

# InterventionLens: A Multi-Agent Framework for Detecting ASD Intervention Strategies in Parent-Child Shared Reading

Anonymous CVPR submission

Paper ID 15

## Abstract

001 *Home-based interventions like parent-child shared reading*  
002 *provide a cost-effective approach for supporting children*  
003 *with autism spectrum disorder (ASD). However, analyzing*  
004 *caregiver intervention strategies in naturalistic home inter-*  
005 *actions typically relies on expert annotation, which is costly,*  
006 *time-intensive, and difficult to scale. To address this chal-*  
007 *lenge, we propose InterventionLens, an end-to-end multi-*  
008 *agent system for automatically detecting and temporally*  
009 *segmenting caregiver intervention strategies from shared*  
010 *reading videos. Without task-specific model training or fine-*  
011 *tuning, InterventionLens uses a collaborative multi-agent*  
012 *architecture to integrate multimodal interaction content and*  
013 *perform fine-grained strategy analysis. Experiments on*  
014 *the ASD-HI dataset show that InterventionLens achieves*  
015 *an overall F1 score of 79.44%, outperforming the base-*  
016 *line by 19.72%. These results suggest that InterventionLens*  
017 *is a promising system for analyzing caregiver intervention*  
018 *strategies in home-based ASD shared reading settings. Ad-*  
019 *ditional resources will be released on the project page.*

## 020 1. Introduction

021 Autism Spectrum Disorder (ASD) is a prevalent neurode-  
022 velopmental disorder that substantially impairs social com-  
023 munication and daily functioning [4, 8]. While early in-  
024 tervention can improve core symptoms [7], many families  
025 lack consistent access to effective services due to long-term  
026 costs and uneven resource distribution [17]. This gap has  
027 motivated scalable, caregiver-deliverable home-based in-  
028 terventions [18, 22, 26, 27], among which shared book read-  
029 ing (SBR) is a particularly effective and accessible routine  
030 [6, 25]. In practice, families often record SBR sessions  
031 at home and submit the videos to speech-language pathol-  
032 ogists (SLPs), who manually conduct video-by-video re-  
033 view to identify caregiver intervention strategies and pro-  
034 vide guidance [2]. However, such expert-driven annota-  
035 tion is labor-intensive, costly, and difficult to scale. There-

fore, a key technical need is to automatically detect inter- 036  
vention strategies from home-based videos. In this work, 037  
we focus on three caregiver intervention strategies used in 038  
parent-implemented ASD interventions: *Modeling*, where 039  
the caregiver demonstrates the target word or phrase for the 040  
child; *Mand-Model*, where the caregiver prompts the child 041  
to respond and then provides a model when needed; and 042  
*Time Delay*, where the caregiver deliberately pauses to cre- 043  
ate an opportunity for the child to initiate or complete the re- 044  
sponse [19–21]. These strategies originate from the Parent- 045  
Implemented Communication Strategies (PiCS) framework 046  
and were later adapted to shared reading and telepractice- 047  
based interventions for children with ASD [1–3]. 048

Achieving this goal is inherently difficult because shared 049  
book reading (SBR) is highly dynamic and interactive, typ- 050  
ically involving a tightly coupled, triadic interaction among 051  
a caregiver, a child, and a book. Caregivers drive this pro- 052  
cess through continuous verbal scaffolding while simulta- 053  
neously coordinating multi-modal social cues —such as 054  
gaze, facial expressions, and gestures —to regulate child 055  
engagement [16]. Fine-grained recognition and reasoning 056  
of these interaction segments become even more challeng- 057  
ing in unstructured, in-the-wild home settings. In these en- 058  
vironments, non-fixed camera viewpoints, ambient noise, 059  
and variable interaction rhythms are compounded by the in- 060  
herent acoustic differences between adult and child voices 061  
[15]. These factors drastically degrade the performance of 062  
mainstream Automatic Speech Recognition (ASR) systems 063  
like Whisper [5], severely undermining the reliability of 064  
subsequent interaction modeling. Finally, these technical 065  
hurdles are further magnified by a severe scarcity of high- 066  
quality annotated data. Although the ASD-HI dataset [16] 067  
partially addresses the gap in home-based SBR, its limited 068  
scale cannot support the massive data demands of modern 069  
data-intensive model training. Consequently, developing a 070  
robust and generalizable automated analysis system under 071  
these compounding constraints remains a core open chal- 072  
lenge in the field. 073

To address these limitations, we propose Intervention- 074  
Lens, a hierarchical multi-agent framework for caregiver in- 075

076 intervention detection without parameter updating. Our con-  
077 tributions are three-fold:

- 078 • We propose a two-layer multi-agent architecture that de-  
079 couples multimodal perception from intervention detec-  
080 tion and segmentation, combining structured transcript  
081 and book-context modeling with coarse-to-fine expert-  
082 based temporal localization.
- 083 • We introduce a progressive knowledge refinement mech-  
084 anism that uses a small amount of labeled data to itera-  
085 tively update decision rules and refine the agent’s struc-  
086 tured guidance. By translating clinical standard operat-  
087 ing procedures (SOPs) into explicit agent roles and deci-  
088 sion constraints, our framework enables structured strat-  
089 egy analysis without data-intensive fine-tuning.
- 090 • We demonstrate strong empirical performance on the  
091 ASD-HI dataset, where InterventionLens substantially  
092 outperforms existing baselines under strict evaluation cri-  
093 teria, improving the overall F1 score from 59.72% to  
094 79.44% and Precision from 50.21% to 82.36%.

