

Do Vision Language Models infer human intention without visual perspective-taking? Towards a scalable "One-Image-Probe-All" dataset

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 At the core of understanding the knowledge grounding of Multimodal Large Language Models (MLLMs) are two key challenges: (1) ensuring fair comparability
 2 across concepts and (2) scaling multimodal datasets to reflect real-world complexity. This paper presents a solution through the **Omni-Perspective** benchmark,
 3 which scales the construction of a 5-level question-context-answers (QCAs) from **1 real-world image**. This benchmark pertains to 3 concepts along the Theory-of-
 4 Mind (ToM) ability hierarchy in humans and is further divided into 10 fine-grained subdifficulties. Through inference tasks, complexity, and ablation analysis, we
 5 evaluate over 2,200 consolidated QCAs on 61 MLLMs. Our findings reveal a key observation: MLLMs mostly follow the human ToM grounding pathway with
 6 exception of level-2 perspective taking. Furthermore, this dataset enables nuanced analysis of how such observations change across varying difficulty levels,
 7 modalities, distractor logic, and prompt types.
 8
 9
 10
 11
 12
 13

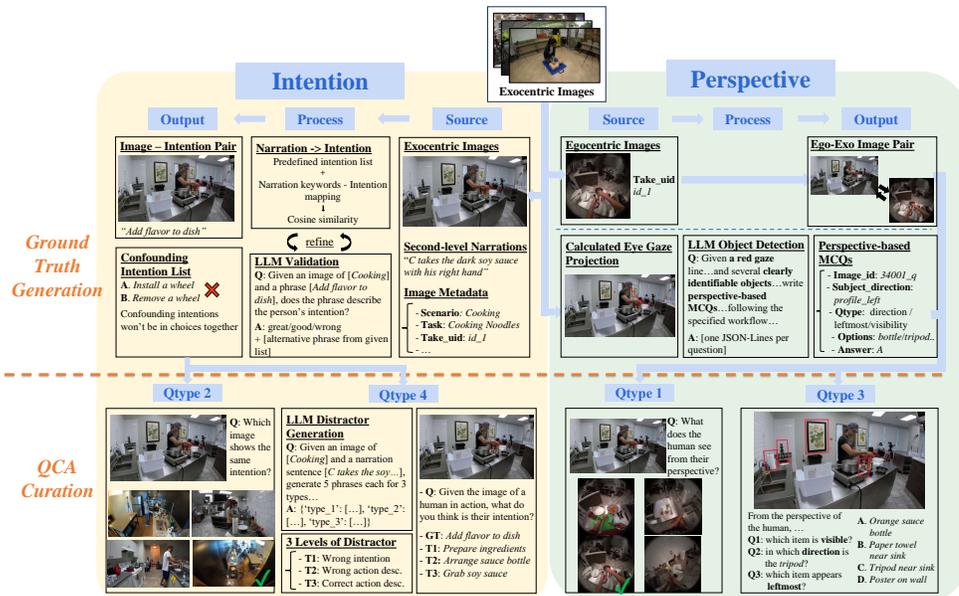


Figure 1: The scalable curation of Omni-perspective dataset

14 1 Introduction

15 Recent advances in Multimodal Large Language Models (MLLMs) have sparked growing interest
16 in evaluating their capacity for complex reasoning grounded in both visual and linguistic inputs.
17 However, rigorous assessment remains challenging due to the absence of scalable, cognitively
18 structured benchmarks that support controlled, hierarchical, and comparative probing across diverse
19 conceptual domains (Li et al., 2025). In this work, we address this gap by introducing a multi-
20 image, hierarchical, and concept-controlled Question-Context-Answer (QCA) generation framework,
21 designed to facilitate systematic evaluation of reasoning abilities across aligned tasks and cognitive
22 levels. This framework enables the use of reusable image-intention pairs, supports fine-grained control
23 over task difficulty, and allows for modular expansion to large-scale multimodal datasets—offering a
24 generalizable solution for cognitively diagnostic evaluation.

25 A key application of this framework is the assessment of visual perspective-taking (VPT) in relation
26 to Theory of Mind (ToM) capabilities (Premack and Woodruff, 1978; Barnes-Holmes et al., 2004;
27 Schaafsma et al., 2015). VPT involves understanding what others see (Level 1, or VPT-1) and how
28 they see it (Level 2, or VPT-2). Understood to be grounded in perspective-taking abilities, ToM
29 entails modeling others’ beliefs, goals, and intentions. These cognitive capacities develop in humans
30 along a staged trajectory (Barnes-Holmes et al., 2004; Barsalou, 2008; Schurz et al., 2021), offering a
31 natural scaffold for probing whether—and how—MLLMs internalize comparable representational
32 structures (Sucholutsky et al., 2023).

