical contexts, and the paired-setting analysis (Table 9), based on few-shot results, provides clearer insight into how these setting differences influence classification. GPT-4 exhibited the most pronounced drop, with macro F1 scores decreasing from 53.4% in the in-person setting to 27.7% in the virtual setting, particularly struggling with English in the virtual context. DeepSeek-V3, while still achieving the highest overall scores, also showed a substantial decline from 56.3% to 28.1%, indicating that it too was affected by setting-related challenges. Qwen2 was the only model to improve in the Virtual setting, rising from 19.7% to 22.6%, with its best performance observed in virtual Mandarin sessions (25.9%). In contrast, Llama3 performed poorly across all conditions, with macro F1 scores falling from 27.6% in-person to just 8.8% in virtual sessions, and without a consistent language-specific trend. These findings emphasize that the clinical setting (in-person or virtual) has a substantial impact on model performance, and that the ability to generalize across settings varies widely by model, language-specific pretraining, and robustness to behavioral variability.

5.2 Effects of Parent Language

Table 6 presents micro and macro F1 scores broken down by the language used to describe parent behaviors in the transcripts: English, Mandarin, or a mix of both. In our dataset, non-verbal parent behaviors are consistently described in English. Code-mixing occurs when parents code-switch during speech or when Mandarin utterances are paired with English descriptions of non-verbal actions. Across most models, macro F1 scores are highest on English or Mixed-language transcripts and lowest on Mandarin-only transcripts, especially in the Virtual setting. GPT-4 demonstrates the most consistent performance across languages, with balanced macro F1 scores for English and Mandarin transcriptions in both settings (e.g., 51.3% vs. 46.7% in-person; 35.4% vs. 34.7% virtual). It achieves the best result in all models on

Mandarin-only data in the Virtual setting.

In contrast, DeepSeek-V3 performs best on English and Mixed-language transcripts, but its macro F1 on Virtual Mandarin (29.9%) is lower than GPT-4s. Qwen2 shows a relatively small performance gap between English and Mandarin, although its overall accuracy remains low, with its highest macro F1 score (27.1%) occurring on Mixed-language Virtual transcripts. Llama3 performs poorly across all conditions, with its lowest macro F1 on Virtual Mandarin transcripts (6.9%). These results demonstrate that transcription language significantly influences model performance, with Mandarin-only transcripts, particularly in Virtual sessions, posing the greatest challenge.

Interestingly, while DeepSeek-V3 outperforms GPT-4 in overall Mandarin session performance (Table 9), GPT-4 surpasses DeepSeek when classifying Mandarin-only behavior transcriptions (Table 6), especially in the Virtual setting (34.7% vs. 29.9%). This discrepancy highlights the distinction between *session language* (i.e., whether the child was tested in Mandarin) and the *transcription language* (i.e., whether the parents behavior was described using Mandarin). DeepSeek-V3 may be more effective at capturing contextual patterns in full Mandarin-language sessions, whereas GPT-4 appears better at parsing Mandarin within isolated transcription entries. These findings underscore the importance of distinguishing between session language and transcription language when evaluating multilingual performance in LLM-based behavioral classification.

5.3 Binary Classification Results

Table 4 presents macro F1 scores for the binary classification task distinguishing Interference and Support behaviors across In-Person and Virtual sessions. This analysis enables a clearer understanding of how LLMs perform when behavior categories are simplified, reducing the impact of class imbalance present in the 8-way classification task.

Across all models, performance is substantially higher in the In-Person setting than in the Virtual setting for both behavior types. GPT-4 achieves the highest overall performance, with macro F1 scores of 87.6% for Interference and 86.7% for Support in the In-Person setting. While its performance drops in the Virtual setting, it still maintains relatively strong accuracy, especially on Support behaviors (90.1%). DeepSeek-V3 closely follows, with similar In-Person scores (85.0% and

In-Person English				In-Person Mandarin			
Label pair	Annotator 1 Acc	Annotator 2 Acc	Overall	Label pair	Annotator 1 Acc	Annotator 2 Acc	Overall
RQ-AI	100.0	75.0	75.0	RQ-AI	60.0	70.0	60.0
AI-RQ	0.0	100.0	0.0	AI-RQ	69.0	100.0	69.0
JC-E	71.4	71.4	57.1	E-RQ	22.2	0.0	0.0
				E-M	40.0	16.6	16.6
				TS-B	100.0	100.0	100.0

Virtual English				Virtual Mandarin			
Label pair	Annotator 1 Acc	Annotator 2 Acc	Overall	Label pair	Annotator 1 Acc	Annotator 2 Acc	Overall
TS-M	82.2	53.3	53.3	TS-M	85.7	90.2	74.5
M-E	75	92.3	68.8	M-E	68.8	61.5	61.5

Table 5: Clinician annotation accuracy based on the misclassified pairs from Figure 14. A-B denotes that A is the true clinician-annotated label, while B represents the GPT prediction. The overall accuracy is calculated as the number of correctly classified behaviors for both clinicians divided by the total number of misclassified behavior pairs. These low values (highlighted in the table) suggest that these misclassified pairs are also somewhat challenging for clinicians to classify accurately.

85.9%) and slightly lower Virtual performance on Interference (29.2%) but the highest score on Support (90.6%). Llama3 and Qwen2 perform considerably worse on Interference behaviors in Virtual settings, with macro F1 scores below 30%. However, their performance on Support behaviors remains relatively strong in Virtual contexts (e.g., 84.7% for Qwen2 FS and 76.4% for Llama3 ZS). This suggests that Support behaviors are more consistently recognized across models, whereas Interference behaviors are more difficult to detect, especially in Virtual sessions where contextual cues may be limited or harder to interpret.

These findings reinforce that LLM performance degrades in virtual environments, particularly for subtle or ambiguous behavioral categories like Interference. However, simplifying the task to binary classification improves overall accuracy and highlights meaningful variation across model architectures and prompting strategies.

5.4 Error Analysis with Human Annotators

While our primary evaluation focused on macro and micro F1 scores (Tables 3 and 4), we conducted a detailed error analysis using item-level accuracy to identify challenging behavior pairs, focusing on one of the top-performing models: GPT-4 ZS for the In-Person dataset, and GPT-4 FS for the Virtual dataset. The confusion matrices in the Appendix Figures 14 identified the misclassified pairs. To examine these errors, we selected the most frequently misclassified pairs for each behavioral category within each dataset (Table 5). These misclassified data were selected based on two criteria: 1) the number of instances in the class is no

less than the average of that dataset, and 2) the prediction accuracy for that class is below 80%. Two novel clinical expert annotators were trained using the same clinical protocol in Table 3 for error analysis. Novice human annotators were not utilized due to the specialized training required for behavioral coding, without such training, they have variability and reliability issues. Two annotators independently selected which option they believe is the correct answer without knowing which class was machine or human annotation. By comparing their responses to the true labels, we identified challenging pairs the clinical experts also struggled to perform classification.

Errors in In-Person Dataset. Several commonly misclassified pairs by GPT-4 can be identified in the top rows from Table 5. Overall, the LLM struggles to distinguish "Repeating Questions (RQ)" from "Analyzing Items" in both English and Mandarin tests. Notably, this distinction is also challenging for clinicians, as the AI-RQ category for the In-Person dataset (English) shows 0.0% accuracy, with significant disagreement among experts. Figure 6 illustrates a sample parent-child interaction transcript for this disagreement between the best-performing model (GPT4-ZS) and two human annotators.

A key insight from our analysis is that LLM errors reveal potential *linguistic ambiguities* in the clinical definition of parent behaviors. GPT-4 mislabels "Analyzing Items (AI)" as "Repeating Questions (RQ)" when parents repeat only key components of a question (e.g., "wash the cat") for their children across both in-person and virtual datasets.

This is likely because GPT-4 lacks the specific details to differentiate by repeating how much of the partial question is considered "Analyzing Items." The errors also appeared in code-mix utterances and could be due to inadequate translation from word-level lexicon to sentence-level utterances. Through the disagreements between human annotators and predictions from LLMs, we gain a deeper insights onimproving both the development of a clinical annotation codebook (Leeson et al., 2019) as well as the prompting strategies for LLMs (Lin et al., 2022; Ranaldi and Pucci, 2023; Lønfeldt et al., 2023).