## 095 2. Related Work

096 **LLM-based Multi-Agent Systems:** Large language  
097 model (LLM)-based multi-agent systems (MAS) have  
098 emerged as an effective paradigm for complex reasoning  
099 by decomposing problems into modular subtasks handled  
100 by specialized agents [11, 23, 29]. Compared with single-  
101 agent systems, which often suffer from context dilution and  
102 error accumulation over long interactions, MAS can bet-  
103 ter isolate subproblems, reduce interference across reason-  
104 ing stages, and enable targeted optimization for heteroge-  
105 neous tasks [9, 12, 14, 24, 30]. While prior MAS research  
106 has mainly focused on general planning, coding, and task-  
107 solving settings, the use of MAS for ASD intervention strat-  
108 egy detection in shared book reading interactions remains  
109 underexplored.

110 **Challenges in Multimodal Perception.** Fine-grained  
111 analysis of parent-child shared book reading requires ac-  
112 curate speech recognition and reliable temporal grounding.  
113 However, child speech remains challenging for traditional  
114 ASR systems because of its substantial acoustic variability  
115 and pronunciation instability [10, 13]. In parallel, recent  
116 multimodal large language models have shown promising  
117 video understanding ability, but they still struggle to provide  
118 precise temporal localization in long-form videos [28, 31],  
119 especially when second-level timestamps are required for  
120 segment-level analysis. These limitations make robust mul-  
121 timodal perception a key bottleneck for automated interven-  
122 tion detection in naturalistic home environments.

123 **ASD Intervention Detection:** In the domain of ASD lan-  
124 guage intervention assessment, the ASD-HI dataset [16]

provides the main benchmark for caregiver intervention  
strategy detection in home-based shared book reading  
videos. The strongest publicly available baseline uses Whisper for speech transcription and a single LLM agent guided by prompts designed by ASD intervention experts for intervention identification, followed by greedy sequential scanning for temporal segmentation. This benchmark establishes an important foundation for the task; however, its overall detection performance still leaves substantial room for improvement because it relies on error-prone transcripts and rigid heuristic boundary search.

## 3. Problem Formulation

We formulate the task as a strategy detection problem over parent-child shared book reading videos. Given an input video, the objective is to temporally localize intervention strategies and classify each detected segment into its corresponding strategy category. Let  $V \in \mathcal{V}$  denote a parent-child shared book reading video, and let  $\mathcal{S} = \{Modeling, Mand-Model, Time Delay\}$  be the intervention strategy space. The goal is to predict a set of temporally localized intervention strategy segments:

$$\hat{\mathcal{Y}} = \{(\hat{t}_m^s, \hat{t}_m^e, \hat{s}_m)\}_{m=1}^M, \quad \hat{s}_m \in \mathcal{S}, \quad (1)$$

where  $M$  is the total number of predicted interventions. For the  $m$ -th predicted segment,  $\hat{t}_m^s$  and  $\hat{t}_m^e$  denote its predicted start and end timestamps, respectively, and  $\hat{s}_m$  represents its assigned strategy label.

InterventionLens first utilizes the OPA module to process the raw video  $V$  into grounded interaction transcript features  $\mathcal{X}$ . Subsequently, the BMA module is applied to extract the book context representation  $\mathcal{B}$  aligned with the current interaction:

$$\mathcal{X} = OPA(V), \quad \mathcal{B} = BMA(\mathcal{X}). \quad (2)$$

Next, the system employs the ICSA module to generate a set of  $K$  intervention candidates along with their routing cues based on the transcript features  $\mathcal{X}$ :

$$\mathcal{C} = \{(c_k, q_k)\}_{k=1}^K = ICSA(\mathcal{X}), \quad (3)$$

where  $c_k = (t_k^s, t_k^e)$  denotes the  $k$ -th intervention candidate segment with its initial coarse start and end times, and  $q_k \in \mathcal{Q} = \{MODEL-like, MAND-like, TD-like\}$  is the corresponding strategy cue used for expert routing.

Each candidate is then routed to the corresponding strategy-specific expert  $Expert_{q_k}$  for verification and boundary refinement. The expert first evaluates the candidate  $c_k$  against its predefined strategy criteria in conjunction with the book context  $\mathcal{B}$  to verify and assign the definitive strategy label  $\hat{s}_k \in \mathcal{S}$ . For accepted candidates, it subsequently

171 adjusts the coarse boundaries  $(t_k^s, t_k^e)$  into precise bound-  
172 aries  $(\hat{t}_k^s, \hat{t}_k^e)$ :

$$173 \hat{y}_k = \text{Expert}_{q_k}(c_k, \mathcal{B}), \quad (4)$$

174 where the output is explicitly defined as  $\hat{y}_k \in$   
175  $\{(\hat{t}_k^s, \hat{t}_k^e, \hat{s}_k), \emptyset\}$ . Here,  $\emptyset$  denotes that the candidate is re-  
176 jected (i.e., classified as background or invalid) during the  
177 strategy-specific verification.

178 The final prediction set is composed of all accepted and  
179 temporally refined candidates:

$$180 \hat{\mathcal{Y}} = \{\hat{y}_k \mid \hat{y}_k \neq \emptyset\}_{k=1}^K. \quad (5)$$

## 181 4. Method

182 As shown in Fig. 1, InterventionLens is a two-layer multi-  
183 agent framework that decouples multimodal perception and  
184 contextual grounding from downstream intervention detec-  
185 tion and boundary refinement. Specifically, the Perception  
186 Layer converts raw shared book reading videos into time-  
187 aligned multimodal interaction transcripts and infers struc-  
188 tured book context from the ongoing session. Conditioned  
189 on these structured representations, the Intervention Detec-  
190 tion and Segmentation (IDS) Layer follows a coarse-to-fine  
191 pipeline that first extracts and routes candidate intervention  
192 segments, and then performs strategy-specific verification  
193 and temporal boundary refinement.

### 194 4.1. Perception Layer

195 The Perception Layer consists of two components: the  
196 Omni-Perception Agent (OPA), which constructs time-  
197 aligned multimodal interaction transcripts, and the Book  
198 Modeling Agent (BMA), which reconstructs a structured  
199 book-context representation.