33 While several benchmarks have explored vision-language reasoning, many are limited in either
34 scope or ecological validity. For example, synthetic datasets such as CLEVR, CATER, and related
35 benchmarks have demonstrated the utility of 3D scene modeling and controlled object manipulation
36 for investigating compositional reasoning (Johnson et al., 2017; Girdhar and Ramanan, 2020).
37 However, these datasets operate in highly idealized environments, characterized by clean object
38 boundaries, minimal perceptual noise, and fully specified symbolic constraints. As a result, they tend
39 to overestimate generalization: models trained and evaluated in these “lab-grade” settings often fail to
40 transfer their reasoning capabilities to real-world scenes, where visual ambiguity, occlusion, temporal
41 dynamics, and social intent are critical (Mitchell and Krakauer, 2023).

42 Benchmarks such as ALPRO and VQA-X expand the modality coverage and include real images
43 or videos, but they often lack hierarchical cognitive task design or do not isolate the compositional
44 demands of ToM-related inference. Moreover, overreliance on language priors can inflate perfor-
45 mance in multimodal benchmarks even when visual inputs are ignored, undermining interpretability
46 (Dongxu Li, 2022; Park et al., 2018).

47 To address these limitations, we propose Omni-Perspective, a cognitively motivated benchmark
48 instantiated from our QCA generation framework. Built upon the rich, multimodal Ego-Exo4D
49 dataset, Omni-Perspective includes over 2,200 curated QCAs structured around a six-level hierarchy
50 that spans low-level spatial awareness to high-level belief reasoning. Each question is grounded
51 in a shared image-intention pair and linked to a cognitive hypothesis, enabling both depth and
52 comparability across reasoning types. Our scalable pipeline combines narration-intention mappings
53 with GPT-4o-assisted refinement, allowing for high-quality annotation at scale without extensive
54 manual labeling.

55 We evaluate 50+ MLLMs of varying modalities, sizes, and pretraining objectives, finding that while
56 many models perform well on spatial reasoning, they falter on belief-based or intention-predictive
57 tasks. This suggests a deviation from the developmental trajectory observed in human ToM, and
58 motivates architectural or training-level interventions to improve grounding and inference capabilities.

59 In summary, this work makes three key contributions:

- 60 1. **A multi-modal probing framework** for scalable, hierarchical, and controlled Question-
61 Context-Answer (QCA) generation, aligned with cognitive theory for systematic evaluation
62 of multimodal reasoning.
- 63 2. **A controlled and hierarchical benchmark**, *Omni-Perspective*, designed to probe Theory
64 of Mind (ToM) and visual perspective-taking abilities using real-world, multimodal visual
65 data from naturalistic scenarios.

66 3. **An empirical analysis** revealing consistent ToM-related failure modes in state-of-the-art
67 MLLMs, offering diagnostic insights and guiding principles for future model and training
68 improvements.

69 2 Related Works

70 2.1 MLLM related

71 2.1.1 Benchark

72 The field of Multi-modal Large Language Models (MLLMs) requires a comprehensive evaluation
73 of their remarkable capabilities to ensure that their development is progressing on a correct and
74 appropriate trajectory. Early benchmarks primarily focused on single tasks, such as VQA (Antol
75 et al., 2015), OK-VQA (Marino et al., 2019), MSCOCO (Lin et al., 2015), OCR (Liu et al., 2023),
76 and GQA (Hudson and Manning, 2019), but have become insufficient for thoroughly assessing the
77 broad multimodal perception and reasoning abilities of LMMs. In response, more holistic evaluations
78 have emerged, such as LAMM (Yin et al., 2024), MM-Vet (Yu et al., 2023), SEED-Bench (Li et al.,
79 2024), and MMBench (Liu et al., 2024c), which cover a wider range of capabilities.

80 2.1.2 Multi-modal Large Language Models

81 Recent advancements in multimodal learning have been largely driven by the unified modeling of
82 visual and textual data using transformers (Li et al., 2019; Xu et al., 2023; Tan and Bansal, 2019;
83 Alayrac et al., 2022; Radford et al., 2021). With the emergence of Large Language Models (LLMs),
84 state-of-the-art (SOTA) Multi-modal Large Language Models (MLLMs) (Liu et al., 2024a; Li et al.,
85 2023a) now integrate open-source LLMs (Touvron et al., 2023; Peng et al., 2023; Jiang et al., 2023),
86 aligning visual features with the embedding space of LLMs (Li et al., 2023b).

