

SIV-Bench: A Video Benchmark for Social Interaction Understanding and Reasoning

Anonymous ACL submission

Abstract

Understanding social interaction, which encompasses perceiving numerous and subtle multimodal cues, inferring unobservable mental states and relations, and dynamically predicting others' behavior, is the foundation for achieving human-machine interaction. Despite rapid advances in Multimodal Large Language Models (MLLMs), the rich and multifaceted nature of social interaction has hindered the development of benchmarks that holistically evaluate and guide their social interaction abilities. Based on social relation theory, which has been widely regarded as a foundational framework for understanding social behavior, we provide SIV-Bench, a novel video benchmark for systematically evaluating MLLMs' capabilities across Social Scene Understanding (SSU), Social State Reasoning (SSR), and Social Dynamics Prediction (SDP). SIV-Bench features 2,792 originally collected video clips and 5,455 meticulously generated question-answer pairs derived from a human-LLM collaborative pipeline. It covers 14 typical relationships, diverse video lengths, genres, presentation styles, and linguistic and cultural backgrounds. Our comprehensive experiments show that leading MLLMs perform relatively well on SSU but remain weak on SSR and SDP, with the systematic confusion in relation inference as a key bottleneck. An in-depth analysis of the reasoning process attributes MLLMs' sub-optimal performance to misalignment with human thoughts and insufficient reasoning depth. Moreover, we find audio and subtitles aid in reasoning-intensive SSR and SDP. Together, SIV-Bench offers a unified testbed to measure progress, expose limitations, and guide future research toward more socially intelligent MLLMs.

1 Introduction

The rapid development of Multimodal Large Language Models (MLLMs) capable of processing text,

images, and video has driven strong performance across tasks such as visual reasoning, video captioning, and multimodal dialogue (Team et al., 2023; Wu et al., 2024; Zhu et al., 2025b; Hurst et al., 2024; Yang et al., 2024). As these capabilities expand, there is a growing need for benchmarks that can evaluate model performance, uncover limitations, and guide future research (Fu et al., 2024; Fang et al., 2024; Zhou et al., 2024; Qiang et al., 2025; Huang et al., 2024). One critical yet underexplored area is the social interaction understanding and reasoning, which is a core aspect of social intelligence that encompasses not only observable behaviors but also implicit mental states and social relationships governing behaviors such as forming bonds, exchanging information, and coordinating actions (Berger et al., 1972; Smith-Lovin and Heise, 1988). However, existing video benchmarks, whether designed for specific tasks like video object segmentation (Ding et al., 2025; Athar et al., 2025), captioning (Chen et al., 2025; Wu et al., 2025), and fine-grained action understanding (Perrett et al., 2025), or for broader video understanding (Wang et al., 2025; Li et al., 2024b,c), still struggle to systematically probe MLLMs' understanding of the multifaceted nature of social interaction.

To address this gap, we *first* decompose the capacity to understand and reason about social interaction into three core, interrelated dimensions. **1) Social Scene Understanding (SSU)** is foundational, enabling the recognition of visible elements such as objects, environments, and socially salient human features like body movements, clothing, and physical appearance. Reliable scene perception is required to ground interpretations in relevant cues. **2) Social State Reasoning (SSR)** is essential for interpreting the unobservable states of interaction, such as emotions, intents, attitudes, and interpersonal relationships, which guide and shape behavior (Strachan et al., 2024; Wu et al., 2020). This capacity allows models to move beyond surface-level

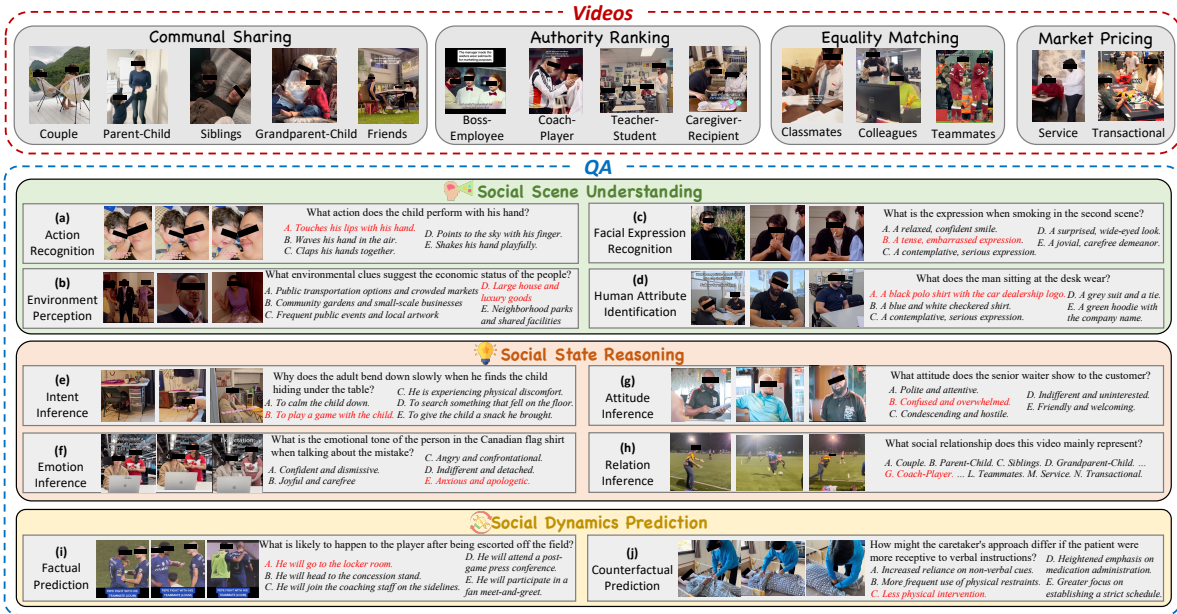


Figure 1: Overview of SIV-Bench, showing its diverse videos spanning various social interactions and sample QAs for three task dimensions: Social Scene Understanding (SSU), Social State Reasoning (SSR), and Social Dynamics Prediction (SDP), along with their fine-grained sub-tasks.

085 features and grasp the underlying states. **3) Social**
 086 **Dynamics Prediction (SDP)** enables the model to
 087 reason about how interactions proceed over time
 088 or under alternative conditions, capabilities essen-
 089 tial for a flexible and human-like understanding of
 090 social scenarios (Ramnani and Miall, 2004; Byrne,
 091 2016). It involves factual prediction and counterfac-
 092 tual prediction. The former anticipates upcoming
 093 actions or changes in social states, while the latter
 094 examines how alterations in social scenes or states
 095 can affect interaction outcome.

096 *Second*, to systematically and comprehensively
 097 evaluate the three dimensions, we introduce **SIV-**
 098 **Bench (Social Interaction Video Benchmark)**
 099 grounded in social relation theory, which has
 100 been widely regarded as a foundational frame-
 101 work for understanding social behavior (Thibaut,
 102 2017; Burkitt, 1997; Hartup, 1989). Specifically,
 103 SIV-Bench is built on Fiske’s Relational Models
 104 Theory (Fiske, 1992), categorizing social interac-
 105 tions via four foundational models (*Communal*
 106 *Sharing*, *Authority Ranking*, *Equality Matching*,
 107 and *Market Pricing*), instantiated through 14 rela-
 108 tion types (e.g., parent-child, friends, colleagues).
 109 This relational context underpins all three dimen-
 110 sions. It conditions how cues are interpreted in
 111 SSU (e.g., a gaze between colleagues vs. lovers),
 112 modulates mental-state inference in SSR (e.g., crit-
 113 icism from a mentor vs. a stranger), and shapes
 114 SDP by constraining future behaviors and coun-

115 terfactuals through relational norms and history
 116 (e.g., siblings vs. business rivals). SIV-Bench
 117 comprises 2,792 video clips sourced from TikTok and
 118 YouTube, and a QA pipeline that combines adver-
 119 sarial filtering with human verification produces
 120 5,455 high-quality questions. To systematically as-
 121 sess the contribution of linguistic cues, SIV-Bench
 122 further provides audio tracks and three subtitle con-
 123 ditions: the original version (origin), a version with
 124 transcribed dialog added (+sub), and a version with
 125 all on-screen text removed (-sub).

126 In the experiments, we evaluate a broad set of
 127 models, including leading commercial MLLMs,
 128 strong open-source MLLMs, and video-specialized
 129 models. Overall, models are comparatively strong
 130 on SSU but weak on high-level reasoning, particu-
 131 larly within SSR, where **relation inference** exhibits
 132 systematic confusions. Our fine-grained analyses
 133 reveal recurring error sources, including selecting
 134 plausible but secondary relations, over-reliance on
 135 contextual cues, insufficient commonsense social
 136 reasoning, and missed perceptual details. SDP
 137 remains challenging as well, with counterfactual
 138 cases being relatively better handled by top models.
 139 We further study textual cues via controlled subtitle
 140 ablations and observe a task-dependent effect. Re-
 141 moving language typically has a limited impact on
 142 basic perception, but can degrade performance on
 143 inference-intensive SSR/SDP cases. Furthermore,
 144 on a hard diagnostic subset with brief explanations

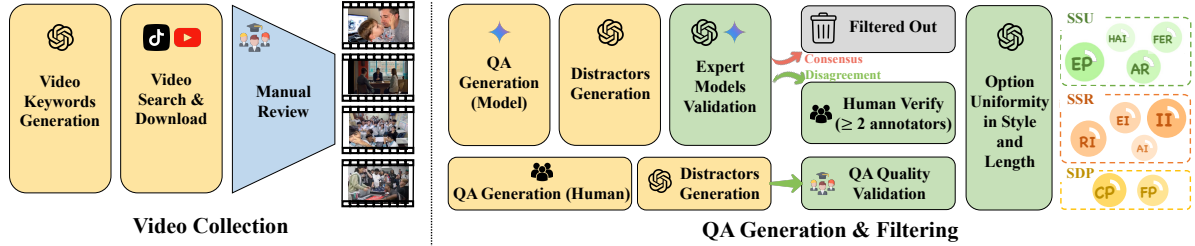


Figure 2: The SIV-Bench construction pipeline, detailing the data collection process (left), and the QA generation & filtering process (right) with human-LLM collaboration. In the diagram, ■ blocks indicate content (like keywords, video and QA) generation steps, ■ blocks represent validation or modification stages.