200 **Omni-Perception Agent (OPA):** The Omni-Perception  
201 Agent (OPA) constructs time-aligned multimodal interac-  
202 tion transcripts from caregiver-child Shared Book Reading  
203 (SBR) videos for downstream intervention analysis. It cap-  
204 tures verbal exchanges between caregivers and children to-  
205 gether with observable actions and facial expressions from  
206 both participants. We use *gemini-3-flash-preview* as the  
207 backbone multimodal audio-video understanding model to  
208 parse dialog and visible interaction cues. For temporal  
209 grounding, we follow an OCR-assisted timestamping strat-  
210 egy in which a red timestamp is overlaid in the bottom-right  
211 corner of each frame and used as an explicit temporal ref-  
212 erence to predict start and end times [28]. For long record-  
213 ings, each reading session is processed as a sequence of 15-  
214 second clips. In parallel, the audio stream is transcribed  
215 using *Whisper*. We then align Gemini-generated caregiver  
216 utterances with Whisper transcriptions under semantic and  
217 temporal consistency constraints. When a Gemini-parsed  
218 caregiver utterance has a semantically matched Whisper

counterpart with closely aligned timestamps, we replace  
the Gemini-predicted temporal boundaries with the cor-  
responding Whisper timestamps. The resulting transcript  
serves as perceptual input to the downstream BMA and IDS  
Layer.

**Book Modeling Agent (BMA):** Shared book reading  
forms a triadic interaction among the caregiver, the child,  
and the book. As a result, caregiver interventions are of-  
ten grounded in the specific context of the book being read  
rather than in dialogue alone. To incorporate this contex-  
tual information, we introduce the Book Modeling Agent  
(BMA), which infers a Structured Book-Context Represen-  
tation from the parent-child interaction transcript. Instead  
of reconstructing the full book content, BMA extracts con-  
textual cues such as recurring sentence patterns and poten-  
tial target words that are relevant to the current reading in-  
teraction. This inferred book context is then used as se-  
mantic grounding for downstream expert modules, enabling  
more accurate intervention detection by aligning caregiver  
utterances with the underlying reading context.

### 4.2. Intervention Detection & Segmentation Layer

The Intervention Detection and Segmentation Layer detects  
and temporally segments caregiver intervention strategies  
by leveraging the enriched interaction transcripts and re-  
constructed book context produced by the Perception Layer.  
It is implemented through task-specific agents whose deci-  
sion logic is derived from the ASD-HI coding manual and  
refined on the training split. Rather than training model pa-  
rameters, we translate the operational definitions of Mod-  
eling, Mand-Model, and Time Delay into reusable decision  
rules.

**Intervention Candidate Scanner Agent (ICSA)** ICSA  
is a module for candidate intervention extraction. Based on  
the caregiver-child interaction transcripts produced by the  
Perception Layer, it identifies interaction segments that may  
constitute intervention events. Specifically, ICSA extracts  
two clinically meaningful candidate patterns: Complete In-  
tervention Loops, which follow a caregiver stimulus, child  
response, and caregiver reinforcement sequence, and Inter-  
vention Attempts, in which a caregiver provides a stimulus  
but the child does not respond. Based on these candidates,  
ICSA further assigns a coarse-grained strategy cue, such as  
MODEL-like, MAND-like, or TD-like, and uses it as a rout-  
ing signal to dispatch each candidate to the corresponding  
strategy-specific expert for subsequent fine-grained verifi-  
cation and boundary refinement.