87 To enhance open-ended conversational abilities, LLaVA (Liu et al., 2024a) introduces a method to
88 distill the conversational capabilities of ChatGPT into MLLMs, resulting in a substantial performance
89 boost. This approach has since become a standard procedure in the field (Wang et al., 2023; Bai
90 et al., 2023; Gemini, 2023; Team, 2024; Sun et al., 2023; Li et al., 2022). As a result, MLLMs
91 have demonstrated competitive performance in complex tasks requiring high-level perception and
92 reasoning (Li et al., 2024; Liu et al., 2024a; Gemini, 2023; Fu et al., 2023; OpenAI, 2023), including
93 spatial reasoning (Chen et al., 2024; Cai et al., 2024), character recognition (Mori et al., 1999),
94 scene understanding (Cordts et al., 2016; Chen et al., 2017), action recognition (Jhuang et al., 2013;
95 Herath et al., 2017), and prediction (Lan et al., 2014; Kong and Fu, 2022), often reaching near-human
96 performance.

97 2.2 Visual perspective taking, Intentionality and Theory-of-Mind

98 The capacity to adopt another individual’s visual perspective is widely recognized as a foundational
99 component of social cognition and is considered a developmental precursor to theory of mind
100 (ToM)—the ability to attribute mental states such as beliefs, intentions, and knowledge to oneself and
101 others (Premack and Woodruff, 1978). While early research emphasized intention inference as central
102 to ToM, more recent accounts have identified visual perspective taking (VPT) as a perceptual substrate
103 supporting the emergence of mental state attribution. VPT is typically differentiated into two levels:
104 Level-1 perspective taking (VPT-1) involves representing *what* another agent can see (i.e., which
105 objects fall within their line of sight) whereas Level-2 perspective taking (VPT-2) entails representing
106 *how* those objects appear from another spatial viewpoint, including their orientation and relative
107 configuration (Kessler and Rutherford, 2010). Because VPT-2 requires mental transformations
108 of one’s egocentric reference frame—often instantiated through embodied simulation or motor
109 imagery—it has been proposed as a particularly robust route to social understanding, even though
110 such simulation is not strictly necessary for theory of mind reasoning in general (Hamilton et al.,
111 2009; Gallese and Goldman, 1998; Barlassina and Gordon, 2017).

112 Beyond these two levels, several developmental models posit a graded trajectory in which perceptual
113 perspective taking scaffolds increasingly abstract forms of social cognition. For example, Barnes-
114 Holmes and colleagues propose a sequence extending from recognition of differing viewpoints to
115 inferential use of perceptual access for epistemic judgments, prediction of actions based on true

116 beliefs, and ultimately the attribution of behavior based on false beliefs (Barnes-Holmes et al., 2004).
 117 Although terminological distinctions vary across frameworks, similar hierarchical structures were
 118 long proposed in traditional Piagetian theories of cognitive development (Piaget and Inhelder, 1969)
 119 and have since been elaborated in contemporary neurocognitive models that integrate perspective
 120 taking, empathy, and mental state attribution along continuous processing gradients (Schurz et al.,
 121 2021). Converging evidence from theoretical analyses suggests that tasks classified as measuring
 122 theory of mind in fact engage a distributed set of perceptual, inferential, and executive systems as
 123 opposed to being targeting a monolithic construct (Schaafsma et al., 2015; Quesque and Rossetti,
 124 2020; Barresi and Moore, 1996). These perspectives collectively support the view that higher-order
 125 social reasoning emerges through the gradual abstraction of perceptual and embodied capacities like
 126 visual perspective taking.

127 This developmental progression aligns with the theoretical framework of grounded cognition, which
 128 posits that high-level cognitive functions are constitutively supported by sensorimotor systems evolved
 129 for real-world interaction (Barsalou, 2008; Gallese, 2007). Accordingly, visual perspective taking
 130 offers a principled pathway through which embodied simulation mechanisms give rise to abstract
 131 representations of others’ mental states, supporting flexible and context-sensitive social inference in
 132 ecologically valid settings.