145 and a human baseline, we observe a large remain- 183
 146 ing gap to humans, underscoring the difficulty of 184
 147 robust, human-aligned social reasoning. 185

148 Our key contributions are as follows: **1)** We pro- 186
 149 pose a novel analytical framework that structurally 187
 150 decomposes the complex task of multimodal soci- 188
 151 al interaction understanding and reasoning into 189
 152 three core, interrelated dimensions, each further 190
 153 detailed into fine-grained sub-tasks. **2)** We intro- 191
 154 duce SIV-Bench, a new video benchmark specifi- 192
 155 cally curated for the analysis and comprehension 193
 156 of complex real-world social interactions. SIV- 194
 157 Bench comprises 2,792 real-world video clips rep- 195
 158 resenting 14 distinct social relationship types, and 196
 159 features 5,455 high-quality question-answer pairs 197
 160 generated through a human-LLM collaborative 198
 161 pipeline. **3)** Our comprehensive experiments on 199
 162 diverse MLLMs reveal the limitations in their cur- 200
 163 rent capacity for deep human social understanding. 201

164 2 SIV-Bench

165 This section presents the construction of SIV- 202
 166 Bench, including video collection (Section 2.1), 203
 167 QA composing (Section 2.2), and a comparison 204
 168 with existing benchmarks (Section 2.3). Figure 1 205
 169 offers some illustrative examples from SIV-Bench. 206
 170 Figure 2 shows the construction pipeline. 207

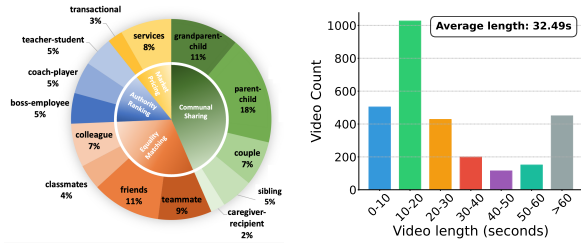
171 2.1 Video Collection

172 Firstly, we utilize GPT-4o-mini (Hurst et al., 2024) 212
 173 to generate comprehensive search keywords for 213
 174 each of the 14 relationship types (the word-clouds 214
 175 are shown in Figure 8), specifically including terms 215
 176 associated with varying degrees of intimacy, both 216
 177 positive (e.g., "love", "encourage") and negative 217
 178 (e.g., "conflict", "fight"). Leveraging these key- 218
 179 words, we conduct targeted searches and down- 219
 180 load initial video candidates from TikTok and 220
 181 YouTube platforms using Python libraries such as 221
 182 TikTokApi and yt-dlp, yielding approximately 222

5000 raw video clips. Each video is then manu- 183
 ally reviewed by the authors to ensure it contains 184
 clearly observable and meaningful social interac- 185
 tions. Videos are excluded if they do not depict 186
 clear social interaction (e.g., a vlogger speak- 187
 ing directly to the camera without interacting with 188
 others), if the interaction context does not fit within a 189
 set of well-defined interpersonal scenarios (e.g., an 190
 interview setting with scripted dialogue), or if the 191
 dominant interaction is difficult to identify due to 192
 the presence of multiple overlapping social dynam- 193
 ics (e.g., a large multi-generational family posing 194
 for a group photo). These criteria are designed to 195
 ensure that each included video primarily features 196
 one interpretable and coherent type of interpersonal 197
 interaction, allowing for more consistent analysis. 198

199 The SIV-Bench comprises **2,792** curated video 200
 clips, with statistics shown in Figure 3. The collec- 201
 tion showcases a rich diversity in social relation- 202
 ships (Figure 3a). Communal Sharing interactions 203
 are the most represented category, reflecting the 204
 naturalistic prevalence and psychological central- 205
 ity of these relationships in daily life (Simão and 206
 Seibt, 2014; Kameda et al., 2005). The other three 207
 relational models are also well represented, ensur- 208
 ing broad interpersonal coverage. In addition to 209
 relationship diversity, the benchmark covers het- 210
 erogeneous genres and filming/editing conventions, 211
 which helps test models under different visual nar- 212
 ratives and interaction realizations. In terms of 213
 duration (Figure 3b), clips average 32.49 seconds. 214
 While most clips are 10–20 seconds, the dataset 215
 spans a wide distribution, including many short 216
 clips (under 10 seconds) and a significant num- 217
 ber over 60 seconds. English predominates, but 218
 SIV-Bench also includes other languages, adding 219
 to multicultural diversity (detailed in Figure 9). De- 220
 tails on video diversity are in Appendix A.2.

221 To evaluate how different forms of textual in- 222
 formation affect MLLMs’ understanding of social



(a) Relation type distribution. (b) Video length distribution.

Figure 3: Video statistics for SIV-Bench.

interactions, we implement specific subtitle processing methods (Figure 4). Many original videos (‘Origin’) contain embedded on-screen text that often serves as scene descriptors or keywords (e.g., the ‘Buy and Sell’ text overlay). To focus on visual and auditory cues, we create a ‘-Subtitle’ version by removing such original textual overlays using video-subtitle-remover (YaoFANGUK, 2025). Conversely, to provide full access to spoken dialogue, we generate a ‘+Subtitle’ version. We employ Whisper-large-v3 (Radford et al., 2022) for audio transcription of the dialogue, and then use GPT to translate these transcriptions into English, ensuring consistent and high-quality subtitles (e.g., ‘He sold us his iPhone 12 with 256GB’).



Figure 4: Illustration of the three subtitle conditions.

2.2 QA Composing

Our Question-Answer (QA) composition pipeline is designed to move beyond simple recognition tasks and ensure a high density of challenging reasoning problems. The process consists of three main stages: scalable generation, adversarial filtering, and human-in-the-loop curation.

Scalable Generation. We begin by generating a large initial pool of diverse QA pairs for each video, leveraging the capabilities of Gemini-2.0-Flash with full video input. The prompts are carefully designed to elicit questions that span our SSU, SSR, and SDP evaluation dimensions. This stage typically yields an average of 10 QA pairs per

video. Subsequently, we employ GPT-4o-mini to generate four distractors for each QA pair, ensuring they are contextually relevant yet clearly distinguishable from the correct answer. Separating distractor generation significantly reduces ambiguity compared to simultaneous generation.

Adversarial Filtering. To ensure SIV-Bench serves as a rigorous diagnostic tool rather than a saturated test set, we implement a strict *Model-Consensus Filtering* mechanism. We utilize three commercial models (Gemini-2.0-Flash, Gemini-2.0-Pro, and GPT-4o-mini) to independently answer all candidate questions. We adopt an adversarial approach: questions correctly answered by all models (indicating trivial reasoning or visual obviousness) are discarded. We focus exclusively on the remaining "hard" samples where model consensus is not achieved, as these disagreements signal the need for deeper social reasoning.

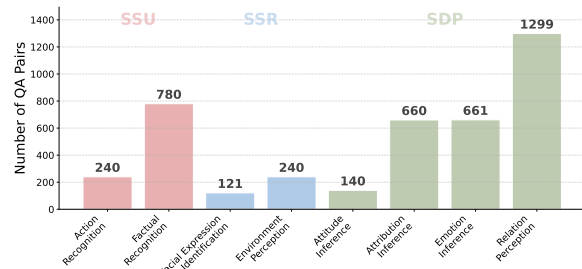


Figure 5: QA Distribution across 10 sub-tasks.

Human Verification and Authoring. The questions surviving the adversarial filter undergo a rigorous human-in-the-loop process to distinguish legitimate challenging samples from noise: (1) We recruit 20 human annotators to verify the non-consensus items. Each QA is reviewed by at least two annotators. Only questions where humans independently agree on the correct answer are retained. This yields 3,096 high-quality, model-challenging QAs. (2) To further expand the benchmark’s ceiling, annotators are instructed to generate novel, complex questions and their answers. These human-authored 2,359 QAs are designed to probe nuances often missed by model generators. Appendix B.2 details the annotation guidelines and interface. For all new human-generated QAs, distractors are created using GPT-4o-mini.

Finally, all curated QA pairs undergo an automated refinement stage to standardize linguistic style and option length, minimizing superficial cues that models might exploit. Full prompts are provided in Appendix B.1. Through this "Filter-

Table 1: Comparison of various benchmarks, including total number of items (**#Items**, representing the number of videos, dialogues, images, etc.), number of QA pairs (**#QAs**), annotation method (**Anno.**, M/A means manually/automatic manner), and tasks (*Task Types*: **SU** for Scene Understanding, **SR** for State Reasoning, **DP** for Dynamics Prediction). Note that *Task Types* here refer to general task categories, which include both social and physical scenarios. The table also shows whether each benchmark includes **Multi-Person Interaction**, covers **Various Relations**, is based on newly collected data (**Original Collection**), and provides subtitle/audio (**S.A.**).

Benchmark	#Items	#QAs	Anno.	Task Types			Multi-Person Interaction	Various Relations	Original Collection	S.A.
				SU	SR	DP				
<i>Social Relation Inference</i>										
DialogRE (Yu et al., 2020)	1,788	-	M	✓	✓	✗	✓	✓	✗	-
PIPA (Sun et al., 2017)	37,107	-	M	✗	✓	✗	✓	✓	✗	-
ViSR (Liu et al., 2019)	8,000	-	M	✗	✓	✗	✓	✓	✓	✗
<i>General Video Understanding and Reasoning</i>										
VideoMME (Fu et al., 2024)	900	2,700	M	✓	✓	✗	✓	✗	✓	✓
MLVU (Zhou et al., 2024)	1,730	3,102	M	✓	✗	✓	✓	✗	✗	✗
VideoVista (Li et al., 2024c)	894	24,906	A	✓	✗	✓	✗	✗	✗	✓
Social-IQ 2.0 (Wilf et al., 2023)	1,000	6,000	M	✓	✓	✗	✓	✓	✓	✗
Social Genome (Mathur et al., 2025)	272	1,486	M	✓	✓	✗	✓	✗	✗	✓
Perception Test (Patraucean et al., 2023)	11,620	38,000	M	✓	✗	✓	✗	✗	✓	✗
MVBench (Li et al., 2024b)	3,641	4,000	A	✓	✗	✓	✗	✗	✗	✗
Video-Bench (Ning et al., 2023)	5,917	17,036	A&M	✓	✓	✗	✓	✗	✗	✗
EgoSchema (Mangalam et al., 2023)	5,063	5,063	A&M	✓	✓	✗	✓	✗	✗	✗
SIV-Bench	2,792	5,455	A&M	✓	✓	✓	✓	✓	✓	✓

then-Verify" pipeline, SIV-Bench ultimately comprises **5,455** high-quality QA pairs. The distribution across core dimensions is illustrated in Figure 5, featuring a strong emphasis on SSR and SDP, enabling rich evaluation of higher-order social intelligence. Appendix B.3 provides detailed statistics confirming structural balance and diversity.