**Strategy-Specific Experts (SSEs)** The Strategy-Specific  
Experts (SSEs) perform strategy-specific verification and

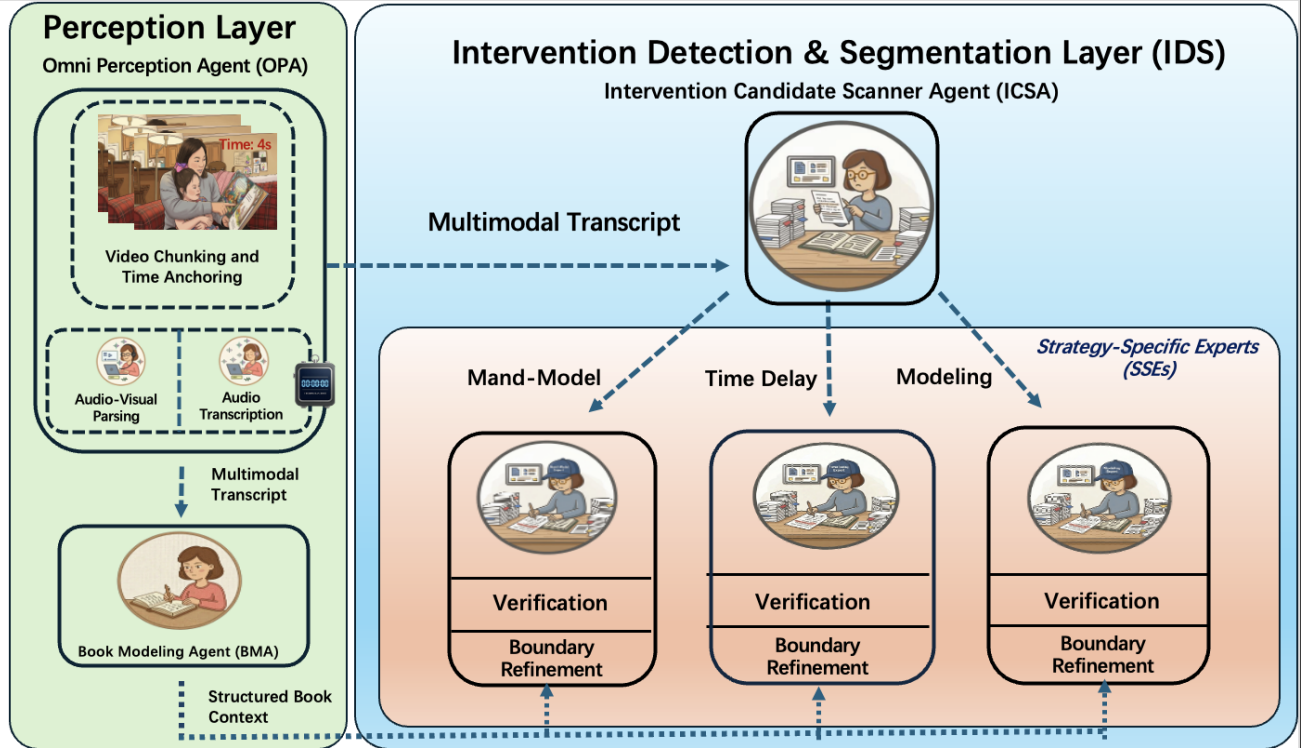


Figure 1. Overview of InterventionLens. The framework consists of a Perception Layer and an Intervention Detection & Segmentation Layer (IDS). The Perception Layer first converts the raw shared book reading video into a time-aligned multimodal transcript and structured book context. The IDS Layer then proposes candidate intervention segments and routes them to strategy-specific experts, each of which performs strategy-specific verification and boundary refinement to produce the final temporally localized intervention segments.

267 temporal boundary refinement for the candidate interven- 290  
 268 tion segments routed from the Intervention Candidate Scanner 291  
 269 Agent (ICSA). To account for the substantial differences 292  
 270 in linguistic patterns and interaction structures among Mod- 293  
 271 eling, Mand-Model, and Time Delay, we design SSE as a 294  
 272 parallel strategy-specific expert architecture. This architec- 295  
 273 ture consists of three independent expert systems, each ded- 296  
 274 icated to one intervention strategy. All three expert systems 297  
 275 follow the same two-stage cascade paradigm, where Stage 1 298  
 276 performs strategy-specific verification and Stage 2 refines 299  
 277 temporal boundaries. 300

278 *Strategy-Specific Verification:* Strategy-Specific Verifica- 301  
 279 tion is the core logical filtering stage for candidate interven- 302  
 280 tion screening. Because the candidate segments proposed 303  
 281 by ICSA still contain noisy or weak matches, this stage ver- 304  
 282 ifies each candidate using strategy-specific constraints indu- 305  
 283 ced from the ASD-HI dataset together with book-context 306  
 284 cues from the Book Modeling Agent (BMA). The verifica- 307  
 285 tion considers utterance form, interaction structure, local 308  
 286 discourse context, and book-related patterns. Only candi- 309  
 287 dates consistent with the target strategy are retained for sub- 310  
 288 sequent boundary refinement. 311

289 *Temporal Boundary Refinement:* After strategy verifica- 312

tion, the retained candidates proceed to the temporal bound-  
 ary refinement stage, where their temporal extent is local-  
 ized more precisely. Guided by local dialogue flow and ex-  
 pert annotation conventions, this stage refines the start and  
 end boundaries of each intervention segment. Its goal is to  
 narrow the relatively broad candidate region from the pre-  
 vious stage into a precise intervention interval by removing  
 adjacent turns before and after the core intervention that do  
 not belong to the intervention itself. As a result, the refined  
 segments are more closely aligned with expert annotation  
 standards and more accurately capture the actual temporal  
 span of the intervention.

### 4.3. Progressive Knowledge Refinement for SSEs

We initialize the agent’s structured guidance using the  
 ASD-HI annotation protocol, which specifies strategy def-  
 initions, temporal boundary judgment, and key interaction  
 cues. We then iteratively refine the guidance on the training  
 split with F1 score as the optimization objective. In each  
 iteration, we run the full system, analyze recurrent failure  
 cases, and update the guidance rules based on the most fre-  
 quent error patterns. We stop the refinement process when  
 the F1 score does not improve for three consecutive itera-  
 tions.

313 **5. Experiments**

314 This section presents the experimental evaluation of Inter-  
 315 ventionLens on the ASD-HI benchmark. To ensure direct  
 316 comparability, we follow the same evaluation protocol as  
 317 prior work. We first introduce the dataset, evaluation met-  
 318 rics, and implementation details. We then compare Inter-  
 319 ventionLens with the strongest reported baseline and con-  
 320 duct an ablation on the Book Modeling Agent (BMA). Fi-  
 321 nally, we provide additional analysis on the unguided pre-  
 322 feedback subset, where caregivers had not yet received ex-  
 323 pert feedback, and discuss the main error patterns observed  
 324 in the benchmark.

325 **5.1. Experimental Setup**

326 **Dataset** We evaluate InterventionLens on the ASD-HI  
 327 [16] benchmark, a multimodal parent-child shared book  
 328 reading dataset collected in home-based intervention set-  
 329 tings for families of children with ASD. The dataset records  
 330 natural caregiver child interactions in real home environ-  
 331 ments. It contains 48 complete shared-reading sessions  
 332 from 3 families, together with expert-provided fine grained  
 333 annotations of caregiver intervention segments. For each in-  
 334 tervention segment, the dataset provides not only the corre-  
 335 sponding intervention strategy label and its temporal bound-  
 336 aries, but also fine-grained annotations of the caregiver-  
 337 child interaction process within the segment, including di-  
 338 alogue content, behavioral actions, and other interaction  
 339 cues. In addition, the samples in the dataset can be fur-  
 340 ther divided into two groups: those collected before care-  
 341 givers received expert feedback, and those collected after  
 342 caregivers received expert feedback.