133 3 Omni-Perspective: A Scalable One-Image-For-All Benchmark From Visual 134 Perspective to Intentionality Understanding

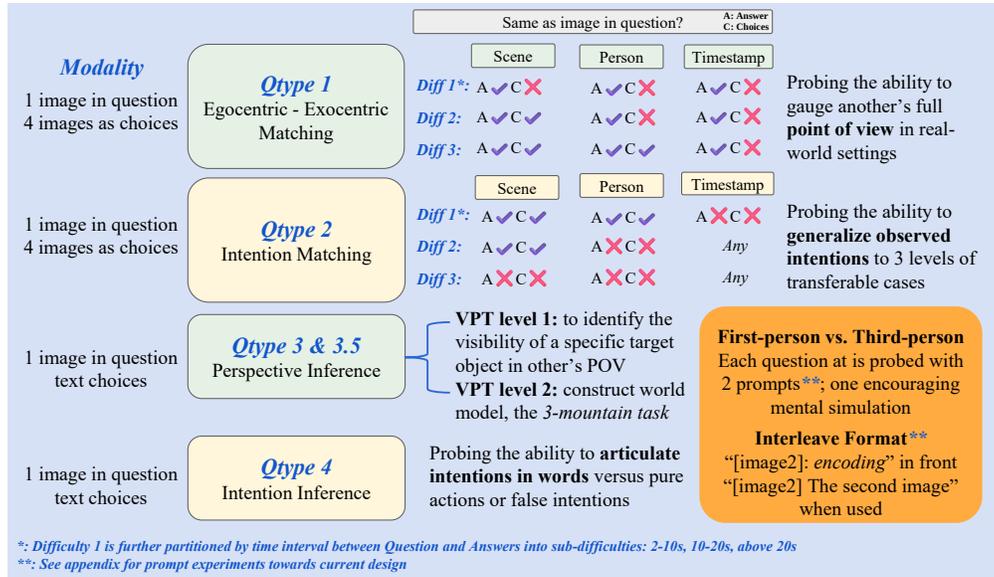


Figure 2: Overview of Omni-Perspective Bench

135 We define four distinct MCQ question types. Each is designed to target specific subskills aligned
 136 with the Theory-of-Mind hierarchy.

137 *Qtype 1 (Multi-image, Egocentric - Exocentric Matching)* - This question type presents the model
 138 with an exocentric image of a human in action and asks it to identify the corresponding view from
 139 four egocentric images. This task primarily probes Level-1 visual perspective-taking, requiring the
 140 model to reason about what the person sees based on spatial alignment and visual cues. Example
 141 prompt: “You are given an exocentric view of a person... Which of the following images best depicts
 142 what the person sees from their perspective?”

143 *Qtype 2 (Multi-image, Intention Similarity)* - In this task, the model is given an exocentric image
 144 of a person in action and asked to select the image depicting the most similar intention from four
 145 exocentric candidates. This question assesses the ability to generalize intention inference across
 146 individuals and scenes, contrasting with Qtype 4, which focuses on discriminating between actions

147 and intentions within a single context. Example prompt: “Given the image of a person performing an
148 action... Which of the following images shows someone with a similar intention?”

149 *Qtype 3 & 3.5 (Single-image, Spatial Perspective Inference)* – The model is shown an exocentric
150 image of a person and asked to determine the visibility or directional relation of an object from
151 that person’s perspective. All listed objects are visible in the scene, ensuring the task cannot be
152 solved through simple object detection or visual salience heuristics. This task is inspired by the
153 classic Piagetian "Three-Mountain Task" paradigm (Piaget and Inhelder, 1969), requiring the model
154 to construct a Level-2 perspective-taking world model—that is, to represent not only what another
155 agent sees, but how the scene is spatially organized from that agent’s viewpoint. The model must
156 perform an egocentric transformation of the scene, shifting reference frames to simulate another’s
157 first-person perspective. This demands an internal representation of spatial layout conditioned on
158 agent pose and orientation. Example prompt: “From the perspective of the woman in the black shirt
159 in the picture, which of the following items appears leftmost compared to the other choices?”

160 *Qtype 4 (Single-image, Intention Inference)* - This question presents a single exocentric image of a
161 person in action and asks the model to choose the most likely intention from four textual options. To
162 scale and control difficulty, distractor options are generated using a large language model (GPT-4o),
163 conditioned on the image and atomic action annotation (See Section A.3). This format targets
164 intention inference, requiring the model to go beyond object recognition. Example prompt: “You are
165 given an image of a human performing an action... What do you think is their intention?”