2.3 Comparison with Existing Benchmarks

Table 1 highlights the key differences between SIV-Bench and existing works. Prior datasets on social interaction, such as DialogRE (Yu et al., 2020) (text-only), PIPA (Sun et al., 2017) (image-only), and ViSR (Liu et al., 2019), are limited to relation recognition, single modalities, and lack task diversity and scalability due to manual annotation. Currently, most general video understanding benchmarks (like VideoVista (Li et al., 2024c) and MVBench (Li et al., 2024b)) lack a focus on social interaction; others (like MLVU (Zhou et al., 2024) and Video-Bench (Ning et al., 2023)), while touching upon it, do not fully cover social relations. Even Social-IQ 2.0 (Wilf et al., 2023), which concentrates on this area, has limitations in task diversity and dynamic reasoning. SIV-Bench is built on original data, combines manual and automatic annotations, and is one of the few benchmarks to provide subtitle and audio information, supporting richer multimodal social reasoning.

3 Experiments

3.1 Experimental Setup

We evaluate a diverse set of closed- and open-source MLLMs on SIV-Bench, including Gemini-2.0/2.5-Flash (Google, 2025a,b), Gemini-2.5-Pro (Doshi, 2025), GPT-4o (Hurst et al., 2024), o4-mini (OpenAI, 2025), Qwen2.5-VL-7B/72B-Instruct (Bai et al., 2025), mPLUG-Owl3 (Ye et al., 2024), InternVL3-8B/78B (Zhu et al., 2025a), LLaVA-OneVision (Li et al., 2024a), and LLaVA-Video (Zhang et al., 2024). Evaluations are conducted using VLMEvalKit (Duan et al., 2024), which provides a unified evaluation interface across MLLMs. All models generate responses under their default inference settings to ensure a fair comparison.

We use a standardized prompt (Figure 18) across all models, asking for both an option letter (e.g., 'A.') and the full answer text. To parse outputs robustly, we first extract a valid option letter; if missing, we match the raw output to all options using text similarity. To enable scalable and model-agnostic evaluation across heterogeneous MLLMs, we cast all tasks into a unified multiple-choice interface and report accuracy as the primary metric.

3.2 Overall Performance

The results shown in Table 2 highlights a clear performance stratification primarily driven by model scale. While the proprietary Gemini-2.5-Pro

Table 2: Evaluation results of MLLMs on SIV-Bench. Accuracy (%) on SSU, SSR, SDP, and Overall under different subtitle settings ('origin', '+sub', '-sub'). Statistical significance tests are reported in Appendix C.4.

Models	Params	Social Scene Understanding			Social State Reasoning			Social Dynamics Prediction			Overall		
		origin	+ sub	- sub	origin	+ sub	- sub	origin	+ sub	- sub	origin	+ sub	- sub
<i>Open-source MLLMs</i>													
mPLUG-Owl3	7B	46.11	45.94	46.34	39.78	39.50	38.13	44.30	46.08	44.35	42.06	42.42	41.15
LLaVA-OneVision	7B	39.20	39.41	40.08	41.95	43.65	38.66	43.79	44.51	39.42	41.97	43.04	39.42
LLaVA-Video	7B	50.22	50.61	50.56	39.33	38.19	36.14	41.60	42.14	39.94	41.09	42.00	38.66
Qwen2.5-VL-7B-Instruct	7B	51.22	50.88	50.22	40.24	38.94	37.66	42.69	43.82	42.24	44.02	44.21	41.65
InternVL3-8B	8B	56.83	56.13	56.50	40.35	40.90	37.92	44.53	45.52	44.40	45.82	46.05	44.56
Qwen2.5-VL-72B-Instruct	72B	75.73	76.24	73.54	<u>52.25</u>	<u>52.75</u>	<u>51.21</u>	59.02	58.40	57.78	<u>58.80</u>	<u>59.63</u>	<u>57.66</u>
InternVL3-78B	78B	71.46	73.66	71.76	51.65	52.39	50.14	55.77	56.28	54.25	55.46	56.32	54.50
<i>Closed-source MLLMs</i>													
o4-mini	-	78.83	79.04	78.13	50.47	51.30	48.99	56.89	56.00	55.26	55.68	55.89	54.54
GPT-4o	-	79.10	79.74	78.06	52.73	53.20	51.79	59.02	<u>60.59</u>	<u>58.60</u>	58.02	58.86	56.99
Gemini-2.0-Flash	-	78.46	78.16	78.34	51.89	52.43	49.78	57.59	58.63	55.70	56.40	57.23	54.64
Gemini-2.5-Flash	-	<u>81.70</u>	<u>82.14</u>	<u>79.71</u>	48.99	50.54	47.60	<u>59.47</u>	59.95	56.88	57.87	58.11	56.05
Gemini-2.5-Pro	-	85.07	85.41	84.94	54.30	54.85	52.32	60.45	61.54	58.83	61.65	62.40	60.22

achieves the state-of-the-art (62.40% w/ subtitles), large-scale open-source models demonstrate exceptional competitiveness; notably, Qwen2.5-VL-72B (59.63%) outperforms other leading proprietary systems like GPT-4o and Gemini-2.5-Flash (>58%). A significant capacity gap is evident, as larger models consistently dominate their smaller counterparts (e.g., Qwen-72B’s 59.63% vs. 7B’s 44.21%), suggesting that sufficient parameter count is a prerequisite for robust social reasoning. We verify that the observed performance gaps are statistically significant using paired bootstrap tests; full details are reported in Appendix C.4.

3.3 Decomposed Performance

Table 3: Relative performance change (%) against the baseline. The audio ablation (w/o Audio) is evaluated on Gemini-2.5-Flash, while subtitle variations (+Sub/-Sub) reflect the average impact across all models.

Condition	SSU	SSR	SDP	Overall
w/o Audio	-0.35	-2.44	-2.33	-1.40
+ Subtitle	+0.05	+0.28	+0.65	+0.15
- Subtitle	-0.97	-2.07	-1.68	-0.90

Subtitle and Audio Influence. To assess the impact of different modalities, we analyze the relative performance changes compared to the original video baseline, as shown in Table 3. The audio ablation is conducted on Gemini-2.5-Flash, while subtitle variations represent the average results across all evaluated models. The results demonstrate a clear hierarchy of modality importance. Removing audio leads to the most substantial performance drop (-1.40%), particularly in reasoning-intensive tasks (SSR: -2.44%), highlighting the necessity of

auditory cues for social interaction understanding and reasoning. For linguistic cues, the impact is notably asymmetric: while adding subtitles provides only marginal gains (+0.15%), their removal significantly hinders performance (-0.90%), especially in SSR and SDP. This suggests that while models may not always benefit from redundant text, the absence of critical textual context creates a significant bottleneck for complex social reasoning.

Task Influence. SSU focuses on recognizing visible elements and is generally the easiest dimension, where models achieve their highest scores. This is not always the case, since LLaVA-OneVision sometimes performs better on SSR or SDP. In SSU, larger models, including closed-source models and high-parameter open-source models, consistently outperform smaller ones, which likely reflects stronger perceptual ability from greater capacity and broader training (Alabdulmohsin et al., 2022). After SSU, SDP is usually the next most tractable, while SSR remains the most difficult because it requires inferring latent social states such as emotions and intentions. Figure 6 further breaks down performance across 10 fine-grained tasks in original videos. Stronger closed-source models such as the Gemini and GPT-4o form the outer ring and outperform smaller open-source models. In SSU, the performance gap is most pronounced in Action Recognition (AR) and Facial Expression Recognition (FER), suggesting advantages in capturing subtle visual cues. SSR is more challenging. Models perform moderately on Intent Inference (II) and Emotion Inference (EI), but they struggle most with Relation Inference (RI), which is often the lowest point on the radar. In

SDP, performance improves on Factual Prediction (FP) and Counterfactual Prediction (CP), which may reflect social commonsense acquired from language data. Most models perform better on CP than FP, possibly because hypothetical framing provides clearer reasoning cues. Representative failures are shown in Figure 19, 20, and 21.

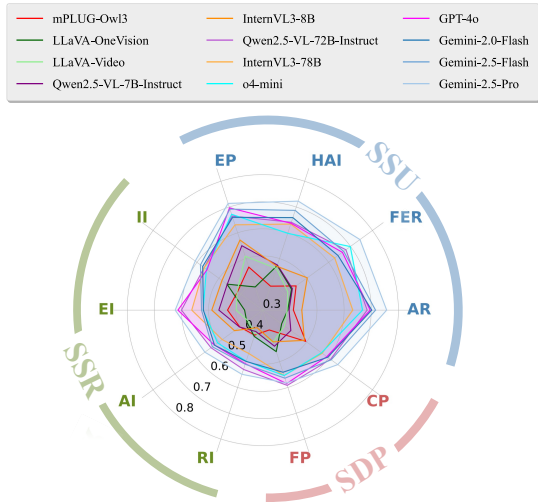


Figure 6: Radar chart of MLLM performance across the 10 fine-grained SIV-Bench sub-tasks.

3.4 In-Depth Analysis of Relation Inference

Since RI emerges as a key bottleneck and provides the relational grounding needed to interpret many higher-level social states, we take a closer look at RI and characterize its error structure quantitatively. Figure 7 aggregates RI predictions across the evaluated models into a confusion matrix, which reveals structured error clusters rather than random noise: a prominent pattern is Authority–Equality confusion, where hierarchical relations such as Boss–Employee and Coach–Player are frequently predicted as their egalitarian counterparts (Colleagues and Teammates). Errors also concentrate within relation families. The red dashed boundaries delineate the four foundational relational models, and off-diagonal clusters around these boundaries highlight systematic confusions between conceptually adjacent relation groups.