Errors in the Virtual Dataset. In the virtual dataset, the common misclassifications are demonstrated in the bottom rows from Table 5. We found the challenging pairs for clinicians to classify are "TS-M" ("Technical Support" vs. "Miscellaneous") for Virtual English and "M-E" ("Miscellaneous" vs. "Encouragement") for Virtual dataset in Mandarin. Figure 7 in the appendix illustrates a sample transcript, showing disagreement across two human annotators and the best-performing model (GPT-4 FS). GPT-4 appears to overgeneralize utterances that do not contain words related to "website" or "computer" as "Technical Support" behaviors. This is indeed due to *behavioral definition ambiguities* originated from the clinician's codebook which informed the prompting. After two annotators conducted the member-checking procedure (Birt et al., 2016) to discuss inconsistent annotations, they identified additional utterances (e.g., parents monitoring children's needs for breaks or snacks) as a new potential category of support behavior in the virtual dataset. This error analysis process highlights the models limitations in analyzing novel information, suggesting that future use of more aligned behavioral coding definitions to inform prompting for LLMs to accurately identify existing categories accurately and recognize novel patterns to enhance behavioral coding process. In a deeper error analysis, we conducted an ablation study to assess the impact of various components of GPT-4's performance, see results in Appendix B.

6 Conclusion

This paper introduces a bilingual dataset for classifying parental behaviors during English-Mandarin child language assessments. This study highlights the potential of LLMs to support bilingual lan-

guage assessment in clinical and school settings. By identifying supportive and interfering parental behaviors, automated classification can inform improvements in family-centered assessment practices with enhanced clinical validity. Additionally, LLMs may help alleviate the shortage of bilingual SLPs by supporting screening processes, expanding access to linguistically appropriate evaluations, and reducing the risk of over-diagnosing language disorders in emergent bilingual children.

While current SOTA LLMs show moderate accuracy, they struggle with Mandarin data, a challenge also faced by human annotators, particularly with virtual data. These difficulties point to potential ambiguities in behavioral definitions, especially in cases involving code-switching and nuanced parental language. This highlights the need to refine the clinical annotation codebook, which may improve consistency across both human and automated coding efforts. This dataset promotes further NLP research for multilingual clinical tasks, advancing the analysis of using multimodal behavioral coding (Yang et al., 2024c) of bilingual datasets (Hlavnova and Ruder, 2023; Weissweiler et al., 2023) during child language assessment in in-person and telehealth contexts.

7 Future Work

Our bilingual dataset also contained prosodic information (e.g., parents emphasize words when analyzing items with different stress patterns) which has been shown to introduce additional complexity in clinical NLP (Black et al., 2013), especially in the bilingual context (Pattichis et al., 2023). Manual transcription by clinicians introduces potential for error and inconsistency. To reduce such errors, future work should continue exploring available LLMs to achieve better performance or integration of multimodal speech (OpenAI Whisper, Llama-Omni, Qwen2-Audio) that can automate audio-to-text transcription for additional technical novelty in this research (Tavabi et al., 2020). By aligning our methodology with clinical annotation guidelines and error analysis, future work may include collaborating with clinicians to refine annotation guidelines and additional error analysis, and aligning model outputs with assessment protocols, and piloting LLM-supported language assessments in clinical settings for bilingual children.

8 Limitations

Our study is constrained by the imbalance between the in-person and virtual datasets, as well as a relatively small sample size with (1,304 parent behavior annotations from 59 sessions), which further limited underrepresented classification categories (e.g., broadcasting). Such a small sample size further affect generalizability and classification reliability; additionally, the imbalance in label distribution across the eight categories contributes to this limitation, affecting overall model performance. For instance, smaller subcategories like "Broadcasting" make up less than 1% compared to larger categories such as "Technical Support" impacting overall model evaluation. Despite the small sample size, our clinically informed annotation guideline could be extended to the ongoing data collection with our clinical partners; it can also inform clinical research in other bilingual populations, such as Spanish-English and Vietnamese-English speaking children and parents (Dam and Pham, 2023; Pratt et al., 2022).

Similar to other closed-source models, our best-performing model GPT-4 poses challenges for replication. One challenge is the semantic alignment at the word and utterance level and how this impacts behavioral classification, especially in bilingual datasets (Cao et al., 2019; Huzaifah et al., 2024). Although our behavioral classifications included categories such as "Encouragement," due to the scope of this paper, we did not explore the use of sentiment analysis (Zhang et al., 2023a) on specific linguistic features in the transcripts to improve accuracy for individual subcategories.

Additionally, error analysis revealed challenges in classifying behaviors related to the "Miscellaneous" category for the best performing model GPT-4 and human annotators. Due to the scope of this paper, we did not conduct additional ablation studies to evaluate more specifically defined behavioral categories. Since this is a bilingual dataset from a novel clinical pilot study, ongoing collaborations between NLP researchers and clinicians may continue to improve the behavioral coding protocol as well as the performance of LLMs.

Furthermore, the lack of improvement in model performance with ICL suggests that the primary bottleneck in this classification task may not arise from challenges in semantic understanding of the text, but rather from capturing the deeper, underlying intent embedded in parental language.

9 Ethical Considerations

Data collection and analysis. Our dataset was collected and processed under university human subject research approval and data sharing agreements. The de-identified text transcripts from the clinical video analysis being used for model evaluation contain no sensitive information about parent-child pairs.

Bias mitigation. When comparing the 16 virtual and 16 in-person parent-child pairs, we considered the effects of childrens age and parents education level, which can influence language abilities and parental behavior during assessments. A detailed parent-child demographic table for the 16 In-person and Virtual pairs are provided in Table 11.

Privacy and trust. Privacy issues are paramount when annotating parent behaviors and utterances. Annotators anonymized sensitive information in the transcripts (e.g., using boys name). This is crucial when applying LLMs to clinical data to protect patient privacy. Implementing LLMs in evaluating clinical data could lead to biases that affect clinical decision-making. Clinicians may rely on models without fully understanding their limitations. We also evaluated GPT-4S misclassifications alongside two human annotators to compare decision-making and address biases, which is essential for building trust and explainable AI in clinical settings.

References

- Ahmed Abbasi, David Dobolyi, John P Lalor, Richard G Netemeyer, Kendall Smith, and Yi Yang. 2021. Constructing a psychometric testbed for fair natural language processing. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiusi Du, Zhe Fu, et al. 2024. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*.
- Linda Birt, Suzanne Scott, Debbie Cavers, Christine Campbell, and Fiona Walter. 2016. Member checking: a tool to enhance trustworthiness or merely