166 3.1 Dataset Overview

167 Ego-Exo4D Dataset

168 We base our evaluation framework on the Ego-Exo4D dataset (Grauman et al., 2024), a large-
169 scale, multimodal, multi-view video corpus featuring humans performing skilled activities such as
170 cooking, bike repair, and COVID-19 self-testing. Each recording session (take) includes synchronized
171 egocentric video from a head-mounted camera and up to four fixed exocentric views, capturing the
172 same activity from multiple viewpoints.

173 The dataset is structured hierarchically across scenarios (e.g., cooking), physical settings (e.g.,
174 kitchen), takes (video sessions), cameras (time synchronized viewpoints), and annotations. Annota-
175 tions include narration (atomic description of actions), procedural keysteps, and expert commentary,
176 making it particularly suited for our use case. Our dataset includes below retrieved distribution of
177 narrated images and goes beyond for prompt, ablation, and question evaluation analysis.

Task Type	Total Count
Cooking	70
Covid Test	101
Bike Repair	29
Total Tasks	$(70 + 101 + 29) \times (3 \times 2 + 2) = 200 \times 8 = 1600$

Table 1: Task counts by type

178 Generalization and Extensibility

179 Our benchmark pipeline is designed to generalize to any dataset offering (1) multi-view video and
180 (2) action-level annotation, e.g. the LEMMA dataset (Jia et al., 2020). This modularity enables the
181 broader application of our framework to evaluate ToM reasoning in multimodal LLMs across diverse
182 environments and tasks.

183 3.2 Benchmark Overview

184 Scalable Ground-Truth Image-Intention Pair

185 We construct a scalable set of image-intention pairs that serve as the foundation for all question
186 types in our benchmark. Four scenarios are selected based on the number of annotated takes and
187 coverage of non-repetitive actions. For each scenario, we define a set of high-level intentions and
188 identify representative image frames by applying a narration-keywords-to-intentions mapping. This

189 mapping is then empirically refined using GPT-4o, which evaluates each image-intention pair and
190 suggests corrections when misaligned. To minimize ambiguity, intentions that are visually similar
191 (e.g. *install a wheel* and *remove a wheel*) or sequentially entailed (e.g. *set up test* and *perform test*) —
192 referred to as confounding distractors — are excluded from co-occurrence within the same question.
193 This iterative process enables scalable generation of high-quality ground-truth image-intention pairs.
194 Refer to Section A.1 for more technical details.

195 **Comparability across Question Types**

196 *Reusing images across question types* - Each image-intention pair links to both egocentric and
197 exocentric views that are time-synchronized within the same take. This allows the same visual
198 context to be used for both perspective and intention questions, minimizing variability arising from
199 differences in scene content.

200 *Consistent question phrasing* - We standardize the linguistic structure of prompts across all question
201 types, avoiding shortcut through language cues. This reduces the risk of models exploiting superficial
202 lexical patterns and promotes a fairer assessment of reasoning capabilities.

203 *Uniform image abstraction level* - All images are sampled from real-world video footage with
204 similar resolution, camera specification, and background complexity. This avoids confounding effects
205 associated with abstraction level — such as those seen when mixing synthetic, staged, or cartoon
206 images with natural scenes — and ensures that all questions have perceptually comparable visual
207 input.

208 **First- and Third-Person Language Query**

209 Each question type is presented in both first-person and third-person point-of-view to distinguish
210 between two levels of perspective-taking. First-person prompts (e.g., “If you were the person in the
211 image, what is in your line of sight?”) encourage the model to take the subject’s role, reflecting a
212 mental simulation of world model and thus Theory-of-Mind reasoning (Barresi and Moore, 1996).
213 Third-person prompts (e.g., “Given the image with a person in action, what is their intention?”) treat
214 the model as an external observer, targeting Level-1 perspective-taking.

215 **Distractors with Multiple Difficulty Levels or Types**

216 Qtype 1 and 2 in our benchmark are presented at three levels of difficulty, defined by the design of
217 distractor choices. Difficulty increases as distractors become visually similar to the correct answer
218 (e.g. comparable objects or spatial arrangements), while easier distractors differ more clearly in
219 object type or environment setting. Qtype 4 does not use fixed difficulty levels but instead includes
220 three semantically distinct distractor types, ranging from low-level action descriptions to high-level
221 intentions. This controlled variation allows us to probe the robustness and granularity of model
222 reasoning under varying cognitive demand.