Qualitatively, we observe four common failure modes. (1) Models fail to differentiate primary vs. secondary relations in multi-relational scenarios, predicting a plausible but less salient relation. (2) Models are misled by scene- and human-induced cues, over-relying on stereotypical settings or surface language without validating the actual inter-

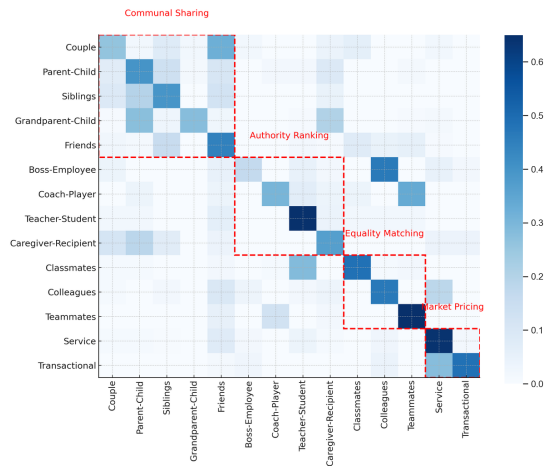


Figure 7: Aggregated Confusion Matrix for the Relation Inference (RI) task across all models. Red dashed lines delineate the four foundational relational models.

action (e.g., classroom → teacher–student). (3) Some errors reflect deficient commonsense social reasoning, where predictions contradict basic social conventions (e.g., two <friends> wearing different sports team uniforms but showing friendly behavior are misclassified as <teammates>, disregarding the fact that such uniforms typically signify competition). (4) Models ignore key perceptual details (e.g., age/role cues) that are decisive for disambiguating closely related relations. More failure cases are provided in Appendix C.2.

3.5 In-depth Analysis of Reasoning Process

To diagnose the cognitive boundaries of current MLLMs, we establish SIV-Bench-Hard, a subset of 200 QAs where models collectively demonstrate the highest error rates. This focused selection enables concentrated failure-mode analysis and facilitates a reliable human baseline, which is prohibitively expensive at full scale. We recruit three independent crowdworkers, distinct from the original annotators, to provide answer selections and free-text explanations. Table 4 highlights a substantial performance gap, with Gemini-3-Pro and GPT-5.1 achieving accuracies of 45.50% and 39.00%, compared to a human baseline of 74.40%.

LLM-judge results show a gap between reasoning presentation and social cognition. Models achieve strong structural scores (Logical Coherence, Relevance >4.0) but substantially lower Alignment and Depth (<3.5). We also find that reasoning similarity to human traces is positively associated with answer correctness (Figure 24, Table 9), proving that robust social interaction understanding

Table 4: Accuracy and LLM-judge scores (1-5) on reasoning quality (Rel=Relevance, Cohe=Logical Coherence, Depth=Depth of Analysis, Align=Alignment with Human, Conc=Conciseness, Ovrl=Overall).

Model	Acc%	Rel	Align	Cohe	Depth	Conc	Ovrl
Human	74.40	-	-	-	-	-	-
Gemini-3-Pro	45.50	4.66	3.30	4.67	3.49	4.87	4.10
GPT-5.1	39.00	4.58	3.29	4.65	3.26	4.88	4.00
Gemini-2.5-Pro	37.00	4.57	3.26	4.65	3.41	4.89	4.05
Gemini-2.5-Flash	32.32	4.48	3.17	4.55	3.22	4.87	3.95
GPT-4o-mini	29.00	4.45	3.20	4.56	3.12	4.91	3.90
Qwen2.5-VL-7B	24.50	4.00	2.89	4.21	3.05	4.45	3.63

requires human-like inference rather than purely fluent justifications. Interestingly, while Gemini-3-Pro maintains a lead in analysis depth over GPT-5.1 (+0.23), this does not improve alignment, indicating that an increased depth of reasoning does not inherently guarantee social accuracy.

A primary failure mode involves the inability to identify core social elements. In 40% to 54% of failures, models prioritize secondary visual cues over the "long-term separation" or "familial bonds" emphasized by humans. Furthermore, 45% of depth-related failures stem from a lack of social dynamic analysis, where models stop at direct causality and ignore power structures or cultural norms. We also observe a unique "excessive caution" bias in GPT-5.1, which frequently refuses to make valid socio-cultural inferences from appearance, resulting in the lowest alignment scores where humans readily use such cues for grounding. These insights suggest that MLLMs require multi-layered social grounding that explicitly partitions inference into observational, affective, and normative layers, enabling the models to bridge the gap between literal perception and human-aligned social intelligence.

4 Related Works

Video Benchmarks for MLLMs. Rapid progress in video MLLMs has motivated a growing suite of benchmarks spanning perception to high-level reasoning. Representative benchmarks cover multiple tasks (e.g., QA and summarization) over diverse video sources (Zhou et al., 2024; Fu et al., 2024; Li et al., 2024b; Fang et al., 2024; Ning et al., 2023). A major line of work emphasizes long-form understanding and temporal reasoning (Rawal et al., 2024; Mangalam et al., 2023; Wang et al., 2024b; Liu et al., 2024). Meanwhile, emerging benchmarks probe new deployment regimes, such as real-time/streaming video reasoning ((Xun et al., 2025)) and egocentric intent understanding with gaze cues

((Peng et al., 2025)). While these benchmarks provide broad coverage of video comprehension, they are not designed to systematically evaluate the fine-grained understanding and reasoning required for complex multi-person social interactions.

Evaluating Social Intelligence in AI Systems.

Beyond general video understanding, a parallel line of work evaluates social intelligence, moving from component-level signals (e.g., actions and affect) toward richer interpersonal reasoning and grounded social cognition. Recent benchmarks provide narrower but deeper probes, including Theory-of-Mind understanding from short films (Villa-Cueva et al., 2025), grounded social reasoning traces from interaction videos (Mathur et al., 2025), and non-verbal social reasoning without speech (Li et al., 2025). Complementary efforts assess emotional and social intelligence in video QA (Zhang et al., 2025) and social intelligence in authentic multi-turn human conversations (Huang et al., 2025), alongside interactive social environments (Zhou et al., 2023; Wang et al., 2024a). However, existing benchmarks largely evaluate social understanding in a task-specific or modality-specific manner. SIV-Bench complements them by providing a unified evaluation of social interaction understanding in real-world videos, spanning perception, social state reasoning, and dynamics prediction.

5 Conclusion

We introduce **SIV-Bench**, a video benchmark for evaluating MLLMs on real-world social interaction understanding and reasoning grounded in foundational social relationship models. SIV-Bench contains 2,792 videos and 5,455 curated QA pairs, organized into three dimensions (SSU, SSR, and SDP) as well as ten fine-grained sub-tasks for diagnosis. Extensive experiments show that current MLLMs perform relatively well on SSU but struggle with high-level reasoning, with Relation Inference exhibiting systematic confusions across relation types. Subtitle ablations suggest that textual signals can affect performance in a task-dependent manner, particularly for complex inference cases. On a challenging hard subset, we further observe a substantial gap to human performance, reinforcing the difficulty of human-aligned social reasoning. We hope SIV-Bench will serve as a unified testbed to track progress and drive more robust social intelligence in MLLMs.

562 Limitations

563 While SIV-Bench offers a structured diagnostic
564 benchmark for social interaction understanding, the
565 current size and coverage still leave room to expand
566 toward broader social settings, cultures, and inter-
567 action styles. Moreover, our primary evaluation
568 remains multiple-choice, which simplifies scoring
569 but does not fully capture open-ended grounding,
570 generation, or interactive behaviors. We include
571 SIV-Bench-Hard with brief explanations to par-
572 tially address this gap, yet it is relatively small
573 and still constrained in evaluation scope. Finally,
574 social interpretation can be subjective and cultur-
575 ally contingent, and some clips may admit mul-
576 tiple plausible readings, suggesting the need for
577 future work on more calibrated human protocols
578 and richer evaluation formats.

579 Ethical Considerations

580 SIV-Bench is built from publicly available videos
581 (e.g., TikTok and YouTube), and we carefully con-
582 sider licensing and privacy risks when creating and
583 releasing the benchmark. To reduce copyright and
584 platform-policy concerns, we will not redistribute
585 raw video files; instead, we will release annota-
586 tions along with video identifiers (e.g., URLs and
587 timestamps) and provide a script for users to re-
588 trieve content directly from the original platforms,
589 following common practice in prior video bench-
590 marks.

591 Because the videos may contain identifiable in-
592 dividuals, we applied manual filtering to exclude
593 content that appears overly private or sensitive, and
594 we encourage responsible downstream use. Our
595 human annotation procedures were designed to be
596 minimal-risk, with annotators compensated and
597 handled anonymously. Finally, while SIV-Bench
598 is intended to support transparent evaluation of
599 MLLMs’ social understanding, we acknowledge
600 the dual-use nature of social inference technologies
601 and recommend that users avoid deployment in
602 high-stakes settings and follow applicable privacy
603 and platform guidelines when using the benchmark.

604 References

605 Ibrahim M Alabdulmohsin, Behnam Neyshabur, and
606 Xiaohua Zhai. 2022. Revisiting neural scaling laws
607 in language and vision. *Advances in Neural Informa-*
608 *tion Processing Systems*, 35:22300–22312.

609 Ali Athar, Xueqing Deng, and Liang-Chieh Chen.

2025. Vicas: A dataset for combining holistic and
610 pixel-level video understanding using captions with
611 grounded segmentation. In *Proceedings of the Com-*
612 *puter Vision and Pattern Recognition Conference*,
613 pages 19023–19035. 614

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-
615 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie
616 Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl
617 technical report. *arXiv preprint arXiv:2502.13923*. 618

Joseph Berger, Bernard P Cohen, and Morris Zelditch Jr.
619 1972. Status characteristics and social interaction.
620 *American sociological review*, pages 241–255. 621

Ian Burkitt. 1997. Social relationships and emotions.
622 *Sociology*, 31(1):37–55. 623

Ruth MJ Byrne. 2016. Counterfactual thought. *Annual*
624 *review of psychology*, 67(1):135–157. 625

Xinlong Chen, Yuanxing Zhang, Chongling Rao,
626 Yushuo Guan, Jiaheng Liu, Fuzheng Zhang, Chengru
627 Song, Qiang Liu, Di Zhang, and Tieniu Tan. 2025.
628 Vidcapbench: A comprehensive benchmark of video
629 captioning for controllable text-to-video generation.
630 *arXiv preprint arXiv:2502.12782*. 631

Henghui Ding, Kaining Ying, Chang Liu, Shuting He,
632 Xudong Jiang, Yu-Gang Jiang, Philip HS Torr, and
633 Song Bai. 2025. Mosev2: A more challenging
634 dataset for video object segmentation in complex
635 scenes. *arXiv preprint arXiv:2508.05630*. 636

Tulsee Doshi. 2025. Build rich, interactive
637 web apps with an updated gemini 2.5 pro.
638 [https://blog.google/products/gemini/
639 gemini-2-5-pro-updates/](https://blog.google/products/gemini/gemini-2-5-pro-updates/). Accessed: 2025-05-
640 09. 641

Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu
642 Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang
643 Zang, Pan Zhang, Jiaqi Wang, and 1 others. 2024.
644 Vlmevalkit: An open-source toolkit for evaluating
645 large multi-modality models. In *Proceedings of the*
646 *32nd ACM International Conference on Multimedia*,
647 pages 11198–11201. 648

Xinyu Fang, Kangrui Mao, Haodong Duan, Xiangyu
649 Zhao, Yining Li, Dahua Lin, and Kai Chen. 2024.
650 Mmbench-video: A long-form multi-shot benchmark
651 for holistic video understanding. *Advances in Neural*
652 *Information Processing Systems*, 37:89098–89124. 653

Alan P Fiske. 1992. The four elementary forms of
654 sociality: framework for a unified theory of social
655 relations. *Psychological review*, 99(4):689. 656

Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li,
657 Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu
658 Zhou, Yunhang Shen, Mengdan Zhang, and 1 oth-
659 ers. 2024. Video-mme: The first-ever comprehensive
660 evaluation benchmark of multi-modal llms in video
661 analysis. *arXiv preprint arXiv:2405.21075*. 662

You will be given a video depicting human social relationships. Your task is to generate **question-answer (QA) pairs** to test the model's social reasoning ability. Each question must be clear and focused, addressing only one aspect at a time.