743	a nod to validation? <i>Qualitative health research</i> ,	800
744	26(13):1802–1811.	801
745	Matthew P Black, Athanasios Katsamanis, Brian R	802
746	Baucom, Chi-Chun Lee, Adam C Lammert, An-	803
747	drew Christensen, Panayiotis G Georgiou, and	804
748	Shrikanth S Narayanan. 2013. Toward automat-	805
749	ing a human behavioral coding system for married	
750	couples interactions using speech acoustic features.	806
751	<i>Speech communication</i> , 55(1):1–21.	807
752	Tom Brown, Benjamin Mann, Nick Ryder, Melanie	808
753	Subbiah, Jared D Kaplan, Prafulla Dhariwal,	809
754	Arvind Neelakantan, Pranav Shyam, Girish Sastry,	810
755	Amanda Askell, Sandhini Agarwal, Ariel Herbert-	
756	Voss, Gretchen Krueger, Tom Henighan, Rewon	811
757	Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu,	812
758	Clemens Winter, Chris Hesse, Mark Chen, Eric	813
759	Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,	814
760	Jack Clark, Christopher Berner, Sam McCandlish,	815
761	Alec Radford, Ilya Sutskever, and Dario Amodei.	816
762	2020. Language models are few-shot learners . In	817
763	<i>Advances in Neural Information Processing Systems</i> ,	818
764	volume 33, pages 1877–1901. Curran Associates,	819
765	Inc.	820
766	Jie Cao, Michael Tanana, Zac E Imel, Eric Poitras,	821
767	David C Atkins, and Vivek Srikumar. 2019. Ob-	822
768	servating dialogue in therapy: Categorizing and	
769	forecasting behavioral codes. <i>arXiv preprint</i>	823
770	<i>arXiv:1907.00326</i> .	824
771	Jiaju Chen, Yuxuan Lu, Shao Zhang, Bingsheng Yao,	825
772	Yuanzhe Dong, Ying Xu, Yunyao Li, Qianwen	826
773	Wang, Dakuo Wang, and Yuling Sun. 2023. Fairy-	827
774	talecqa: Integrating a commonsense knowledge	
775	graph into children’s storybook narratives. <i>arXiv</i>	828
776	<i>preprint arXiv:2311.09756</i> .	829
777	Zhuohao Chen, Nikolaos Flemotomos, Zac E Imel,	830
778	David C Atkins, and Shrikanth Narayanan. 2022.	831
779	Leveraging open data and task augmentation to au-	832
780	tomated behavioral coding of psychotherapy conver-	833
781	sations in low-resource scenarios. <i>arXiv preprint</i>	834
782	<i>arXiv:2210.14254</i> .	835
783	Hyung Won Chung, Le Hou, Shayne Longpre, Bar-	836
784	ret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi	837
785	Wang, Mostafa Dehghani, Siddhartha Brahma, et al.	
786	2022. Scaling instruction-finetuned language mod-	838
787	els. <i>arXiv preprint arXiv:2210.11416</i> .	839
788	Elizabeth R Crais, Vicky Poston Roy, and Karen Free.	840
789	2006. Parents’ and professionals’ perceptions of the	841
790	implementation of family-centered practices in child	842
791	assessments.	843
792	Quynh Diem Dam and Giang T Pham. 2023. Re-	844
793	mote first-language assessment: Feasibility study	845
794	with vietnamese bilingual children and their care-	846
795	givers. <i>Language, Speech, and Hearing Services in</i>	847
796	<i>Schools</i> , 54(2):618–635.	848
797	J Damico. 1985. Clinical discourse analysis: a func-	849
798	tional approach to language assessment. <i>Communi-</i>	850
799	<i>cation skills and classroom success</i> .	851
	Yao Du, Li Sheng, and Katie Salen Tekinbas. 2020.	852
	"try your best" parent behaviors during adminis-	853
	tration of an online language assessment tool for	854
	bilingual mandarin-english children. In <i>proceedings</i>	855
	<i>of the interaction design and children conference</i> ,	856
	pages 409–420.	857
	Lesley Edwards-Gaither, Ovetta Harris, and Valencia	
	Perry. 2023. Viewpoint telepractice 2025: Explor-	
	ing telepractice service delivery during covid-19 and	
	beyond. <i>Perspectives of the ASHA Special Interest</i>	
	<i>Groups</i> , 8(2):412–417.	
	Mohammad Amin Ghanizadeh and Mohammad Javad	
	Dousti. 2025. Towards data-efficient language mod-	
	els: A child-inspired approach to language learning.	
	<i>arXiv preprint arXiv:2503.04611</i> .	
	Kyle Gorman, Steven Bedrick, Géza Kiss, Eric Morley,	
	Rosemary Ingham, Metrah Mohammad, Katina Pa-	
	padakis, and Jan PH van Santen. 2015. Automated	
	morphological analysis of clinical language samples.	
	In <i>Proceedings of the conference. Association for</i>	
	<i>Computational Linguistics. North American Chap-</i>	
	<i>ter. Meeting</i> , volume 2015, page 108. NIH Public	
	Access.	
	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	
	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	
	Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-	
	ten, Alex Vaughan, et al. 2024. The llama 3 herd of	
	models. <i>arXiv preprint arXiv:2407.21783</i> .	
	Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie,	
	Kai Dong, Wentao Zhang, Guanting Chen, Xiao	
	Bi, Yu Wu, YK Li, et al. 2024. Deepseek-	
	coder: When the large language model meets	
	programming—the rise of code intelligence. <i>arXiv</i>	
	<i>preprint arXiv:2401.14196</i> .	
	Ester Hlavnova and Sebastian Ruder. 2023. Em-	
	powering cross-lingual behavioral testing of nlp	
	models with typological features. <i>arXiv preprint</i>	
	<i>arXiv:2307.05454</i> .	
	Muhammad Huzaifah, Weihua Zheng, Nattapol Chan-	
	paist, and Kui Wu. 2024. Evaluating code-	
	switching translation with large language models. In	
	<i>Proceedings of the 2024 Joint International Con-</i>	
	<i>ference on Computational Linguistics, Language</i>	
	<i>Resources and Evaluation (LREC-COLING 2024)</i> ,	
	pages 6381–6394.	
	Baris Karacan, Ankit Aich, Avery Quynh, Amy	
	Pinkham, Philip Harvey, Colin Depp, and Natalie	
	Parde. 2024. Towards comprehensive language an-	
	alysis for clinically enriched spontaneous dialogue. In	
	<i>Proceedings of the 2024 Joint International Con-</i>	
	<i>ference on Computational Linguistics, Language</i>	
	<i>Resources and Evaluation (LREC-COLING 2024)</i> ,	
	pages 16457–16472.	
	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-	
	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-	
	guage models are zero-shot reasoners. <i>Advances in</i>	
	<i>neural information processing systems</i> , 35:22199–	
	22213.	