223 **4 Experiment**

224 **4.1 Setup**

225 **Inference:** With the curated QCA-prompt, we assessed an extensive collection of models spanning a
226 wide spectrum of architectures, parameter scales, and training methodologies. Our study encompassed
227 a total of 61 MLLMs. The selection included prominent proprietary models such as those from
228 the ChatGPT and Claude families, chosen for their established performance and widespread use.
229 The open-source cohort featured state-of-the-art models, including InternVL, the Qwen series, and
230 the recently released DeepSeek models, which have received increasing attention for their strong
231 performance in multimodal tasks. The open-source models under evaluation ranged in size from 1
232 billion to 110 billion parameters, enabling detailed performance analysis across scales. Proprietary
233 models were evaluated through API calls on standard personal computers. For open-source models,
234 we performed inference locally on a compute cluster equipped with 8×NVIDIA A100 80GB GPUs. In
235 practice, models under 13B parameters were typically executed on a single GPU, models between 13B
236 and 32B required two GPUs, those between 32B and 70B utilized four GPUs, and models exceeding
237 70B ran across all eight GPUs. We adhered closely to the official inference codebases provided by
238 model developers to ensure reproducibility and preserve model-specific inference optimizations. To
239 further ensure consistency and correctness in handling multimodal inputs, we developed a unified
240 evaluation toolkit capable of parsing and validating model responses across varying input formats.

241 **Evaluation:** To determine correctness, the model’s selected option is compared against the ground
 242 truth, with any instance labeled as FAIL in the matching process automatically marked incorrect.
 243 Specifically: 1) Template aatching is attempted first, using a set of pre-defined output formats to
 244 map the model’s response to one of the answer choices. 2) If template matching fails, the instance is
 245 passed to LLM matching, where a large language model—Llama-3.1-70B-Instruct(Grattafiori et al.,
 246 2024)—acts as a semantic judge to infer the intended answer choice.

247 To reduce the influence of answer-position bias, we adopt circular evaluation (Liu et al., 2024b).
 248 In this method, the multiple-choice options for each question are rotated across all possible posi-
 249 tions. The model must correctly answer all k permutations of a k-choice question to be considered
 250 accurate—ensuring that its success is not due to token position or randomness.

251 **4.2 Main Results**

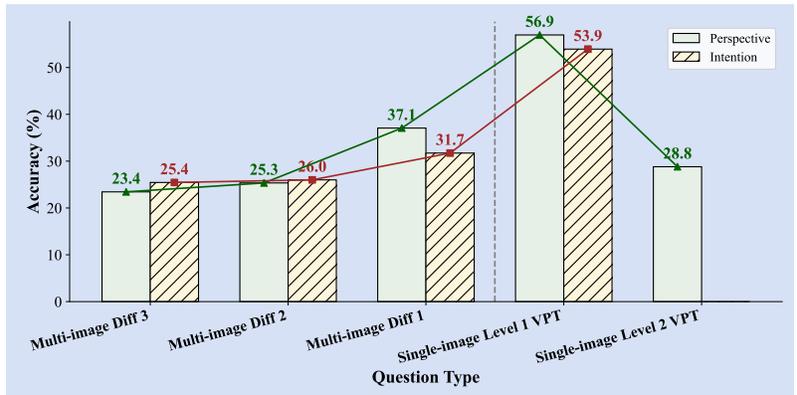


Figure 3: Comparative result between perspective taking and intention understanding across different difficulty levels and input types.

252 **Visual Perspective Grounding in Multi-Modal Large Language Models** We present comparative
 253 results (perspective vs. intention) across different difficulty levels (difficulty 1, 2 and 3) and input
 254 settings (single v.s. multi-image) in Figure 3. Several expected observations validate our benchmark
 255 design: 1. As difficulty increases from left to right (in the left section of the dashed line), both
 256 perspective and intention performance improve. 2. Performance on single-image tasks is consistently
 257 higher than on the three levels of multi-image tasks (to the right vs. left of the dashed line), largely
 258 due to the limited ability of MLLMs to process multi-image inputs.

259 Surprisingly, except for difficulty-3, where perspective is on par with intention, all other comparisons
 260 (difficulty-2, difficulty-1, and single-image) show better performance in perspective taking than
 261 in intention understanding. This contrasts with prior work Gao et al. (2025); Li et al. (2025). To
 262 further explore this distinction, we evaluate performance on level-2 perspective taking, specifically
 263 the three-mountain task (rightmost bar in Figure 3). In a fair comparison (both single-image), the
 264 three-mountain task performs lower than intention understanding, which aligns with previous findings
 265 Gao et al. (2025); Li et al. (2025). This suggests that the discrepancy between intention and level-2
 266 perspective taking is not due to a lack of visual perspective-taking ability, but rather factors such as
 267 limited spatial reasoning in the current MLLMs.