Question Requirements

- The questions should be challenging, requiring complex reasoning, an understanding of social norms, or real-world knowledge. Avoid surface-level observations.
- Generate **8-10 QA pairs**, distributed as follows:
 - Descriptive: Focus on different aspects of the scene, but you don't need to cover all of them. You can choose what you find most important:
 - Verbal characteristics (e.g., tone, pitch, speech style, etc.)
 - Non-verbal characteristics (e.g., facial expressions, body language, etc.)
 - Environmental features (e.g., location, setting, background details, etc.)
 - Human interaction elements (e.g., emotion, feelings, attitude, reactions, etc.)
- The question should be specific and unambiguous. Avoid questions like "describe the non-verbal characteristics of the video."
- Explanatory: Questions that explore the reasons behind observed behaviors.
- Predictive: A question that asks what is likely to happen next.
- Counterfactual: A question that explores hypothetical changes. For example, consider how the interaction might change if the individuals in the video had a stronger or weaker relationship, or if their relationship were of a different nature."

Answer Requirements

- Give a confident answer. Do not use words like "could" or "might." Provide only the essential information without additional explanations. Avoid making assumptions about relationships or identities in the QA. For example, DO NOT use descriptions like "the father and son in the video."

The output should be formatted in **JSON** as follows:

```
```json
{ "qa_pairs": [{ "category": "Descriptive", "question": "...", "answer": "...", }, { "category": "Explanatory", "question": "...", "answer": "...", },
{ "category": "Predictive", "question": "...", "answer": "...", }, { "category": "Counterfactual", "question": "...", "answer": "...", }] }
```
```

Figure 12: The prompt used to guide Gemini for the initial generation of question-answer pairs.

Given the following question and answer, generate four distractors that are reasonable and distinguishable from the correct answer. The sentence length, language style and grammar should be consistent with the answer.\n\n"
 f"Question: {qa_pair.get("question", "none")}\n"
 f"Answer: {qa_pair.get("answer", "none")}\n\n"
 "Provide the output in the following JSON format: \n"
 "{\n \"distractors\": [\n <distractor 1>\n ,\n <distractor 2>\n ,\n <distractor 3>\n ,\n <distractor 4>\n]\n}"

Figure 13: The prompt used to generate distractors for each QA pair.

Annotator Guidelines. To maintain consistency and high quality in human annotation, particularly for the creation of new challenging Question-Answers (QAs) for Subset 2, annotators are provided with a comprehensive set of guidelines, as illustrated in Figure 15. These instructions detail the process for responding to existing multiple-choice questions and, crucially, guide annotators in formulating one new challenging question per video. For this QA creation task, annotators are directed to ensure their questions primarily test one of our three core assessment dimensions: Social Scene Understanding (SSU), Social State Reasoning (SSR), or Social Dynamics Prediction (SDP). The guidelines, including example questions, also emphasize that new QAs must be demanding, clearly worded, and require a deep understanding of the video content rather than superficial observation.

Annotation Process. All human annotation

tasks for SIV-Bench, including the review of existing Question-Answers (QAs) and the generation of new QAs, are conducted using a custom-designed web interface. Figure 16 provides a representative example of this interface, which allows annotators to view the social interaction video, respond to provided multiple-choice questions, and submit their own newly authored questions and answers based on the video content and provided guidelines. This platform ensures a standardized environment for all human annotation contributions.

Quality Control and Inter-Annotator Agreement. To quantitatively address label quality, we compute the Inter-Annotator Agreement (IAA) using Fleiss' Kappa. For the initial set of approximately 27,000 QA pairs generated by our three LLM experts, the resulting Fleiss' Kappa is 0.52, indicating "moderate" agreement. For the questions that the LLMs disagreed, we measure the

1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027

For the understanding and reasoning of a social interaction video, it needs to be examined from the following three aspects of ability:

SSU (Social Scene Understanding): Action Recognition, Facial Expression Recognition, Human Attribute Identification, Environment Perception;

SSR (Social State Reasoning): Relation Inference, Emotion Inference, Intent Inference, Attitude Inference;

SDP (Social Dynamics Prediction): Factual Prediction, Counterfactual Prediction;

Now you need to analyze and categorize a multiple-choice question, output which level of test it belongs to and what aspects it specifically tests.

Here is the content:

Question: "{question}"

Options: "{options}"

Your output should be in the following format:

SSR: Relation Inference.

No other text or explanations are permitted.

Figure 14: Prompt used to instruct LLM for the final classification of Question-Answer pairs into one of the 10 fine-grained sub-tasks under SSU, SSR and SDP.

agreement between our human annotators, yielding a Fleiss' Kappa of 0.68. According to established standards (Landis and Koch, 1977), this value represents "substantial" agreement. Achieving this on the most difficult portion of our dataset confirms that humans can establish a reliable ground truth and justifies our methodology of retaining only those questions on which all assigned annotators unanimously agreed for the final benchmark.

B.3 QA statistics

To ensure the quality and balance of our SIV-Bench Question-Answer (QA) pairs, we conducted a statistical analysis of their structural properties, as summarized in Figure 17. This analysis examines several aspects: the average word count per multiple-choice option (Figure 17a) is relatively consistent across options A through E, minimizing length-based cues. The distribution of correct answers (Figure 17b) is fairly uniform across the five options, preventing positional bias. Furthermore, the overall question length (Figure 17c) peaks at around 11 words but shows a broad range, indicating variability in question complexity. Finally, an analysis of the first word in questions (Figure 17d) reveals a diverse set of interrogative types, led by 'what,' 'why,' and 'which,' reflecting a variety of reasoning challenges posed to the models.

B.4 Quality and Diversity Analysis

To further validate the quality of SIV-Bench, we conduct two analyses to ensure our benchmark encourages genuine comprehension over fitting to superficial patterns.

B.4.1 Analysis of Template Pattern Exploitation

SIV-Bench organizes tasks into SSU, SSR, and SDP to evaluate distinct dimensions of social cognition. To directly test whether models exploit surface-level cues, we analyze the properties of questions that a representative model (Gemini-2.0-Flash) answers correctly versus incorrectly. As shown in Table 5, while statistically significant differences in length and word count exist, the small effect sizes (Cohen's d) indicate no meaningful structural separation between the two groups. Furthermore, Table 6 shows a high cosine similarity between the embeddings of the correct and incorrect QA sets, along with similar internal distributions (intra-similarity and variance). This close alignment suggests that model performance is unlikely to rely on superficial statistical or semantic patterns.

B.4.2 Linguistic Diversity Analysis

The majority of questions in SIV-Bench are generated by large language models or written by human

Multimodal Social Interaction Understanding and Reasoning Test

Task Introduction

This task evaluates your ability to understand and reason about various aspects of social interactions depicted in videos, including the observable scene, the participants' potential mental states, and how the interaction unfolds or might change.

Instructions for Answering and Creating Questions

Each video is approximately **40 seconds** long on average. After watching each video, you will answer **1-3 multiple-choice questions**. Please select the answer you believe is most correct based on the video content.

Additionally, you will be asked to create **one challenging new question** about the video. Your question should encourage deep thinking about the social interaction. Please aim for your question to primarily test one of the following broad areas of understanding:

- 1. Social Scene Understanding (SSU):** Questions focusing on the observable aspects of the interaction.
(e.g., "What are people doing or saying explicitly? What are their visible expressions or key attributes? What does the environment tell us about the context?")
- 2. Social State Reasoning (SSR):** Questions requiring reasoning about the unobservable mental states of the individuals or the social relationship between them.
(e.g., "What might a person be feeling, intending, or believing? What is their attitude? What is the nature of their relationship?")
- 3. Social Dynamic Prediction (SDP):** Questions about how the interaction evolves, predicting future events, or considering hypothetical changes to the scenario.
(e.g., "What is likely to happen next as a result of this interaction? If a key element were different, how might the outcome change?")

Ensure your newly created question is challenging, clearly worded, and requires careful consideration of the video content, not just superficial observation.

Example Questions You Might Be Asked or Could Create:

- What clues in the setting or attire suggest this might be a formal event? (Tests SSU)
- Why do the two individuals exhibit contrasting emotional expressions after the announcement? (Tests SSR)
- If the person in blue had not intervened, what would likely have been the immediate consequence for the group? (Tests SDP)


Name:

Figure 15: Screenshot of the guidelines provided to human annotators, detailing the tasks of answering multiple-choice questions and creating new challenging questions about Social Scene Understanding (SSU), Social State Reasoning (SSR), and Social Dynamics Prediction (SDP).

Table 5: Length and word count comparison between correctly and incorrectly answered questions.

| Metric | Mean Length (chars) | Word Count |
|-----------|---------------------|---------------|
| Correct | 251.28 ± 93.88 | 42.13 ± 15.50 |
| Wrong | 228.96 ± 94.49 | 38.21 ± 15.78 |
| t-stat | 9.05 | 9.54 |
| p-value | 2.36e-19 | 2.63e-21 |
| Cohen's d | 0.2373 | 0.2516 |

1082 annotators, rather than using rigid templates. To
1083 objectively measure our benchmark's linguistic di-
1084 versity, we compared it against several prominent
1085 video QA benchmarks using two standard metrics:
1086 **Mean Semantic Distance** (average pairwise co-
1087 sine distance of sentence embeddings) and **Vec-
1088 tor Variance** (average variance across embedding
1089 dimensions). As shown in Table 7, SIV-Bench



Question 1: What is the tone of voice of the person speaking to the baby?