858	Yanis Labrak, Adrien Bazoge, Emmanuel Morin,	Fnu Neha and Deepshikha Bhati. A survey of deepseek	913
859	Pierre-Antoine Gourraud, Mickael Rouvier, and	models . <i>Authorea Preprints</i> .	914
860	Richard Dufour. 2024. Biomistral: A collection of		
861	open-source pretrained large language models for	OpenAI. 2023. Gpt-4 technical report . <i>Preprint</i> ,	915
862	medical domains. <i>arXiv preprint arXiv:2402.10373</i> .	arXiv:2303.08774.	916
863	Antonio Laverghetta Jr and John Licato. 2023. Gen-	Rebecca Pattichis, Dora LaCasse, Sonya Trawick, and	917
864	erating better items for cognitive assessments using	Rena Cacoullos. 2023. Code-switching metrics us-	918
865	large language models. In <i>Proceedings of the 18th</i>	ing intonation units. In <i>Proceedings of the 2023</i>	919
866	<i>Workshop on Innovative Use of NLP for Building Ed-</i>	<i>Conference on Empirical Methods in Natural Lan-</i>	920
867	<i>ucational Applications (BEA 2023)</i> , pages 414–428.	<i>guage Processing</i> , pages 16840–16849.	921
868	William Leeson, Adam Resnick, Daniel Alexander,	Mathijs Pellemans, Salim Salmi, Saskia Mérelle, Wilco	922
869	and John Rovers. 2019. Natural language process-	Janssen, and Rob van der Mei. 2024. Automated be-	923
870	ing (nlp) in qualitative public health research: a	havioral coding to enhance the effectiveness of moti-	924
871	proof of concept study. <i>International Journal of</i>	vatational interviewing in a chat-based suicide preven-	925
872	<i>Qualitative Methods</i> , 18:1609406919887021.	tion helpline: secondary analysis of a clinical trial.	926
873	Xiaonan Li, Kai Lv, Hang Yan, Tianyang Lin, Wei Zhu,	<i>Journal of medical internet research</i> , 26:e53562.	927
874	Yuan Ni, Guotong Xie, Xiaoling Wang, and Xipeng	Yifan Peng, Bradley A Malin, Justin F Rousseau, Yan-	928
875	Qiu. 2023. Unified Demonstration Retriever for In-	shang Wang, Zihan Xu, Xuhai Xu, Chunhua Weng,	929
876	Context Learning . In <i>Proceedings of the 61st An-</i>	and Jiang Bian. 2025. From gpt to deepseek: Signif-	930
877	<i>annual Meeting of the Association for Computational</i>	icant gaps remains in realizing ai in healthcare.	931
878	<i>Linguistics (Volume 1: Long Papers)</i> , pages 4644–		
879	4668, Toronto, Canada. Association for Computa-	Verónica Pérez-Rosas, Rada Mihalcea, et al. 2021.	932
880	tional Linguistics.	Evaluating automatic speech recognition quality and	933
881	Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu	its impact on counselor utterance coding. In <i>Pro-</i>	934
882	Wang, Shuohui Chen, Daniel Simig, Myle Ott, Na-	<i>ceedings of the Seventh Workshop on Computational</i>	935
883	man Goyal, Shruti Bhosale, Jingfei Du, et al. 2022.	<i>Linguistics and Clinical Psychology: Improving Ac-</i>	936
884	Few-shot learning with multilingual generative lan-	<i>cess</i> , pages 159–168.	937
885	guage models. In <i>Proceedings of the 2022 Con-</i>	Kinga Pozniak, Peter Rosenbaum, and Elaine	938
886	<i>ference on Empirical Methods in Natural Language</i>	Yuen Ling Kwok. 2024. Tasks performed by	939
887	<i>Processing</i> , pages 9019–9052.	parents to enable telepractice for children with	940
888	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,	communication disorders: an interview study with	941
889	Bochao Wu, Chengda Lu, Chenggang Zhao,	clinicians and parents. <i>Disability and Rehabilita-</i>	942
890	Chengqi Deng, Chenyu Zhang, Chong Ruan, et al.	<i>tion</i> , 46(8):1547–1558.	943
891	2024. Deepseek-v3 technical report. <i>arXiv preprint</i>	Amy S Pratt, Jissel B Anaya, Michelle N Ramos, Gi-	944
892	<i>arXiv:2412.19437</i> .	ang Pham, Miriam Muñoz, Lisa M Bedore, and Eliz-	945
893	Nicole Nadine Lønfeldt, Sneha Das, Flavia Dalia Fru-	abeth D Peña. 2022. From a distance: Compari-	946
894	mosu, Anna-Rosa Cecilie Mora-Jensen, Anne Ka-	son of in-person and virtual assessments with adult-	947
895	trine Pagsberg, and Line Clemmensen. 2023.	child dyads from linguistically diverse backgrounds.	948
896	Scaling-up behavioral observation with computa-	<i>Language, Speech, and Hearing Services in Schools</i> ,	949
897	tional behavior recognition.	53(2):360–375.	950
898	Tobias Mayer, Neha Warikoo, Amir Eliassaf, Dana	Emily Prud’hommeaux and Brian Roark. 2015. Graph-	951
899	Atzil-Slonim, and Iryna Gurevych. 2024. Predicting	based word alignment for clinical language evalua-	952
900	client emotions and therapist interventions in psy-	tion. <i>Computational Linguistics</i> , 41(4):549–578.	953
901	chotherapy dialogues. In <i>Proceedings of the 18th</i>	Leonardo Ranaldi and Giulia Pucci. 2023. Does	954
902	<i>Conference of the European Chapter of the Associa-</i>	the english matter? elicit cross-lingual abilities of	955
903	<i>tion for Computational Linguistics (Volume 1: Long</i>	large language models. In <i>Proceedings of the 3rd</i>	956
904	<i>Papers)</i> , pages 1463–1477.	<i>Workshop on Multi-lingual Representation Learning</i>	957
905	Shrikanth Narayanan and Panayiotis G Georgiou. 2013.	(MRL), pages 173–183.	958
906	Behavioral signal processing: Deriving human be-	Nicholas Riina, Likhitha Patlolla, Camilo Hernandez	959
907	havioral informatics from speech and language. <i>Pro-</i>	Joya, Roger Bautista, Melissa Olivar-Villanueva,	960
908	<i>ceedings of the IEEE</i> , 101(5):1203–1233.	and Anish Kumar. 2024. An evaluation of english to	961
909	Mir Tafseer Nayeem and Davood Rafiei. 2024. Kidlm:	spanish medical translation by large language mod-	962
910	Advancing language models for children–early in-	els. In <i>Proceedings of the 16th Conference of the As-</i>	963
911	sights and future directions. <i>arXiv preprint</i>	<i>sociation for Machine Translation in the Americas</i>	964
912	<i>arXiv:2410.03884</i> .	(Volume 2: User Track), pages 222–236.	965

966	Ohad Rubin, Jonathan Herzig, and Jonathan Berant.	Saladi, Alan Schelten, Ruan Silva, Eric Michael	1024
967	2022. Learning To Retrieve Prompts for In-Context	Smith, Ranjan Subramanian, Xiaoqing Ellen Tan,	1025
968	Learning . In <i>Proceedings of the 2022 Conference of</i>	Binh Tang, Ross Taylor, Adina Williams, Jian Xi-	1026
969	<i>the North American Chapter of the Association for</i>	ang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov,	1027
970	<i>Computational Linguistics: Human Language Tech-</i>	Yuchen Zhang, Angela Fan, Melanie Kambadur,	1028
971	<i>nologies</i> , pages 2655–2671, Seattle, United States.	Sharan Narang, Aurelien Rodriguez, Robert Sto-	1029
972	Association for Computational Linguistics.	jnic, Sergey Edunov, and Thomas Scialom. 2023b.	1030
973	Victor Sanh, Albert Webson, Colin Raffel, Stephen H	Llama 2: Open Foundation and Fine-Tuned Chat	1031
974	Bach, Lintang Sutawika, Zaid Alyafeai, Antoine	Models . <i>Preprint</i> , arxiv:2307.09288.	1032
975	Chaffin, Arnaud Stiegler, Teven Le Scao, Arun	Maria Valentini, Jennifer Weber, Jesus Salcido, Téa	1033
976	Raja, et al. 2021. Multitask prompted training en-	Wright, Eliana Colunga, and Katharina Kann. 2023.	1034
977	ables zero-shot task generalization. <i>arXiv preprint</i>	On the automatic generation and simplification of	1035
978	<i>arXiv:2110.08207</i> .	children’s stories. <i>arXiv preprint arXiv:2310.18502</i> .	1036
979	Li Sheng, Danyang Wang, Caila Walsh, Leah Heisler,	Betty Van Aken, Sebastian Herrmann, and Alexander	1037
980	Xin Li, and Pumpki Lei Su. 2021. The bilingual	Löser. 2021. What do you see in this patient? behav-	1038
981	home language boost through the lens of the covid-	ioral testing of clinical nlp models. <i>arXiv preprint</i>	1039
982	19 pandemic. <i>Frontiers in Psychology</i> , 12:667836.	<i>arXiv:2111.15512</i> .	1040
983	Xin Sun, Jiahuan Pei, Jan de Wit, Mohammad Alian-	Danyang Wang, Alexander Choi-Tucci, Anita Mendez-	1041
984	nejadi, Emiel Krahmer, Jos TP Dobber, and Jos A	Perez, Ronald B Gillam, Lisa M Bedore, and Eliz-	1042
985	Bosch. 2024. Eliciting motivational interviewing	abeth D Peña. 2024. Where to start: Use of the	1043
986	skill codes in psychotherapy with llms: A bilingual	bilingual multidimensional ability scale (b-mas) to	1044
987	dataset and analytical study. In <i>Proceedings of the</i>	identify developmental language disorder (dld) in	1045
988	<i>2024 Joint International Conference on Computa-</i>	bilingual children. <i>International Journal of Speech-</i>	1046
989	<i>tional Linguistics, Language Resources and Evalu-</i>	<i>Language Pathology</i> , pages 1–17.	1047
990	<i>ation (LREC-COLING 2024)</i> , pages 5609–5621.	Yilei Wang, Jingyuan Tian, Deniz S Ones, Richard N	1048
991	Leili Tavabi, Kalin Stefanov, Larry Zhang, Brian Bor-	Landers, et al. 2022. Using natural language pro-	1049
992	sari, Joshua D Woolley, Stefan Scherer, and Moham-	cessing and machine learning to replace human con-	1050
993	mad Soleymani. 2020. Multimodal automatic cod-	tent coders. <i>Psychological Methods</i> .	1051
994	ing of client behavior in motivational interviewing.	Yiyi Wang, Emily PrudHommeaux, Meysam Asgari,	1052
995	In <i>Proceedings of the 2020 International Conference</i>	and Jill Dolata. 2020. Automated scoring of clin-	1053
996	<i>on Multimodal Interaction</i> , pages 406–413.	ical expressive language evaluation tasks. In <i>Pro-</i>	1054
997	Serena RL Tomlinson, Nick Gore, and Peter McGill.	<i>ceedings of the conference. Association for Compu-</i>	1055
998	2018. Training individuals to implement applied be-	<i>tational Linguistics. Meeting</i> , volume 2020, page	1056
999	havior analytic procedures via telehealth: A system-	177. NIH Public Access.	1057
1000	atic review of the literature. <i>Journal of Behavioral</i>	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu,	1058
1001	<i>Education</i> , 27:172–222.	Adams Wei Yu, Brian Lester, Nan Du, Andrew M.	1059
1002	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	Dai, and Quoc V. Le. 2021. Finetuned Language	1060
1003	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	Models are Zero-Shot Learners. In <i>International</i>	1061
1004	Baptiste Rozière, Naman Goyal, Eric Hambro,	<i>Conference on Learning Representations</i> .	1062
1005	Faisal Azhar, Aurelien Rodriguez, Armand Joulin,	Leonie Weissweiler, Valentin Hofmann, Anjali Kan-	1063
1006	Edouard Grave, and Guillaume Lample. 2023a.	tharuban, Anna Cai, Ritam Dutt, Amey Hengle,	1064
1007	LLaMA: Open and Efficient Foundation Language	Anubha Kabra, Atharva Kulkarni, Abhishek Vi-	1065
1008	Models . <i>Preprint</i> , arxiv:2302.13971.	jayakumar, Haofei Yu, et al. 2023. Counting the	1066
1009	Hugo Touvron, Louis Martin, Kevin Stone, Peter	bugs in chatgpt’s wugs: A multilingual investiga-	1067
1010	Albert, Amjad Almahairi, Yasmine Babaei, Niko-	tion into the morphological capabilities of a large	1068
1011	lay Bashlykov, Soumya Batra, Prajjwal Bhargava,	language model. <i>arXiv preprint arXiv:2310.15113</i> .	1069
1012	Shruti Bhosale, Dan Bikel, Lukas Blecher, Cris-	Liang Xu, Anqi Li, Lei Zhu, Hang Xue, Chang-	1070
1013	tian Canton Ferrer, Moya Chen, Guillem Cucurull,	tai Zhu, Kangkang Zhao, Haonan He, Xuanwei	1071
1014	David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin	Zhang, Qiyue Kang, and Zhenzhong Lan. 2023. Su-	1072
1015	Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami,	perclue: A comprehensive chinese large language	1073
1016	Naman Goyal, Anthony Hartshorn, Saghar Hos-	model benchmark . <i>Preprint</i> , arXiv:2307.15020.	1074
1017	seini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,	1075
1018	Kerkez, Madian Khabsa, Isabel Kloumann, Artem	Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan	1076
1019	Korenev, Punit Singh Koura, Marie-Anne Lachaux,	Li, Dayiheng Liu, Fei Huang, Guanting Dong, Hao-	1077
1020	Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai	ran Wei, Huan Lin, Jialong Tang, Jialin Wang,	1078
1021	Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov,	Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin	1079
1022	Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew	Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze	1080
1023	Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan		

Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024a. [Qwen2 technical report](#). *Preprint*, arXiv:2407.10671.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. 2024b. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

Jian Yang, Hongcheng Guo, Yuwei Yin, Jiaqi Bai, Bing Wang, Jiaheng Liu, Xinnian Liang, Linzheng Cahi, Liqun Yang, and Zhoujun Li. 2024c. m3p: Towards multimodal multilingual translation with multimodal prompt. *arXiv preprint arXiv:2403.17556*.

Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian Xie, Ziyang Kuang, and Sophia Ananiadou. 2023. Towards interpretable mental health analysis with large language models. *arXiv preprint arXiv:2304.03347*.

Ruochen Zhang, Samuel Cahyawijaya, Jan Christian Blaise Cruz, Genta Indra Winata, and Alham Fikri Aji. 2023a. Multilingual large language models are not (yet) code-switchers. *arXiv preprint arXiv:2305.14235*.

Xiang Zhang, Senyu Li, Bradley Hauer, Ning Shi, and Grzegorz Kondrak. 2023b. Don’t trust chatgpt when your question is not in english: A study of multilingual abilities and types of llms. *arXiv preprint arXiv:2305.16339*.

Yiming Zhang, Shi Feng, and Chenhao Tan. 2022. [Active Example Selection for In-Context Learning](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9134–9148, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Jie Zhu, Junhui Li, Yalong Wen, and Lifan Guo. 2024. [Benchmarking large language models on cflue – a chinese financial language understanding evaluation dataset](#). *Preprint*, arXiv:2405.10542.

A Appendix A

A.1 User Interface of MERLS

Figure 2 illustrates the interaction of a sample test item in the English subtest. When children interact with MERLS, they first hear an audio recording (e.g., "The chicken is hugged by the penguin.") and then select the corresponding picture that matches the audio to demonstrate their understanding of the sentence in English. Parents may or may not demonstrate a behavior depending on their observation of the item and how children respond to the item.

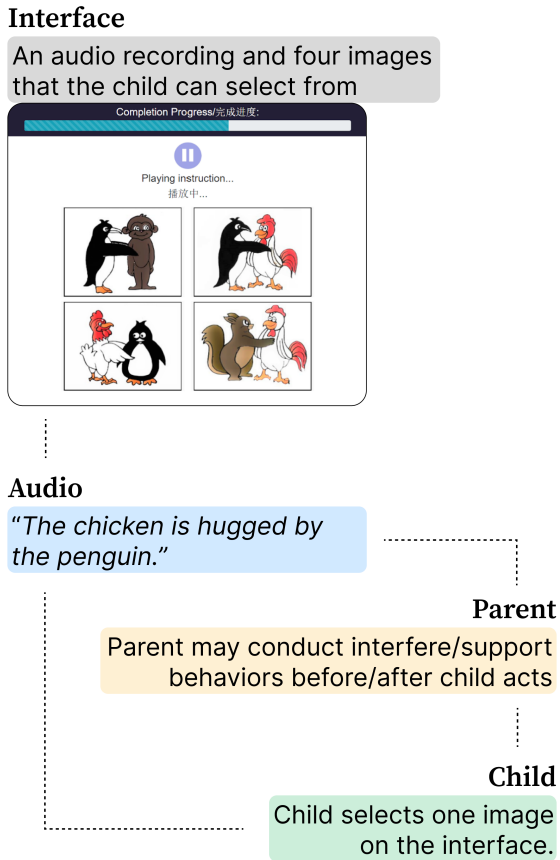


Figure 2: MERLS English test item "The chicken is hugged by the penguin."

A.2 Results of item level accuracy

Here we also provide the accuracy in item level. Table 7 shows the ITEMACC for the entire In-Person dataset (n=16) and the entire Virtual dataset (n=43) partitioned upon question languages, whereas Table 8 shows the ITEMACC for the entire In-person dataset (n=16) and the subset of the Virtual dataset (n=16) that is matched in terms of child age and parent education level.

B Appendix B: Ablation Study

B.1 Removal of role play description in prompting

This experiment tests whether understanding the role and scenario is necessary, or if simple instructions based on 'annotation of parental behavior' suffice, by removing the 'role play' description and retaining only the basic instructions to evaluate both datasets. We tested GPT-4's performance by using a prompt based on our original zero-shot test in Figure 3, but without the role description (e.g., "Assume you are a video analyst classifying transcribed text conversation...") (see prompt in Figure 12). The corresponding results are summarized in Table 10. The maximum variation in accuracy was approximately 3%, and the maximum variation in F1 score was about 4%. Therefore, we conclude that the 'role play' descriptor does not significantly impact the overall prediction results.

B.2 Chain of Thoughts Experiment

To further evaluate the factors contributing to performance variance, we then conducted a Chain of Thoughts (CoT) experiment with GPT-4 to assess the best model performance for in-person data (using zero-shot prompts) and virtual data (using few-shot prompts) after retaining "role play" scenario. We then conducted an error analysis across all categories to compare GPT-4's reasoning with that of human expert annotators. We designed the prompts following a similar two-step procedure as outlined in (Kojima et al., 2022). In the first prompt, we described the classification task as in the original paper, with the addition of a trigger sentence: Lets think step by step. This encourages GPT-4 to generate a step-by-step reasoning process as output (see Figure 8, 9 and 10 for the first prompt input in the Appendix). In the second prompt (see Figure 11 in the Appendix), we combined the original task description with GPT-4's analysis from the first step and include an answer extraction instruction, such as: Therefore, among all the categories, please respond with the category name only. After running the experiment with prompts that include CoT, we then analyzed the initial responses in greater depth by examining the step-by-step reasoning provided. Here we provide one example where GPT-4 makes a false prediction as demonstrated in Figure 13. In this example, GPT-4 falsely predicts the parent behavior to be "Repeat Questions", however the human

expert annotates "Technical Support" because the parent guides the child to make a selection and request a system repetition independently. GPT-4 incorrectly interpreted the parent behavior verbatim as "Repeat Questions" because the utterance itself seems to request repetition; in reality, the parent did not repeat any of the test items and did not interfere with the child. This difference in interpretation provides evidence that GPT models lack specific knowledge of this clinical context, and adding more examples (e.g., few-shot) may not necessarily help model performance.