268 **Does prompting for Mental Simulation help?** Encouraging mental simulation (putting oneself in
 269 another’s shoes) is discussed to potentially benefit both visual perspective taking and intention under-
 270 standing ability, raising an intriguing question: Does explicitly prompting MLLMs to perform mental
 271 simulation improve performance on these tasks (Barlassina and Gordon, 2017)? A drill down into
 272 single image-prompt pairs (less confounded by distractor selection methods) shows that prompting
 273 MLLMs with first-person phrasing significantly improves performance on perspective-taking tasks (p
 274 $= 0.0321$) on spatial reasoning, while remaining inconclusive for intention understanding.

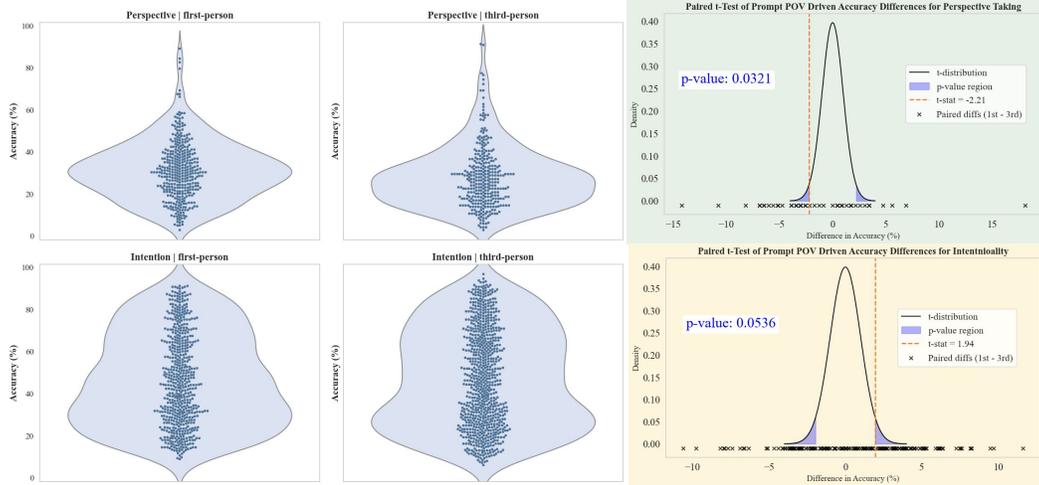


Figure 4: **Left:** Distribution of accuracy partitioned by probing concept and point-of-view of prompt; **Right:** Paired-T test results of single-image question for 2 types of prompts

275 **4.3 Distractor Ablation Tests**

276 For Qtype 4 - where distractors differ semantically (e.g. action descriptions versus high-level
 277 intentions) - we randomly select and mix choices from all three types for 200 questions. We then
 278 construct an additional ablation set of 95 randomly selected questions, each replicated into three
 279 versions containing distractors exclusively from one type. All other variables, including the image,
 280 prompt wording, and correct answer, remain constant for controlled comparison.

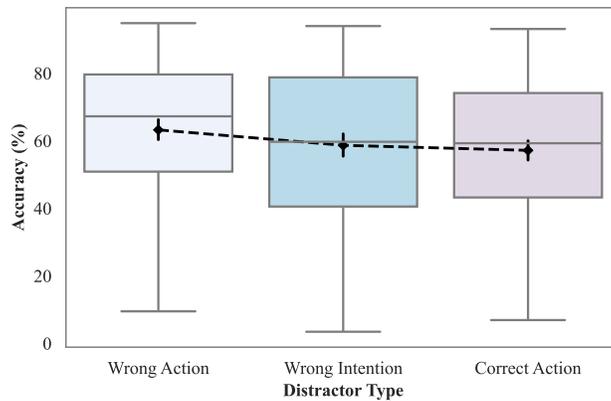


Figure 5: Accuracy by distractor type in Qtype 4 Ablation Test where the distractor type is controlled

281 Figure 5 reveals that average model accuracy varies across distractor types. Compared to the original
 282 Qtype 4 setup with an average accuracy of 53.9% (Figure 3), the ablation set yields consistently
 283 higher performance. This improvement likely stems from the reduced semantic variability, allowing
 284 models to exploit language-based shortcuts. Among the distractor types, wrong action results in the
 285 highest accuracy, which may be attributed to its double-layered deviation from the correct answer:
 286 it involves low-level action or object recognition rather than high-level intention inference, and the
 287 action described is itself incorrect, limiting the model’s ability to rely on object-centric heuristics.