A Serious and commanding.

B Playful and teasing.

C Calm and soothing.

D Angry and frustrated.

E Indifferent and detached.

Please write the most challenging question you can think of for this video, along with its answer, separated by a line break.

Input the question and answer here...

Figure 16: Example of the web-based interface used by human annotators for watching videos, answering provided multiple-choice questions, and authoring new Question-Answer pairs for SIV-Bench.

Table 6: Semantic similarity analysis between correctly and incorrectly answered question sets.

| Metric | Value |
|---|----------|
| Cosine Similarity (TF-IDF) | 0.9765 |
| Cosine Similarity (SentenceTransformer) | 0.9759 |
| Correct Intra-similarity | 0.1689 |
| Wrong Intra-similarity | 0.1711 |
| Correct Embedding Variance | 0.002164 |
| Wrong Embedding Variance | 0.002157 |

exhibits high semantic diversity, ranking among
the top benchmarks. These results support that
SIV-Bench's questions are varied and not limited
to shallow templates, thereby promoting genuine
semantic understanding.

C Experimental Details

C.1 Settings Details

This section provides further details on our experi-
mental setup for evaluating MLLMs on SIV-Bench.
Figure 18 displays the standardized prompt tem-
plates employed for model evaluations, with dis-
tinct versions tailored to models based on their in-
put capabilities. For models that process sequences
of images, the "PROMPT for frames input" (Figure
18, Top) informs the MLLM that it will receive

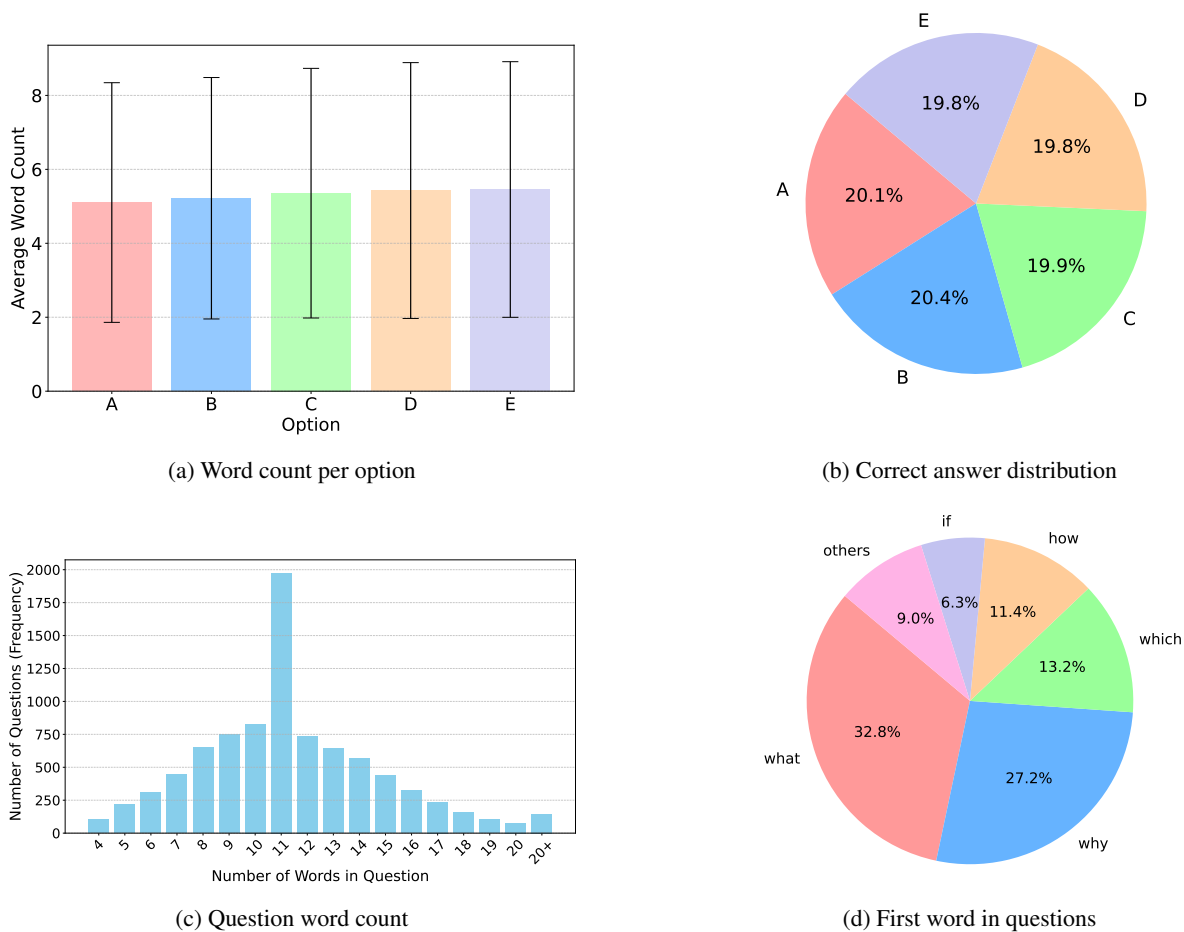


Figure 17: Statistical analysis of SIV-Bench question-answer (QA) pairs. (a) Average word count consistency across answer options, (b) distribution of correct answers among options, (c) question length distribution by word count, and (d) frequency of common first words in questions.

Table 7: Semantic diversity comparison across video QA benchmarks.

| Benchmark | Mean Dist. \uparrow | Variance \uparrow |
|------------------|-----------------------|---------------------|
| SIV-Bench | <u>0.8321</u> | <u>0.0022</u> |
| Video-Bench | 0.8811 | 0.0023 |
| Perception_Test | 0.7677 | 0.0020 |
| VideoVista | 0.7604 | 0.0020 |
| Social-IQ 2.0 | 0.7260 | 0.0019 |
| MVBench | 0.7124 | 0.0019 |
| EgoSchema | 0.5411 | 0.0014 |

a set of uniformly sampled frames from a video in chronological order. For models capable of direct video processing, the "PROMPT for videos input" (Figure 18, Bottom) is used. Both prompts clearly instruct the MLLM on its role, the task of answering multiple-choice questions based on the provided visual input, the expected JSON-like for-

mat for organizing answers (providing the exact text of the chosen option for each question), and a strict directive to avoid any extraneous text such as explanations or conversational remarks. This standardized, yet input-adaptive, prompting approach ensures consistency in task presentation across different model architectures.

For all evaluations, the specific inference parameters used for each model—such as temperature, top-p, or maximum new tokens—are adopted from their default configurations as provided within the VLMEvalKit (Duan et al., 2024) framework. This adherence to default settings aims to reflect the out-of-the-box capabilities of these models and ensure fair comparability. The experiments are conducted on two primary compute clusters. Cluster 1, utilized for evaluating the largest open-source models (Qwen2.5-VL-72B-Instruct and InternVL-78B), is equipped with an AMD EPYC 7642 48-Core Processor and 4x NVIDIA A100 GPUs. The total

runtime for the reported experiments on this cluster is approximately 3 days. Cluster 2, used for the remaining models, consists of an Intel(R) Xeon(R) Platinum 8369B CPU @ 2.90GHz and 8x NVIDIA RTX 3090 GPUs. The cumulative runtime for experiments on this cluster is approximately 2 days. It should be noted that the overall research project, including preliminary testing on earlier dataset versions and exploratory experiments not included in the final results, involved a greater amount of compute time than the specific durations reported for the final benchmark evaluations.

C.2 Failure Cases

This section presents several illustrative failure cases. We focus on examples from Gemini-2.0-Flash, a strong closed-source model, to highlight that even advanced MLLMs can fail on nuanced aspects of social perception, prediction, and reasoning. These examples are categorized by our primary assessment dimensions: SSU, SSR, and SDP, and are intended to offer concrete instances for future research and model development.

Figure 19 illustrates instances where Gemini-2.0-Flash fails on SSU tasks, which require accurate perception of explicit visual elements. **(a) In Action Recognition**, the model incorrectly identifies the man’s gesture as "He crosses his arms tightly" instead of the correct "He raises one eyebrow slightly", missing a subtle but distinct facial action. **(b) For Environment Perception**, when asked about the weather, the model failed to capture the details of the characters in the scene wearing thick scarves and down jackets to infer that the correct answer was "cold", but instead wrongly chose "wet". **(c) In Facial Expression Recognition**, the model describes the expression as "A mischievous smile" rather than the correct "A stoic glare", misinterpreting the nuanced facial expression display. **(d) For Human Attribute Identification**, concerning the child’s clothing, the model selects "A dress" instead of the correct "A set of pajamas", failing to correctly identify common apparel.

Figure 20 presents failure cases of Gemini-2.0-Flash on SSR tasks, which involve inferring unobservable mental states and relationships. **(e) In Intent Inference**, when a woman says "do you understand?" to a boy who bullies her son in an angry tone, it is to teach him a lesson and warn him not to bully her son again, not for "discourage any defiance", because in fact no child much younger than her can form defiance against her. **(f) For Emo-**

tion Inference, this employee is happy instead of scared after leaving because he successfully deceives the boss into giving him a vacation. **(g) In Attitude Inference**, the coworker is dissatisfied and disappointed with the cashier’s nervousness, panic and even physical reactions when seeing female customers. This could also be seen from his subsequent warning to the cashier not to do so anymore. **(h) For Relation Inference**, we present case studies on the failure patterns of the four common models listed in the main text in this task.

Figure 21 highlights errors made by Gemini-2.0-Flash in SDP tasks, which require predicting future events or reasoning about hypothetical scenarios. **(i) In Factual Prediction**, when asked if the person in the black shirt would be satisfied with the workers’ work, since the two of them have already reached an agreement with smiles at the end of the video, it could be inferred that the answer was "yes", but the model chooses another answer. **(j) For Counterfactual Prediction**, the video shows the dance interaction between a mother and her son. The question raised is what would happen if one of them had more dance experience. This can be inferred from the positive and relaxed interaction between the two in the video. The most likely answer is "The more experienced dancer leads and adapts movements". For example, a son leads his mother to learn dancing happily, rather than "The less experienced dancer hesitates and struggles to keep up". They have a good relationship, and the probability of negative performance like "hesitates" and "struggles" is lower.