Micro F1 / Macro F1 (%) on:	In-Person (n=16)			Virtual (n=43)		
	English	Mandarin	Mixed	English	Mandarin	Mixed
# Examples	233	167	178	226	389	111
Llama3 ZS	39.9/31.2	42.5/22.0	53.9/39.6	44.8/27.5	21.1/18.1	36.0/16.5
Llama3 FS	35.6/25.7	40.1/20.8	50.0/33.1	32.7/12.1	14.7/6.91	32.4/14.6
GPT-4 ZS	54.2/47.1	49.1/45.9	59.0/52.1	55.3/35.1	41.9/30.8	55.0/ 32.5
GPT-4 FS	58.8/51.3	56.9/ 46.7	57.9/52.4	56.2/35.4	43.7/ 34.7	47.7/29.7
DeepSeek-V3 ZS	59.2/49.4	58.7/45.7	59.0/50.6	65.9/ 40.0	42.4/30.7	49.5/30.1
DeepSeek-V3 FS	63.9/ 52.9	61.1/45.5	65.7/ 56.6	66.4/39.9	46.3/29.9	50.5/27.1
Qwen2 ZS	34.5/18.6	30.2/17.0	20.4/16.6	41.0/15.2	39.7/18.8	39.6/12.9
Qwen2 FS	27.4/17.9	23.8/18.0	20.6/16.6	44.5/21.9	37.3/19.7	53.0/27.1

Table 6: Micro/Macro F1 results broken down by the language in the transcript of the parent’s behavior (either English, Mandarin, or a mix of both languages). ZS = zero-shot, FS = few-shot. The best-performing values for each metric are highlighted.

ITEMACC(%) on:	In-Person (n=16)			Virtual (n=43)		
	English	Mandarin	All	English	Mandarin	All
Llama3 ZS	39.8	28.0	31.7	26.7	24.3	25.3
Llama3 FS	34.5	27.6	29.8	20.3	19.8	20.0
GPT-4 ZS	61.9	49.2	53.2	36.9	45.3	41.6
GPT-4 FS	61.1	40.8	47.1	43.3	46.5	45.1
DeepSeek-V3 ZS	54.2	44.5	47.5	42.8	43.6	43.3
DeepSeek-V3 FS	57.5	50.9	53.0	45.5	44.9	45.1
Qwen2 ZS	20.4	16.0	17.4	29.4	27.2	28.1
Qwen2 FS	11.5	15.6	14.3	28.9	28.0	28.4

Table 7: Main results for item level accuracy. ZS = zero-shot, FS = few-shot. The best-performing values for each metric are highlighted.

ITEMACC(%) on:	In-Person (n=16)			Virtual (n=16)		
	English	Mandarin	All	English	Mandarin	All
Llama3 ZS	39.8	28.0	31.7	23.2	23.7	23.4
Llama3 FS	34.5	27.6	29.8	15.9	12.9	14.3
GPT-4 ZS	61.9	49.2	53.2	34.1	37.6	36.0
GPT-4 FS	61.1	40.8	47.1	40.2	40.9	40.6
DeepSeek-V3 ZS	54.2	44.5	47.5	34.1	39.8	37.1
DeepSeek-V3 FS	57.5	50.9	53.0	39.0	38.7	38.9
Qwen2 ZS	20.4	16.0	17.4	28.0	20.4	24.0
Qwen2 FS	11.5	15.6	14.3	28.0	23.7	25.7

Table 8: Results with the entire In-person dataset and the subset of the Virtual dataset consisting of 16 patients who are matched with the 16 In-person patients in terms of child age and parent education level in item accuracy. The data layout follows the same procedure described in Table 7.

Micro F1 / Macro F1 (%) on:	In-Person (n=16)			Virtual (n=16)		
	English	Mandarin	Overall	English	Mandarin	Overall
Llama3 ZS	48.7/37.4	43.6/33.7	45.0/34.6	30.4/ 21.7	33.7/16.6	32.1/18.7
Llama3 FS	45.6/31.8	39.8/26.2	41.3/27.6	21.4/8.7	19.8/8.5	20.5/8.8
GPT-4 ZS	68.4/ 61.6	50.0/43.8	55.0/50.9	44.7/19.8	44.2/26.4	44.4/24.3
GPT-4 FS	63.3/58.6	56.0/51.1	58.0/53.4	43.4/19.0	47.7/31.4	45.6/27.7
DeepSeek-V3 ZS	61.7/53.5	58.5/51.2	59.4/52.4	44.0/19.5	51.7/ 34.5	48.0/ 28.9
DeepSeek-V3 FS	64.6/54.9	63.5/ 56.2	63.8/ 56.3	50.9/21.1	49.4/30.7	50.2/28.1
Qwen ZS	31.4/19.6	28.3/19.6	29.2/20.0	41.2/15.9	36.2/19.0	38.6/17.3
Qwen FS	20.1/15.2	25.6/21.5	24.1/19.7	45.7/16.8	38.6/25.9	42.0/22.6

Table 9: Results with the entire In-Person dataset and the subset of the Virtual dataset consisting of 16 patients who are matched with the 16 In-Person patients in terms of child age and parent education level. The data layout follows the same procedure described in Table 3.

Assume you are video analyst classifying transcribed text conversation shown in <Parent Behavior> from parents, who are supervising their bilingual children to answer language comprehension tests in English and Mandarin. A <Voiceover> is the system output sound which reads the picture on the website that the child has to choose. A <Child Behavior> is the children's behavior. You need to classify <Parent Behavior> as described in <task>.

Please respond the category name only.

<Task>

Based on <Parent Behavior>, please determine which type of behavior it is: 'Repeat Questions', 'Answer Questions', 'Analyze Items', 'Judging', 'Encouragement', 'Technical Support', 'Broadcasting', 'Miscellaneous'. Definition for each category is shown in <Definitions>

<Definitions>

- Repeat Questions: Repeating the <Voiceover> audio before and/or during the process of a child selecting the picture on the web.
- Answer Questions: Using verbal or gestural cues to suggest or select a correct answer for the child.
- Analyze Items: Elaborating on the critical linguistic components by labeling objects and actions, making emphasis via prosodic cues, or breaking down complex sentences from <Voiceover>.
- Judging of Correctness: Verbally evaluating the child's response as "correct" or "incorrect".
- Encouragement: Showing verbal and/physical affirmation for the child to continue, saying "good job/excellent" to reinforce the child's selection, expressing empathy (e.g., "it's okay") on struggled items.
- Technical Support: Offering verbal and/or physical assistance to the child related to interacting with the website and the computer.
- Broadcasting: After the child makes a selection, describing the selection via a word, a phrase, or a sentence.
- Miscellaneous: Initiating and/or responding to events that redirected a child's attention, sharing personal opinions about test procedures and stimuli, or other verbal and nonverbal behaviors that were out of the child's view.

Figure 3: The full zero-shot prompt used in our experiments.

	In-Person (n=16)			Virtual (n=43)		
BEHAVACC(%) / F1(%) on:	English	Mandarin	Overall	English	Mandarin	Overall
GPT-4 ZS	67.1/62.1	59.8/56.3	61.8/58.3	41.7/22.8	46.3/34.3	44.4/31.5

Table 10: Prediction results without role play description: BEHAVACC/F1 on assessment language for In-Person and Virtual dataset.