343 Due to variations in home environments and recording  
 344 conditions, camera placement is not consistent across fam-  
 345 ilies. Moreover, caregivers differ in prompting style, and  
 346 children vary in engagement and responsiveness, result-  
 347 ing in substantial cross-family variability and distribution  
 348 shift. The dataset provides 239 annotated segments for train-  
 349 ing and 120 annotated segments for evaluation. We note that  
 350 3 segments in the evaluation split contain incomplete anno-  
 351 tations, specifically missing strategy labels or missing tem-  
 352 poral boundaries. We therefore conduct our final evaluation  
 353 on 117 valid test segments. All reported results follow the  
 354 same evaluation split used by the original dataset authors.

355 **Evaluation Protocol** We employ Precision, Recall, and  
 356 F1-score as the primary evaluation metrics. A predicted in-  
 357 tervention is counted as a true positive only if it simultane-  
 358 ously satisfies three core criteria: (1) Strategy Correctness,  
 359 where the predicted strategy category must exactly match  
 360 the expert-annotated ground truth; (2) Strategy Complete-  
 361 ness, where the prediction must capture the complete care-  
 362 giver child interaction process underlying the intervention,

Table 1. Overall results on ASD-HI

Method	Prec.	Rec.	F1
Baseline	50.21	73.68	59.72
Ours	82.36±3.44	76.75±0.49	79.44±1.84
Gain	+32.15	+3.07	+19.72

including all necessary interactional details, rather than only  
 a fragmented or isolated portion of it; and (3) Temporal Ac-  
 curacy, where the predicted time interval must align with  
 the annotated reference boundaries within a reasonable tol-  
 erance window of 1.0 second. To assess the stability of  
 model generation, we conduct three independent runs on  
 the same evaluation split. For each run, we compute the  
 results for the three strategy categories Mand-Model, Mod-  
 eling, and Time Delay and report the macro-averaged per-  
 formance. The final results are presented as the mean  $\pm$   
 standard deviation across the three runs.

**Implementation Details** In the perception layer, we  
 use *Whisper-1* for audio transcription and *text-embedding-*  
*3-small* to support semantic alignment between Gemini  
 parsed caregiver utterances and Whisper transcription seg-  
 ments. Specifically, alignment is performed by computing  
 the cosine distance in the embedding space. We replace  
 the Gemini-predicted temporal boundaries with the corre-  
 sponding Whisper boundaries only when the cosine dis-  
 tance is below 0.1 and the boundary discrepancy is within  
 1.0 second; otherwise, the original Gemini timestamps are  
 retained. For multimodal understanding and downstream  
 reasoning, all agent components use *gemini-3-flash-preview*  
 as the backbone model, with the temperature set to 1.0.

During the Progressive Knowledge Refinement for SSEs  
 stage, all rule induction and prompt refinement are per-  
 formed on the subset collected after caregivers received ex-  
 pert feedback, consisting of 197 annotated segments. We  
 use this subset because intervention behaviors at this stage  
 are more standardized, with more stable strategy boundaries  
 and interaction patterns, which makes it better suited for ex-  
 tracting reusable decision rules. In contrast, the interactions  
 collected before expert feedback are more natural and vari-  
 able, and are used in our subsequent analysis. We conduct  
 nine rounds of refinement and stop the process when the  
 F1 score fails to improve for three consecutive iterations,  
 yielding a final training-set F1 score of 0.774.

362 **5.2. Main Results**

We first compare InterventionLens with the baseline re-  
 ported in the original ASD-HI paper under the same eval-  
 uation protocol. Since the original ASD-HI paper reports  
 only overall performance without strategy-wise Precision,

Table 2. Strategy wise results on ASD-HI.

Strategy	Prec.	Rec.	F1
Mand-Model	91.87±1.19	87.88±0.75	89.82±0.31
Modeling	57.34±12.72	26.67±6.67	36.28±8.25
Time Delay	65.03±6.44	72.84±5.66	68.65±5.63

Note: The evaluation contains 117 annotated intervention strategies, including 74 *Mand-Model*, 27 *Time Delay*, and 16 *Modeling* instances.

405 Recall, or F1 scores, we restrict our comparison to the over-  
406 all metrics. As shown in Table 1, InterventionLens sub-  
407 stantially improves overall performance over the reported  
408 baseline [16]. Specifically, our method increases Precision  
409 from 50.21 to 82.36, Recall from 73.68 to 76.75, and F1  
410 from 59.72 to 79.44. The most notable gain comes from the  
411 large improvement in **Precision (+32.15)**, while **Recall** is  
412 also slightly improved (**+3.07**). Consequently, the overall  
413 **F1** score increases by (**+19.72**), indicating a substantially  
414 better balance between prediction accuracy and coverage.  
415 This suggests that InterventionLens does not achieve better  
416 performance simply by producing more aggressive predic-  
417 tions; instead, it substantially reduces false positives while  
418 maintaining strong coverage of true intervention segments.

419 Table 2 further reports the per-strategy performance of  
420 InterventionLens. Our method achieves the strongest results  
421 on *Mand-Model*, with an F1 of 89.82, suggesting that this  
422 strategy can be detected reliably under our framework. Per-  
423 formance on *Time Delay* remains competitive, reaching an  
424 F1 of 68.65. In contrast, *Modeling* is the most challenging  
425 category, with an F1 of 36.28. This gap suggests that differ-  
426 ent intervention strategies pose substantially different levels  
427 of difficulty, and that *Modeling* remains the major bottle-  
428 neck in the current benchmark.

429 In addition, as a supplementary analysis, we further ex-  
430 amine the performance of InterventionLens on the unguided  
431 interaction subset collected before caregivers received ex-  
432 pert feedback. On this subset, the system achieves an over-  
433 all F1 of 85.09±2.24, suggesting that the rules induced from  
434 post-feedback data remain applicable to more natural fam-  
435 ily shared-reading interactions. However, this result should  
436 not be over interpreted, since the evaluation subset contains  
437 only 24 annotated segments and is highly imbalanced across  
438 strategy categories, with only 15, 6, 3 instances for Mand-  
439 Model, Time Delay, and Modeling, respectively. We there-  
440 fore treat this result as a supplementary observation on sys-  
441 tem applicability, rather than as a definitive validation of  
442 cross-distribution generalization.