C.3 Analysis of Chain-of-Thought Prompting

To investigate the impact of explicit reasoning on model performance, we conducted a preliminary experiment using a Chain-of-Thought (CoT) prompting strategy. We prepended the instruction "Let’s think step by step. First, output your reasoning process, and then output the final answer." to our standard evaluation prompt. The overall accuracy on the origin videos, with and without CoT, is presented in Table 8.

The results indicate that applying a generic CoT prompt did not yield significant performance improvements for most models. For several smaller open-source models (e.g., mPLUG-Owl3, LLaVA-Video), it resulted in a notable performance decrease. We observed that this is often because these models struggle to consistently adhere to the more complex two-stage output format (i.e., providing

PROMPT for frames input

You are an AI assistant responsible for answering questions about videos.

You will be provided with {} separate frames uniformly sampled from a video, \ the frames are provided in chronological order of the video.

Please analyze these images and provide the answers to the \ following multiple-choice questions about the video content.

If multiple questions are provided (with indices Q1, Q2, Q3, ...), \ you should organize your answers in the following json format:

- [Exact text of the chosen option for question 1],
- [Exact text of the chosen option for question 2],
- ...

Do NOT add any explanations, introductions, or concluding remarks.

PROMPT for videos input

You are evaluating a video based on the multiple-choice questions provided below.,

For each numbered question, select the best answer from the options listed.,

Your response MUST strictly follow this format:

- [Exact text of the chosen option for question 1],
- [Exact text of the chosen option for question 2],
- ... and so on for all questions.

Do NOT add any explanations, introductions, or concluding remarks.,

--- QUESTIONS ---

...

Figure 18: Standardized prompt templates used for evaluating MLLMs on SIV-Bench. Separate prompts are shown for models that process (Top) uniformly sampled frames and (Bottom) direct video input.

| Task | Video | Summary | QA ✓: Right answer ✗: MLLM's wrong answer |
|---------------------------------------|-------|--|---|
| (a)
Action Recognition | | This short video shows two female employees in an office environment filming themselves with their "big boss" standing behind them. The text overlay reads "tried this trend with our big boss." As a song plays, the two employees look at each other, seemingly performing a social media trend. They both, along with the boss, end up laughing, suggesting a lighthearted and fun interaction where the boss participated good-naturedly. | What gesture does the man make when the two people in front of him turn around?
A. He turns away and walks off.
B. He raises one eyebrow slightly.
C. He extends his arms with palms up.
D. He shakes his head slowly.
E. He crosses his arms tightly. |
| (b)
Environment Perception | | This video is a montage showing a student's final moments and memories from high school. It includes clips of eating with friends, school activities in the gym and classroom, commuting, and spending time outdoors, capturing everyday life before graduation. | What is the overall weather condition in the scene where the individual is walking outside in the evening?
A. Windy. B. Cold. C. Wet. D. Windy. E. Foggy. |
| (c)
Facial Expression Recognition | | This video is a meme. It shows a scene, likely from the movie "Identity Thief," where Melissa McCarthy's character abruptly punches Jason Bateman's character in the throat while in a hospital hallway. The text overlay, "What you really want to do to some coworkers that test your nerves on a daily basis," frames this act of violence as a humorous, exaggerated representation of the intense frustration one might feel towards annoying colleagues. | What is the facial expression of the person in the bottom center as they look at the camera?
A. An astonished gasp. B. A blank stare. C. A mischievous smile.
D. A pained frown. E. A stoic glare. |
| (d)
Human Attribute Identification | | In this funny video, a defiant young daughter sits on the kitchen counter and argues with her dad who is trying to get her to listen. She talks back assertively ("standing on business") about various things like not getting candy and her birthday. | What item of clothing does the child wear?
A. A dress. B. A winter coat. C. A set of pajamas. D. A baseball cap.
E. A pair of jeans. |

Figure 19: Examples of failure cases in Social Scene Understanding (SSU) tasks, including errors in Action Recognition, Environment Perception, Facial Expression Recognition, and Human Attribute Identification.

| Task | Video | Summary | QA \checkmark : Right answer \times : MLLM's wrong answer |
|------------------------|--|---|--|
| (e) Intent Inference |  | The video shows a blonde woman aggressively confronting a boy (presumably her son's bully) who now has a nosebleed. She threatens him, implying she caused the injury as retaliation for him breaking her son's glasses. When the bully's mother arrives, concerned, the blonde woman pretends the nosebleed was from an accidental fall, acting helpful while subtly warning the bully to be more careful. | Why does the woman ask the red-haired person "Do you understand?" in an angry tone?
A. To demonstrate their authority and discourage any defiance. B. To show concern for a potential misunderstanding and encourage dialogue. C. To emphasize the importance of their message and ensure compliance. D. To express frustration over a miscommunication and seek clarification. E. To provoke a reaction and assert dominance in the conversation. |
| (f) Emotion Inference |  | An employee asked for leave from a boss who was using a computer with headphones on the pretext of having an eye infection, and the boss readily agreed. But as soon as the employee stepped out of the door, he took off his eyes, gave a sly smile and quickly ran down the stairs. | What emotions does the person labeled 'employee' display after walking away?
A. Fear. B. Confusion. C. Happiness. D. Sadness. E. Anger. |
| (g) Attitude Inference |  | In a sports store, a young male employee becomes flustered and overgenerous when a woman trades in pink hockey gear. After she leaves, he regrets not offering more. An older coworker points out that the real issue isn't his awkwardness, but his visible erection. Embarrassed, the younger man calls himself a "horny monster." The older man tells him that sexual attraction is natural but urges him to learn control, instructing him to tuck it into his waistband and follow him for guidance. | What is the coworker's attitude to the cashier's behavior?
A. He is indifferent to it. B. He feels angry about it. C. He is disappointed. D. He is amused by the situation. E. He expresses his support. |
| (h) Relation Inference |  | In this funny video, a defiant young daughter sits on the kitchen counter and argues with her dad who is trying to get her to listen. She talks back assertively ("standing on business") about various things like not getting candy and her birthday. | Which kind of relationship plays a dominant role in this video?
Boss-Employee Colleagues
<i>Failure Type: Lack of Primary-Secondary Relation Differentiation.</i> |
| |  | In the video, there is a couple. The woman said to the camera, "if you want him, you have to get past me first," then kissed the boy. After that, she joked, "The best brother in the world." The boy was amused and said that this was not reasonable at all, and that we even came from different races. | Couple Siblings
<i>Failure Type: Scene and Human-Induced Misleading Cues.</i> |
| |  | A shop assistant in purple splashed a basin of water at the customers outside the store, then quickly hid in the store and closed the door. Angry customers kept knocking on the door to demand an explanation, but at this moment, the shop assistant began to take selfies, creating the illusion that customers were eager to open the store for epic deals, thereby indicating the popularity of the event. | Transactional Colleagues
<i>Failure Type: Deficiency in Commonsense Reasoning.</i> |
| |  | A couple in school uniforms were shooting a video of a gesture dance. They first made the four letters "L, O, V, E" with their hands, and then combined them to make a heart shape. | Couple Classmates
<i>Failure Type: Ignoring Key Perceptual Details.</i> |

Figure 20: Examples of failure cases in Social State Reasoning (SSR) tasks, highlighting difficulties in Intent Inference, Emotion Inference, Attitude Inference, and Relation Inference.



| Task | Video | Summary | QA \checkmark : Right answer \times : MLLM's wrong answer |
|-------------------------------|---|---|--|
| (i) Factual Prediction |  | The video is a comedic skit showing a contractor quoting \$5,000 for a job, only to be aggressively haggled down by a skeptical client who disputes the material costs and claims he could do the work himself. As the client challenges the price, referencing cheaper options like Temu, the contractor rapidly drops his quote until he accepts the client's firm counteroffer of \$3,000 with enthusiasm, highlighting the absurdity of extreme price negotiations. | Will the man in the black shirt be happy with the quality of the work?
A. It is unlikely he will care. D. No, he will be disappointed.
B. He might feel indifferent about it. E. Yes, they reached an agreement.
C. There is a chance he will be upset. |
| (j) Counterfactual Prediction |  | The person in black is showing off her proficient dance moves in front of the camera. Then, the man in white joins the dance from the right side of the frame. Their dance movements are very synchronized, and there are also interactions such as holding hands during the process. | How would the dance performance change if one person had more dance experience than the other?
A. Both dancers perform at the same skill level regardless of experience.
B. The more experienced dancer leads and adapts movements.
C. The less experienced dancer hesitates and struggles to keep up.
D. The music tempo slows down to accommodate both dancers.
E. The entire choreography is simplified for easier execution. |

Figure 21: Examples of failure cases in Social Dynamics Prediction (SDP) tasks, covering both Factual Prediction and Counterfactual Prediction.

reasoning before the final answer in the required format), leading to failures in our answer parsing logic.

The primary goal of SIV-Bench is to establish a fair, consistent, and reproducible evaluation of baseline model capabilities. Our standardized prompting strategy, which aligns with widely used toolkits like VLMEvalKit, ensures this fairness. While techniques like CoT are powerful for eliciting maximum performance from certain capable models (e.g., Gemini-2.5-Pro), introducing them as a default can create a confounding variable. Such a setup might shift the evaluation from testing inherent social reasoning to testing complex instruction-following abilities. Therefore, our main experiments use a direct-answering prompt to maintain a level playing field. Nonetheless, these findings suggest that developing more specialized reasoning methods tailored to social intelligence is a valuable direction for future work that builds upon this benchmark.

Table 8: Comparison of Overall Accuracy (%) on the origin videos with and without Chain-of-Thought (CoT) prompting.

| Model | Acc (Origin) | Acc (CoT) |
|-------------------------|--------------|-----------|
| mPLUG-Owl3 | 42.06 | 41.79 |
| LLaVA-OneVision | 41.97 | 41.22 |
| LLaVA-Video | 41.09 | 40.13 |
| Qwen2.5-VL-7B-Instruct | 44.02 | 44.76 |
| InternVL-8B | 45.82 | 44.99 |
| Qwen2.5-VL-72B-Instruct | 58.80 | 59.13 |
| InternVL-78B | 55.46 | 55.25 |
| o4-mini | 55.68 | 55.79 |
| GPT-4o | 58.02 | 57.91 |
| Gemini-2.0-Flash | 56.40 | 56.53 |
| Gemini-2.5-Flash | 57.87 | 57.71 |
| Gemini-2.5-Pro | 61.65 | 61.77 |

C.4 Statistical Significance Analysis

To validate the reliability of our comparative claims, we conducted McNemar’s tests on the full dataset ($N = 5,455$). This paired non-parametric test is appropriate for comparing the performance of two classifiers on the same dataset.