Child ID (In Person)	Age (Year; Month)	Gender	Parent Education	Child ID (Virtual)	Age	Gender	Parent Education
P1	3;8	F	Master	V1	4;0	M	Master
P2	3;8	M	Master	V2	4;1	F	Master
P3	4;1	F	PhD	V3	4;2	M	Master
P4	5;5	F	Master	V4	5;5	F	Master
P5	5;9	M	Master	V5	5;6	M	PhD
P6	5;11	F	Master	V6	5;9	F	Master
P7	6;1	M	PhD	V7	5;9	F	PhD
P8	6;1	M	Master	V8	5;10	F	Master
P9	6;3	M	PhD	V9	5;11	F	PhD
P10	6;4	F	Master	V10	6;3	F	PhD
P11	6;6	M	PhD	V11	6;5	M	Bachelor
P12	6;7	M	PhD	V12	6;11	M	High School
P13	6;7	M	PhD	V13	7;0	M	Bachelor
P14	7;9	M	Vocational School	V14	7;0	F	Bachelor School
P15	8;2	M	Bachelor	V15	8;2	M	Vocational
P16	8;6	M	PhD	V16	8;2	F	PhD

Table 11: Demographic information of gender, age, and parent education for in-person (P1–P16) and virtual study participants (V1–V16). These two groups of children were matched with the comparable level of parent education (except P12 with a parental education of PhD degree, and V12 with a parental education for high school). Childrens ages were also matched with no more than 9 months differences to ensure they are comparable in age for similar language abilities.

Assume you are video analyst classifying transcribed text conversation shown in <Parent Behavior> from parents, who are supervising their bilingual children to answer language comprehension tests in English and Mandarin. A <Voiceover> is the system output sound which reads the picture on the website that the child has to choose. A <Child Behavior> is the children's behavior. You need to classify <Parent Behavior> as described in <task>.

Please respond the category name only.

<Task>
Based on <Parent Behavior>, please determine which type of behavior it is: 'Repeat Questions', 'Answer Questions', 'Analyze Items', 'Judging', 'Encouragement', 'Technical Support', 'Broadcasting', 'Miscellaneous'.

Definition for each category is shown in <Definitions>

<Definitions>
- Repeat Questions: Repeating the <Voiceover> audio before and/or during the process of a child selecting the picture on the web. For example,

<Voiceover>
"the black cat is drinking water"

<Parent Behavior>
"the black cat is drinking water"

<Classification>
Repeat Questions

- Answer Questions: Using verbal or gestural cues to suggest or select a correct answer for the child. For example,

<Voiceover>
"What is the cat drinking?"

<Parent Behavior>
"Drinking water."

<Classification>
Answer Questions

- Analyze Items: Elaborating on the critical linguistic components by labeling objects and actions, making emphasis via prosodic cues, or breaking down complex sentences from <Voiceover>. For example,

<Voiceover>
"the black cat is drinking water"

<Parent Behavior>
"This is the one with a black cat."

<Classification>
Analyze Items

Figure 4: The few-shot prompt used in our experiments, part 1 of 2.

- Judging of Correctness: Verbally evaluating the child's response as "correct" or "incorrect". For example,

<Voiceover>

<Parent Behavior>
"This is not right."

<Classification>
Judging of Correctness

- Encouragement: Showing verbal and/physical affirmation for the child to continue, saying "good job/ excellent" to reinforce the child's selection, expressing empathy (e.g., "it's okay") on struggled items. For example,

<Voiceover>

<Parent Behavior>
"it's fine you are trying your best."

<Classification>
Encouragement

- Technical Support: Offering verbal and/or physical assistance to the child related to interacting with the website and the computer. For example,

<Voiceover>

<Parent Behavior>
"Select the picture to continue."

<Classification>
Technical Support

- Broadcasting: After the child makes a selection, describing the selection via a word, a phrase, or a sentence. For example,

<Voiceover>

<Parent Behavior>
"I selected the picture."

<Classification>
Broadcasting

- Miscellaneous: Initiating and/or responding to events that redirected a child's attention, sharing personal opinions about test procedures and stimuli, or other verbal and nonverbal behaviors that were out of the child's view. For example,

<Voiceover>
<Parent Behavior>
"My child needs to use the bathroom."

<Classification>
Miscellaneous

Figure 5: The few-shot prompt used in our experiments, part 2 of 2.

Timestamp	Speaker	Test Item	Annotator (YD)	Annotator (YY)	GPT4-ZS	Final Agreement
0:20:48	Voiceover	The rabbit that is washing the cat is wearing a ribbon.	n/a	n/a	n/a	n/a
	Child	<attempted to make a selection>	n/a	n/a	n/a	n/a
	Parent	The rabbit washing a cat is wearing a robin.	RQ	RQ	RQ	RQ
	Parent	哪个 rabbit 在 wash the cat? Rabbit 在 wash the cat, 然后它还有 robin. Rabbit.	AI	AI	RQ	AI
	Child	<moved her hand on top of another answer>	n/a	n/a	n/a	n/a
0:21:18	Voiceover	The bird that is singing is pushing the turtle.	n/a	n/a	n/a	n/a
	Child	<attempted to make a selection>	n/a	n/a	n/a	n/a
	Parent	"The bird that is singing is pushing a turtle."	RQ	RQ	RQ	RQ
	Parent	哪个 bird 在 push turtle? 哪个 bird 在 singing?	AI	AI	RQ	AI
	Child	<moved her hand on top of another answer>	n/a	n/a	n/a	n/a
	Parent	" 哎对了。"	JC	JC	JC	JC

Figure 6: Sample annotation 1 from the In-person dataset (English item). N/A is used for transcript content that did not have a behavioral coding.

Timestamp	Speaker	Test Item	Annotator (YD)	Annotator (YY)	GPT4-FS	Final Agreement
20:11	Voiceover	"The fox points at the tree."	n/a	n/a	n/a	n/a
	Child	<points to answer>	n/a	n/a	n/a	n/a
	Parent	Wow! Good job!	E	E	E	E
		你要吃个 snack 吗还是你 okay 我們 keep going?	M	M	TS	M
		我們 keep going 了, 好嗎?	TS	E	E	TS
		okay, let's keep going	TS	E	E	TS
20:41	Voiceover	"The corn is between the bottles."	n/a	n/a	n/a	n/a
	Child	"That one."	n/a	n/a	n/a	n/a
	Parent	"一二還是三。"	TS	TS	TS	TS
	Child	"一。"	n/a	n/a	n/a	n/a

Figure 7: Sample annotation 2 from the Virtual dataset (English item). N/A is used for transcript content that did not have a behavioral coding.

<p>Assume you are a video analyst classifying transcribed text conversation shown in <Parent Behavior> from parents, who are supervising their bilingual children to answer language comprehension tests in English and Mandarin. A <Voiceover> is the system output sound which reads the picture on the website that the child has to choose. A <Child Behavior> is the children's behavior. You need to classify <Parent Behavior> as described in <task>. Let's think step by step.</p> <p><Task> Based on <Parent Behavior>, please determine which type of behavior it is: 'Repeat Questions', 'Answer Questions', 'Analyze Items', 'Judging', 'Encouragement', 'Technical Support', 'Broadcasting', 'Miscellaneous'. Definition for each category is shown in <Definitions>.</p> <p><Definitions> - Repeat Questions: Repeating the <Voiceover> audio before and/or during the process of a child selecting the picture on the web. - Answer Questions: Using verbal or gestural cues to suggest or select a correct answer for the child - Analyze Items: Elaborating on the critical linguistic components by labeling objects and actions, making emphasis via prosodic cues, or breaking down complex sentences from <Voiceover>. - Judging of Correctness: Verbally evaluating the child's response as correct or incorrect. - Encouragement: Showing verbal and/physical affirmation for the child to continue, saying good job/excellent to reinforce the child's selection, expressing empathy (e.g., it's okay) on struggled items. - Technical Support: Offering verbal and/or physical assistance to the child related to interacting with the website and the computer. - Broadcasting: After the child makes a selection, describing the selection via a word, a phrase, or a sentence. - Miscellaneous: Initiating and/or responding to events that redirected a child's attention, sharing personal opinions about test procedures and stimuli, or other verbal and nonverbal behaviors that were out of the child's view.</p>
--

Figure 8: The first zero-shot prompt used in Chain of Thoughts experiment.

Assume you are video analyst classifying transcribed text conversation shown in <Parent Behavior> from parents, who are supervising their bilingual children to answer language comprehension tests in English and Mandarin. A <Voiceover> is the system output sound which reads the picture on the website that the child has to choose. A <Child Behavior> is the children's behavior. You need to classify <Parent Behavior> as described in <task>. Let's think step by step.