### 443 5.3. Ablation Study

444 **Ablation on the Book Modeling Agent (BMA).** To evalu-  
445 ate the contribution of book-conditioned semantic ground-  
446 ing, we remove the Book Modeling Agent (BMA) while

Table 3. Overall results under BMA ablation.

Method	Prec.	Rec.	F1
InterventionLens	82.36±3.44	76.75±0.49	79.44±1.84
w/o BMA	70.22±3.49	67.79±1.28	68.97±2.35
Drop	-12.14	-8.96	-10.47

Table 4. Strategy-wise F1 under BMA ablation.

Strategy	InterventionLens	w/o BMA	$\Delta$ F1
Mand-Model	89.82±0.31	80.27±2.90	-9.55
Modeling	36.28±8.25	16.96±9.06	-19.32
Time Delay	68.65±5.63	64.93±2.17	-3.72

keeping all other components and experimental settings un- 447  
changed. As shown in Table 3, removing BMA leads to a 448  
clear overall performance drop: Precision decreases from 449  
82.36 to 70.22, Recall decreases from 76.75 to 67.79, and 450  
F1 decreases from 79.44 to 68.97. This result indicates that 451  
explicit book modeling is a key component of the full sys- 452  
tem rather than an optional auxiliary module. 453

454 Table 4 further shows that the impact of removing BMA 455  
is not uniform across strategy categories. In particular, the 456  
F1 of Modeling drops from 36.28 to 16.96, while Mand- 457  
Model also decreases from 89.82 to 80.27. In contrast, 458  
Time Delay is relatively less affected, with F1 decreasing 459  
from 68.65 to 64.93. These results suggest that BMA con- 460  
tributes broadly to intervention strategy detection, and may 461  
be particularly important for strategies that require aligning 462  
caregiver utterances with book content and target-word se- 463  
mantics. We note that the Modeling category has relatively 464  
few instances in the evaluation split, and its category-level 465  
results should therefore be interpreted with caution; how- 466  
ever, the clear degradation in both the overall results and 467  
Mand-Model still consistently indicates that explicit book 468  
modeling provides critical contextual grounding for the sys- 469  
tem. Overall, this ablation supports our central assumption 470  
that shared book reading should be modeled as a triadic in- 471  
teraction among the caregiver, the child, and the book.

### 5.4. Error Analysis 472

A notable observation is that the detection performance of 473  
*Modeling* is substantially lower than that of the other two 474  
strategy categories. We believe that this result is influenced 475  
by at least two factors. First, the number of *Modeling* in- 476  
stances in the evaluation split is very small, with only 16 477  
samples, making the corresponding metrics more sensitive 478  
to fluctuations caused by a few individual cases. Second, 479  
we observe that there may be some inconsistency in la- 480  
beling criteria between the training and evaluation splits. 481  
Specifically, when the same caregiver reads the same story 482

483 book, *Brown Bear, Brown Bear, What Do You See?*, utter- 533  
484 ances that are highly similar in both form and semantics are 534  
485 not always annotated consistently across splits. For exam- 535  
486 ple, in the training split, repeated target-word expressions 536  
487 such as “black sheep . . . black sheep” are labeled as *Model-* 537  
488 *ing*, whereas in the evaluation split, similarly structured ex- 538  
489 pressions from the same caregiver, such as “red bird . . . red 539  
490 bird”, are not labeled as *Modeling*. This suggests that the 540  
491 *Modeling* category is challenged not only by limited sample 541  
492 size, but also potentially by ambiguous labeling boundaries 542  
493 or inconsistent annotation criteria. 543

## 494 6. Conclusion and Limitations 544

495 In this paper, we presented InterventionLens, a hierarchi- 545  
496 cal multi-agent framework for detecting and temporally 546  
497 segmenting caregiver intervention strategies in home-based 547  
498 shared reading videos. Unlike approaches that rely on task- 548  
499 specific parameter updating, we decompose the task into a 549  
500 sequence of interconnected subprocesses, including mul- 550  
501 timodal perception, book-context modeling, intervention 551  
502 candidate discovery, strategy-specific verification, and tem- 552  
503 poral boundary refinement. This task decomposition allows 553  
504 the system to organize the overall reasoning pipeline from 554  
505 perception to decision-making in a modular manner, while 555  
506 enabling fine-grained temporal localization of intervention 556  
507 segments on the current benchmark. 557

508 Experiments on the ASD-HI benchmark show that In- 558  
509 terventionLens substantially outperforms the strongest re- 559  
510 ported baseline, achieving clear gains in overall Precision 560  
511 and F1. Further ablation results demonstrate that explicit 561  
512 book-conditioned semantic grounding is a key component 562  
513 of the full system rather than an optional auxiliary design. 563  
514 Modeling shared reading as a triadic interaction among the 564  
515 caregiver, the child, and the book provides important con- 565  
516 textual support for intervention detection, and is particularly 566  
517 helpful for strategies that depend on aligning caregiver ut- 567  
518 terances with recurring story patterns and target-word se- 568  
519 mantics. 569

520 It should also be noted that InterventionLens was devel- 570  
521 oped based on the ASD-HI training set and its correspond- 571  
522 ing task formulation. As a result, the method may still ex- 572  
523 hibit a certain degree of dependence on the benchmark’s 573  
524 annotation conventions, strategy boundaries, and data dis- 574  
525 tribution. Although the experimental results show that 575  
526 InterventionLens achieves substantial improvements under 576  
527 the current evaluation protocol, its transferability and ro- 577  
528 bustness across datasets, family settings, and annotation 578  
529 schemes still require more systematic validation. 579

## 530 7. Ethical Considerations 580

531 This study utilizes the ASD-HI benchmark dataset, which 581  
532 was collected through a prior Institutional Review Board 582