Model Ranking. We verified the leadership of our SOTA model. The performance difference between the top-performing **Gemini-2.5-Pro** (61.65%) and the second-best model **Qwen2.5-VL-72B** (58.80%) is highly statistically significant ($p < 0.001$), confirming the robustness of

the leaderboard rankings.

Task Difficulty. We confirmed that the performance stratifications across our three core dimensions are not due to chance. For Gemini-2.5-Pro, the performance gaps between **SSU** (85.07%) and **SDP** (60.45%), as well as between **SDP** and **SSR** (54.30%), are all highly statistically significant ($p < 0.001$).

Subtitle Influence. We performed significance testing on the subtitle conditions (Table 4). Our analysis reveals that the minor *Overall* improvement from adding subtitles (+sub) is not statistically significant ($p = 0.12$). However, the impact is task-dependent: the negative impact of removing text (-sub) is statistically significant for the **SSR** task ($p < 0.05$), and the benefit of added subtitles (+sub) is significant for the **SDP** task ($p < 0.05$).

D SIV-Bench-Hard Details

To rigorously evaluate the upper limits of current MLLMs, establish a robust human baseline, and probe the reasoning processes beyond simple accuracy, we curated and analyzed a challenging subset of our dataset, termed **SIV-Bench-Hard**. This section details the setup, human performance, and a multi-dimensional analysis of model reasoning quality on this subset.

D.1 Experimental Setup

Dataset Curation. We selected a subset of 200 questions from the human-generated portion of SIV-Bench. These questions were specifically chosen for their complexity and reliance on deep social understanding, filtering out items that could be solved via superficial visual cues.

Task Definition. Unlike the standard multiple-choice evaluation, this study required both human annotators and MLLMs to provide: (1) the selected answer option, and (2) a free-text *reasoning explanation* justifying their choice. The prompt is shown in Figure 22 and 23. This allows for a deeper examination of the cognitive process.

Participants. We recruited 3 independent human annotators to perform this task to establish a human baseline. We evaluated a suite of state-of-the-art MLLMs, including Gemini-3-Pro, GPT-5.1, Gemini-2.5-Pro, Gemini-2.5-Flash, Qwen2.5-VL-7B, and GPT-4o-mini.

Please watch the video and answer the following multiple-choice question.

Question: {question}

Options:
{options_text}

Please provide your answer strictly in the following format:

Answer: [Your chosen option letter, e.g., A/B/C/D/E]
Explanation: [Brief explanation in less than 50 words]

Important:

1. You must choose one option from the given choices
2. Keep your explanation concise and clear
3. Strictly follow the format above

Figure 22: Model prompt for SIV-Bench-Hard.

D.2 Quantitative Analysis of Reasoning

To quantify the divergence observed in the qualitative scores, we performed an embedding-based analysis. We encoded all reasoning texts (both human and model) using the paraphrase-multilingual-MiniLM-L12-v2 model.

Distinguishability. We trained a Random Forest classifier to distinguish between human and model explanations based on their embeddings. The classifier achieved an accuracy of **92.17%** (compared to a 50% random baseline). This high classification accuracy indicates that the latent semantic features of model reasoning are fundamentally distinct from those of humans.

Statistical Significance. We further validated this difference using Mann-Whitney U tests across the embedding dimensions. The tests confirmed that the distributions of human and model embeddings are significantly different ($p < 0.05$) on 8 out of 10 principal dimensions. Collectively, these results provide quantitative evidence that current MLLMs, despite their linguistic fluency, employ reasoning processes that are statistically distinguishable from human social cognition.

The SIV-Bench-Hard subset, along with the human reasoning annotations and analysis code, will be released to facilitate future research into bridging this gap.

You are an expert evaluator assessing the quality of AI model reasoning compared to human reasoning.

Task Context:
Question: {question}
Options: {options}

Human References (3 human reasonings):
Human 1: {human1_reason}
Human 2: {human2_reason}
Human 3: {human3_reason}

Model Response to Evaluate:
Model: {model_name}
Model's Answer: {model_answer}
Model's Reasoning: {model_reason}

Evaluation Instructions:
Please evaluate the model's reasoning quality across the following dimensions. Remember, we assume human-centered AI development, so human reasoning processes are the gold standard to approximate, even if a human's choice is incorrect.

Rate each dimension on a scale of 1-5:
- 1: Very Poor
- 2: Poor
- 3: Fair
- 4: Good
- 5: Excellent

Evaluation Dimensions:

1. **Relevance ()**: Does the reasoning address the key aspects of the question and video content?
2. **Alignment with Human Reasoning ()**: How similar is the model's reasoning process to the human references? Does it consider similar factors and perspectives?
3. **Logical Coherence ()**: Is the reasoning logically sound and internally consistent?
4. **Depth of Analysis ()**: Does the reasoning show deep understanding of social relations, emotions, and context?
5. **Conciseness ()**: Is the reasoning clear and concise without unnecessary verbosity?

Response Format (JSON only):

```
{
  "relevance": {
    "score": <1-5>,
    "justification": "<brief explanation>"
  },
  "alignment with human": {
    "score": <1-5>,
    "justification": "<brief explanation>"
  },
  "logical coherence": {
    "score": <1-5>,
    "justification": "<brief explanation>"
  },
  "depth of analysis": {
    "score": <1-5>,
    "justification": "<brief explanation>"
  },
  "conciseness": {
    "score": <1-5>,
    "justification": "<brief explanation>"
  },
  "overall_score": <average of 5 dimensions>,
  "overall_comment": "<2-3 sentences summarizing the model's reasoning quality>"
}
```

Respond ONLY with valid JSON. Do not include any text before or after the JSON.

Figure 23: LLM-judge prompt for scoring explanations in SIV-Bench-Hard.

D.3 Correlation Between Answer Correctness and Reasoning Fidelity

To rigorously verify that model performance is grounded in genuine social reasoning rather than superficial visual shortcuts (e.g., background cues or object co-occurrence), we analyze the relationship between the correctness of a model’s answer and the semantic quality of its reasoning trace. We hypothesize that if models were merely "guessing" via shortcuts, their generated explanations would lack alignment with human cognitive processes even when they fortuitously select the correct option.

We perform an embedding-based analysis on the *SIV-Bench-Hard* subset using the paraphrase-multilingual-MiniLM-L12-v2 model. For each of the six evaluated models, we calculate the cosine similarity between the model’s generated reasoning and the human expert ground truth. The results are stratified based on whether the model answers the multiple-choice question correctly or incorrectly.

Table 9: Quantitative comparison of reasoning similarity to human ground truth between correct and incorrect responses. The **Gap** (Δ) represents the increase in similarity when the model answers correctly. Significance is calculated using the Mann-Whitney U Test.

| Model | Similarity (Correct) | Similarity (Incorrect) | Gap (Δ) | Significance (p -value) |
|------------------|----------------------|------------------------|------------------|----------------------------|
| GPT-5.1 | 0.572 | 0.506 | +0.066 | < 0.001 (***) |
| Qwen2.5-VL-7B | 0.535 | 0.474 | +0.061 | < 0.01 (**) |
| GPT-4o-mini | 0.571 | 0.513 | +0.058 | < 0.01 (**) |
| Gemini-2.5-Pro | 0.541 | 0.498 | +0.043 | < 0.05 (*) |
| Gemini-2.5-Flash | 0.537 | 0.504 | +0.033 | < 0.05 (*) |
| Gemini-3-Pro | 0.527 | 0.494 | +0.033 | 0.067 (n.s.) |

As illustrated in Figure 24 and summarized in Table 9, we observe a statistically significant positive gap in similarity scores across 5 out of the 6 models. Notably, GPT-5.1 demonstrates the strongest effect, with a similarity gap of 0.066 ($p < 0.001$).

These findings provide empirical evidence that correctness in SIV-Bench is strongly correlated with human-like social reasoning. The significant degradation in reasoning alignment during failure cases suggests that models do not rely on "guessing" via superficial cues; rather, successful performance necessitates a cognitive process that mirrors human social understanding.

E Broader Impact

SIV-Bench is designed to foster positive advancements in artificial social intelligence, potentially

leading to more empathetic, context-aware, and collaborative AI systems for beneficial applications such as assistive technologies, improved human-AI teaming, and richer content understanding. However, enhancing AI’s grasp of social dynamics also presents risks. These capabilities could be misused for sophisticated manipulation, disinformation, or invasive surveillance, and unaddressed biases in data could be amplified, leading to inequitable outcomes. We offer SIV-Bench as a research tool to transparently assess MLLM capabilities and limitations in the social domain, thereby encouraging the community to proactively consider these ethical challenges and develop robust safeguards alongside continued innovation in social AI.

F Use of Ai Assistants

The LLM was utilized exclusively to aid and polish the writing, including for tasks such as improving grammar and clarity, refining sentence structure, and ensuring stylistic consistency

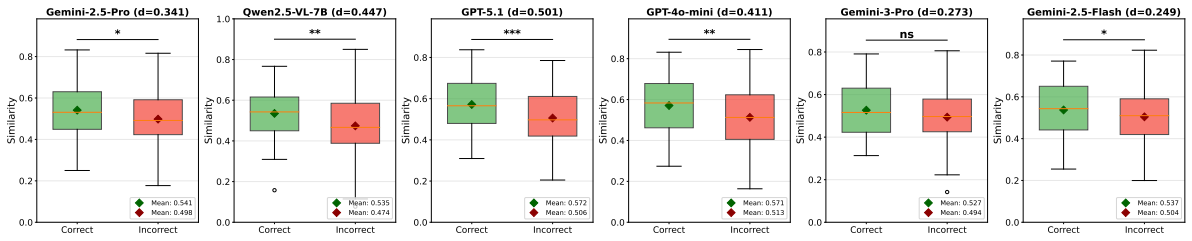


Figure 24: Box plots illustrating the distribution of semantic similarity scores between model reasoning and human ground truth, stratified by answer correctness. We observe a consistent trend where reasoning traces for correct answers exhibit significantly higher alignment with human social cognition compared to incorrect answers. Statistical significance is denoted by * ($p < 0.05$), ** ($p < 0.01$), and *** ($p < 0.001$).