<Task>
Based on <Parent Behavior>, please determine which type of behavior it is: 'Repeat Questions', 'Answer Questions', 'Analyze Items', 'Judging', 'Encouragement', 'Technical Support', 'Broadcasting', 'Miscellaneous'.
Definition for each category is shown in <Definitions>

<Definitions>
- Repeat Questions: Repeating the <Voiceover> audio before and/or during the process of a child selecting the picture on the web. For example,

<Voiceover>
"the black cat is drinking water"

<Parent Behavior>
"the black cat is drinking water"

<Classification>
Repeat Questions

- Answer Questions: Using verbal or gestural cues to suggest or select a correct answer for the child. For example,

<Voiceover>
"What is the cat drinking?"

<Parent Behavior>
"Drinking water."

<Classification>
Answer Questions

- Analyze Items: Elaborating on the critical linguistic components by labeling objects and actions, making emphasis via prosodic cues, or breaking down complex sentences from <Voiceover>. For example,

<Voiceover>
"the black cat is drinking water"

<Parent Behavior>
"This is the one with a black cat."

<Classification>
Analyze Items

Figure 9: The first few-shot prompt used in Chain of Thoughts experiment, part 1 of 2.

- Judging of Correctness: Verbally evaluating the child's response as "correct" or "incorrect". For example,

<Voiceover>

<Parent Behavior>
"This is not right."

<Classification>
Judging of Correctness

- Encouragement: Showing verbal and/physical affirmation for the child to continue, saying "good job/ excellent" to reinforce the child's selection, expressing empathy (e.g., "it's okay") on struggled items. For example,

<Voiceover>

<Parent Behavior>
"it's fine you are trying your best."

<Classification>
Encouragement

- Technical Support: Offering verbal and/or physical assistance to the child related to interacting with the website and the computer. For example,

<Voiceover>

<Parent Behavior>
"Select the picture to continue."

<Classification>
Technical Support

- Broadcasting: After the child makes a selection, describing the selection via a word, a phrase, or a sentence. For example,

<Voiceover>

<Parent Behavior>
"I selected the picture."

<Classification>
Broadcasting

- Miscellaneous: Initiating and/or responding to events that redirected a child's attention, sharing personal opinions about test procedures and stimuli, or other verbal and nonverbal behaviors that were out of the child's view. For example,

<Voiceover>
<Parent Behavior>
"My child needs to use the bathroom."

<Classification>
Miscellaneous

Figure 10: The first few-shot prompt used in Chain of Thoughts experiment, part 2 of 2.

Therefore, among all the categories, please respond the category name only.

Figure 11: The second prompt used in Chain of Thoughts experiment for both zero-shot and few-shot.

This is a text transcript consisting of <Parent Behavior> from bilingual parents who are supervising their bilingual children to answer language assessment questions in English and Mandarin. A <Voiceover> is the system output sound which reads the picture for the accurate answer the child has to choose. A <Child Behavior> is the children's behavior. Please classify <Parent Behavior> as described in <Task> and respond to the category name only.

<Task>
Based on <Parent Behavior>, please determine which type of behavior it is: 'Repeat Questions', 'Answer Questions', 'Analyze Items', 'Judging', 'Encouragement', 'Technical Support', 'Broadcasting', 'Miscellaneous'. Definition for each category is shown in <Definitions>.

<Definitions>
- Repeat Questions: Repeating the <Voiceover> audio before and/or during the process of a child selecting the picture on the web.
- Answer Questions: Using verbal or gestural cues to suggest or select a correct answer for the child.
- Analyze Items: Elaborating on the critical linguistic components by labeling objects and actions, making emphasis via prosodic cues, or breaking down complex sentences from <Voiceover>.
- Judging of Correctness: Verbally evaluating the child's response as correct or incorrect.
- Encouragement: Showing verbal and/physical affirmation for the child to continue, saying good job/excellent to reinforce the child's selection, expressing empathy (e.g., it's okay) on struggling items.
- Technical Support: Offering verbal and/or physical assistance to the child related to interacting with the website and the computer.
- Broadcasting: After the child makes a selection, describing the selection via a word, a phrase, or a sentence.
- Miscellaneous: Initiating and/or responding to events that redirected a child's attention, sharing personal opinions about test procedures and stimuli, or other verbal and nonverbal behaviors that were out of the child's view.

Figure 12: The full zero-shot prompt used in non role play test.

Parent behavior: <in response to child> "" <speaking towards the screen> "

GPT-4 with CoT response:
"let's analyze the provided <parent behavior> step by step in order to classify it:

1. the first part of the <parent behavior> is: ""<in response to child> '""". this is in response to the child's action of moving the parent's hand to the screen. the parent's response, """" or ""you click it yourself"", seems to encourage the child to perform the action by themselves instead of relying on the parent. this can be classified under 'encouragement', as it encourages independence in the child.
2. the second part of the <parent behavior> is: ""<speaking towards the screen> '""", which translates to ""can you say that again, please?"" directed towards the screen. this part indicates that the parent is requesting the system to repeat the <voiceover> audio. this action fits the 'repeat questions' category, as it involves asking for the question or statement from <voiceover> to be repeated. given the two different actions, if we consider the more prominent behavior from the parent's actions, the emphasis appears to be on encouraging the child to engage with the task independently. however, the explicit category from the definitions provided that matches an exact behavior is 'repeat questions', due to the direct request made to the system. <classification> repeat questions"

Figure 13: The first zero-shot prompt used in Chain of Thoughts experiment.

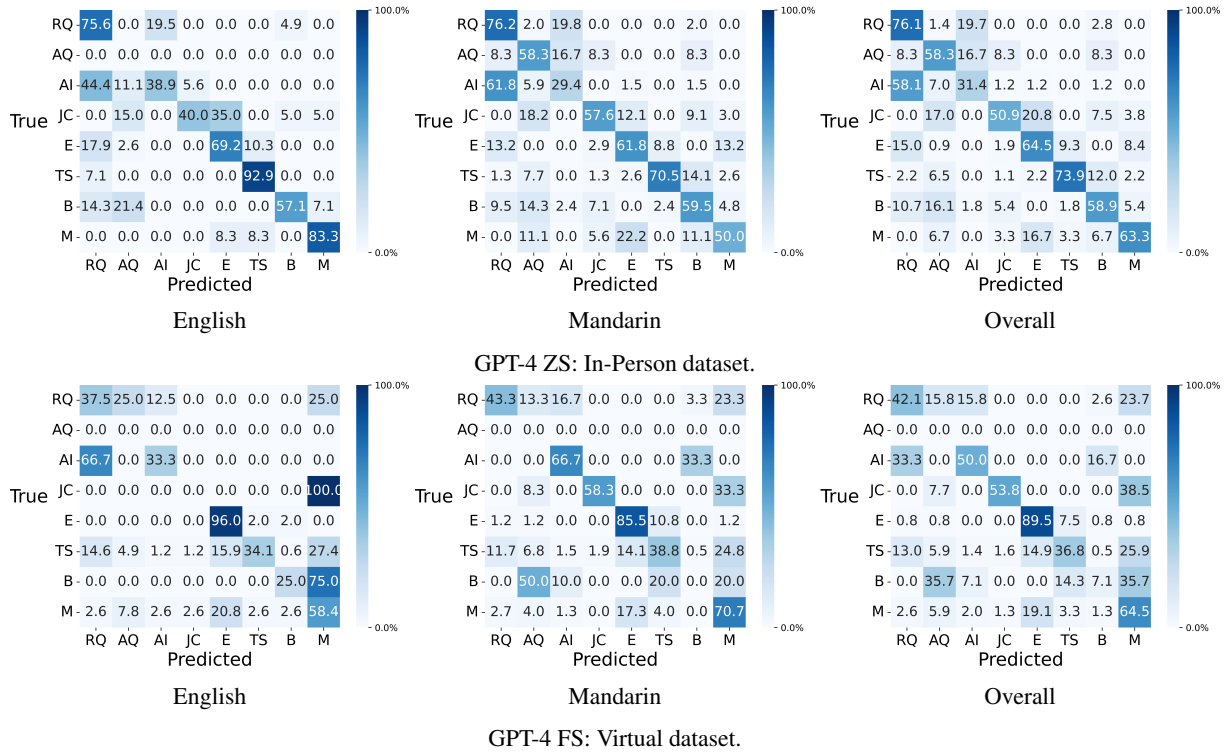


Figure 14: Normalized confusion matrices (in percentages) for GPT-4 ZS (In-Person) and GPT-4 FS (Virtual) datasets. Each row is normalized to sum to 100% within each matrix, representing the percentage distribution of predictions across classes.