(IRB)-approved study. The current research does not in- 533  
volve any new human-subject data collection; the released 534  
dataset is used strictly for academic research purposes in 535  
line with the approved data use guidelines. 536

537 For multimodal analysis, we accessed the Gemini model 537  
538 through Google Cloud Vertex AI. Customer data is not used 538  
539 to train Google’s AI models as stated in their terms of use, 539  
540 and all API communications are protected via encryption in 540  
541 transit using TLS. Furthermore, we have strictly limited the 541  
542 use of data to research analysis, with no attempts to iden- 542  
543 tify or re-identify participants, and no redistribution of raw 543  
544 videos or identifiable materials. 544

## 545 References 545

- 546 [1] Yusuf Akemoglu and Kimberly R Tomeny. A parent- 546  
547 implemented shared-reading intervention to promote com- 547  
548 munication skills of preschoolers with autism spectrum dis- 548  
549 order. *Journal of Autism and Developmental Disorders*, 51 549  
550 (8):2974–2987, 2021. 1 550
- 551 [2] Yusuf Akemoglu, Vanessa Hinton, Dayna Laroue, and 551  
552 Vanessa Jefferson. A parent-implemented shared reading in- 552  
553 tervention via telepractice. *Journal of Early Intervention*, 44 553  
554 (2):190–210, 2022. 1 554
- 555 [3] Yusuf Akemoğlu, Dayna Laroue, Carolina Kudesey, and 555  
556 Mary Stahlman. A module-based telepractice intervention 556  
557 for parents of children with developmental disabilities. *Jour- 557  
558 nal of Autism and Developmental Disorders*, 52(12):5177– 558  
559 5190, 2022. 1 559
- 560 [4] American Psychiatric Association et al. *Diagnostic and sta- 560  
561 tistical manual of mental disorders*. American psychiatric 561  
562 association, 2013. 1 562
- 563 [5] Ahmed Adel Attia, Jing Liu, Wei Ai, Dorottya Demszky, and 563  
564 Carol Espy-Wilson. Kid-whisper: Towards bridging the per- 564  
565 formance gap in automatic speech recognition for children 565  
566 vs. adults. In *Proceedings of the AAAI/ACM Conference on 566  
567 AI, Ethics, and Society*, pages 74–80, 2024. 1 567
- 568 [6] Susannah A Boyle, David McNaughton, and Shelley E 568  
569 Chapin. Effects of shared reading on the early language and 569  
570 literacy skills of children with autism spectrum disorders: A 570  
571 systematic review. *Focus on Autism and Other Developmen- 571  
572 tal Disabilities*, 34(4):205–214, 2019. 1 572
- 573 [7] Annette Estes, Jeffrey Munson, Sally J Rogers, Jessica 573  
574 Greenson, Jamie Winter, and Geraldine Dawson. Long-term 574  
575 outcomes of early intervention in 6-year-old children with 575  
576 autism spectrum disorder. *Journal of the American Academy 576  
577 of Child & Adolescent Psychiatry*, 54(7):580–587, 2015. 1 577
- 578 [8] Elizabeth A Fuller and Ann P Kaiser. The effects of early 578  
579 intervention on social communication outcomes for children 579  
580 with autism spectrum disorder: A meta-analysis. *Journal 580  
581 of autism and developmental disorders*, 50(5):1683–1700, 581  
582 2020. 1 582
- 583 [9] Mingyan Gao, Yanzi Li, Banruo Liu, Yifan Yu, Phillip Wang, 583  
584 Ching-Yu Lin, and Fan Lai. Single-agent or multi-agent sys- 584  
585 tems? why not both? *arXiv preprint arXiv:2505.18286*, 585  
586 2025. 2 586

- 587 [10] Matteo Gerosa, Diego Giuliani, and Fabio Brugnara. Acous- 644  
588 tic variability and automatic recognition of children’s speech. 645  
589 *Speech communication*, 49(10-11):847–860, 2007. 2 646
- 590 [11] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, 647  
591 Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xian- 648  
592 gliang Zhang. Large language model based multi-agents: 649  
593 A survey of progress and challenges. *arXiv preprint* 650  
594 *arXiv:2402.01680*, 2024. 2 651
- 595 [12] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, 652  
596 Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, 653  
597 Steven Ka Shing Yau, Zijuan Lin, et al. Metagpt: Meta pro- 654  
598 gramming for a multi-agent collaborative framework. In *The* 655  
599 *twelfth international conference on learning representations*, 656  
600 2023. 2 657
- 601 [13] Koharu Horii, Naohiro Tawara, Atsunori Ogawa, and Shoko 658  
602 Araki. Why is children’s asr so difficult? analyzing chil- 659  
603 dren’s phonological error patterns using ssl-based phoneme 660  
604 recognizers. In *Proc. Interspeech*, 2025. 2 661
- 605 [14] Jiangping Huang, Wenguang Ye, Weisong Sun, Jian Zhang, 662  
606 Mingyue Zhang, and Yang Liu. Tracecoder: A trace-driven 663  
607 multi-agent framework for automated debugging of llm- 664  
608 generated code. *arXiv preprint arXiv:2602.06875*, 2026. 2 665
- 609 [15] Rishabh Jain, Andrei Barcovschi, Mariam Yiwere, Peter 666  
610 Corcoran, and Horia Cucu. Adaptation of whisper models to 667  
611 child speech recognition. *arXiv preprint arXiv:2307.13008*, 668  
612 2023. 1 669
- 613 [16] Zhaohui Li, Yusuf Akemoglu, Jincheng Lyu, Qingxiao 670  
614 Zheng, and Jinjun Xiong. Asd-hi: A parent-child interac- 671  
615 tion dataset for automated assessment of home intervention. 672  
616 In *International Conference on Artificial Intelligence in Ed- 673*  
617 *ucation*, pages 48–62. Springer, 2025. 1, 2, 5, 6 674
- 618 [17] Yingying Lin, Guozhi Chen, Huaxiang Lu, Rongfei Qin, 675  
619 Jinsheng Jiang, Weiwei Tan, Caibin Luo, Ming Chen, Qin 676  
620 Huang, Liangliang Huang, et al. Inequality and heterogene- 677  
621 ity in medical resources for children with autism spectrum 678  
622 disorders: a study in the ethnic minority region of southern 679  
623 china. *BMC Public Health*, 25(1):1677, 2025. 1 680
- 624 [18] Catherine Lord, Tony Charman, Alexandra Havdahl, Paul 681  
625 Carbone, Evdokia Anagnostou, Brian Boyd, Themba Carr, 682  
626 Petrus J De Vries, Cheryl Dissanayake, Gauri Divan, et al. 683  
627 The lancet commission on the future of care and clinical re- 684  
628 search in autism. *The Lancet*, 399(10321):271–334, 2022. 685  
629 1 686
- 630 [19] Hedda Meadan, Lori E Meyer, Melinda R Snodgrass, and 687  
631 James W Halle. Coaching parents of young children with 688  
632 autism in rural areas using internet-based technologies: A 689  
633 pilot program. *Rural Special Education Quarterly*, 32(3): 690  
634 3–10, 2013. 1 691
- 635 [20] Hedda Meadan, Maureen E Angell, Julia B Stoner, and 692  
636 Marcus E Daczewitz. Parent-implemented social-pragmatic 693  
637 communication intervention: A pilot study. *Focus on Autism* 694  
638 *and Other Developmental Disabilities*, 29(2):95–110, 2014. 695
- 639 [21] Hedda Meadan, Melinda R Snodgrass, Lori E Meyer, Kim W 696  
640 Fisher, Moon Y Chung, and James W Halle. Internet- 697  
641 based parent-implemented intervention for young children 698  
642 with autism: A pilot study. *Journal of Early Intervention*, 699  
643 38(1):3–23, 2016. 1 700
- [22] World Health Organization et al. Training for caregivers of 701  
children with developmental disabilities, including autism. 702  
*Retrieved January, 17:2023*, 2022. 1 703
- [23] Khanh-Tung Tran, Dung Dao, Minh-Duong Nguyen, Quoc- 704  
Viet Pham, Barry O’Sullivan, and Hoang D Nguyen. Multi- 705  
agent collaboration mechanisms: A survey of llms, 2025. 706  
*URL https://arxiv.org/abs/2501.06322*, 2025. 2 707
- [24] Xiao Wang, Lu Dong, Sahana Rangasrinivasan, Ifeoma 708  
Nwogu, Srirangaraj Setlur, and Venugopal Govindaraju. Au- 709  
tomisty: a multi-agent llm framework for automated code 710  
generation in the misty social robot. In *2025 IEEE/RSJ In-* 711  
*ternational Conference on Intelligent Robots and Systems* 712  
*(IROS)*, pages 9194–9201. IEEE, 2025. 2 713
- [25] Marleen F Westerveld, Rachele Wicks, and Jessica Paynter. 714  
Investigating the effectiveness of parent-implemented shared 715  
book reading intervention for preschoolers with asd. *Child* 716  
*Language Teaching and Therapy*, 37(2):149–162, 2021. 1 717
- [26] Amy M Wetherby, Juliann Woods, Whitney Guthrie, Abi- 718  
gail Delehanty, Jennifer A Brown, Lindee Morgan, Renee D 719  
Holland, Christopher Schatschneider, and Catherine Lord. 720  
Changing developmental trajectories of toddlers with autism 721  
spectrum disorder: Strategies for bridging research to com- 722  
munity practice. *Journal of Speech, Language, and Hearing* 723  
*Research*, 61(11):2615–2628, 2018. 1 724
- [27] Paul Wai-Ching Wong, Yan-Yin Lam, Janet Siu-Ping Lau, 725  
Hung-Kit Fok, and WHO CST Team Servili Chiara 4 Sa- 726  
lomone Erica 5 6 Pacione Laura 4 5 Shire Stephanie 6 Brown 727  
Felicity 7 8. Adapting and pretesting the world health orga- 728  
nization’s caregiver skills training program for children with 729  
autism and developmental disorders or delays in hong kong. 730  
*Scientific Reports*, 12(1):16932, 2022. 1 731
- [28] Yongliang Wu, Xinting Hu, Yuyang Sun, Yizhou Zhou, 732  
Wenbo Zhu, Fengyun Rao, Bernt Schiele, and Xu Yang. 733  
Number it: Temporal grounding videos like flipping manga. 734  
In *Proceedings of the Computer Vision and Pattern Recogni-* 735  
*tion Conference*, pages 13754–13765, 2025. 2, 3 736
- [29] Yingxuan Yang, Qiuying Peng, Jun Wang, Ying Wen, and 737  
Weinan Zhang. Llm-based multi-agent systems: Techniques 738  
and business perspectives. *arXiv preprint arXiv:2411.14033*, 739  
2024. 2 740
- [30] Yusen Zhang, Ruoxi Sun, Yanfei Chen, Tomas Pfister, Rui 741  
Zhang, and Sercan Arik. Chain of agents: Large lan- 742  
guage models collaborating on long-context tasks. *Advances* 743  
*in Neural Information Processing Systems*, 37:132208– 744  
132237, 2024. 2 745
- [31] Zijia Zhao, Haoyu Lu, Yuqi Huo, Yifan Du, Tongtian Yue, 746  
Longteng Guo, Bingning Wang, Weipeng Chen, and Jing 747  
Liu. Needle in a video haystack: A scalable synthetic evalua- 748  
tor for video mllms. *arXiv preprint arXiv:2406.09367*, 2024. 749  
2 750