288 **4.4 Benchmark Results**

289 **Stronger Models Exhibit Greater Differentiability on Easier Tasks** Accuracy varies widely
 290 across models at lower difficulty levels, with top-performing models such as *llava-video-72b-*

Model	Qtype 1 <i>Ego-Exo Match</i>			Qtype 2 <i>Intention Match</i>			Qtype 3 <i>Perspective Inference</i>	Qtype 4 <i>Intention Inference</i>
	Diff1	Diff2	Diff3	Diff1	Diff2	Diff3		
GPT-4o	97.24%	46.09%	28.09%	75.87%	36.28%	30.60%	31.37%	59.35%
deepseek-v1.5-small	40.57%	41.98%	41.36%	71.81%	73.47%	75.93%	57.08%	43.45%
Qwen2.5-VL-72B-Instruct	95.99%	45.05%	35.75%	79.26%	34.95%	32.41%	41.27%	61.38%
LLaVA-Video-72B-Qwen2_multi_frame	98.35%	42.69%	29.91%	68.62%	35.20%	37.96%	46.93%	59.23%
LLaVA-Video-7B-Qwen2_multi_frame	95.28%	38.44%	17.99%	67.55%	35.46%	41.67%	48.11%	51.73%
VILA1.5-40b	96.46%	32.78%	31.78%	56.91%	29.34%	23.15%	35.38%	75.68%
Mantis-8B-Idetics2	75.88%	39.25%	28.39%	66.57%	37.18%	32.33%	32.08%	59.85%
Llama-3-LongVILA-8B-256Frames	26.18%	29.72%	26.87%	59.04%	58.67%	58.33%	35.14%	73.88%
llava_next_interleave_7b	67.25%	26.55%	21.73%	49.71%	27.56%	26.72%	34.20%	64.38%
Llama-3-VILA1.5-8B	72.17%	28.30%	21.96%	40.43%	23.72%	23.15%	35.38%	60.93%
Ovis1.6-Gemma2-9B	69.50%	30.44%	25.88%	31.10%	25.64%	28.45%	44.34%	46.15%
Janus-Pro-1B	24.76%	26.18%	25.23%	43.09%	52.55%	56.48%	23.82%	32.50%
Vintern-3B-beta	44.88%	24.48%	25.88%	30.23%	25.51%	26.29%	35.38%	57.45%
InternVL2-4B	28.38%	24.09%	26.63%	37.79%	24.36%	23.71%	41.75%	51.00%

Table 2: Accuracy by model on each Qtype subtask. Best cells are bold and both best and second-best are shaded.

291 *qwen2_multi_frame* achieving near-perfect scores (98% on Qtype 1 Difficulty 1), while many others
292 remain below 30%. This variance diminishes as task difficulty increases: the standard deviation
293 in accuracy drops from nearly 20% at Difficulty 1 to under 5% at Difficulty 3. This pattern is
294 most evident among stronger, higher-capacity models, which show clear separation on simpler tasks
295 but converge to similarly low accuracy as complexity rises. Weaker models, by contrast, perform
296 consistently poorly across all levels with limited differentiation.

297 **Model Series Show Consistent Performance Trends** Certain model series consistently outperform
298 others. The *qwen2_5_vl_series* and *llava_video_multiframe_series* perform especially well at larger
299 scales, often scoring above 50% across tasks. Conversely, the *eagle_series_x4* and *x5* models
300 underperform broadly; even the 13B variant *eagle-x4-13b-plus* averages below 20%, suggesting
301 potential limitations in architecture, pretraining, or fine-tuning strategies.

302 **Scaling Model Size Yields Diminishing Returns Beyond a Point** Larger models generally out-
303 perform their smaller counterparts. For instance, in the *vila_series*, *vila1.5-40b* achieves a mean
304 accuracy of 48%, outperforming *vila1.5-13b* (39%) and *vila1.5-3b* (33%). However, some series
305 show marginal benefits from scaling: *llava-video-72b-qwen2_multi_frame* only slightly outperforms
306 its 7B counterpart (52% vs. 50%), and within *internvl2_series*, the jump from 2B to 40B offers
307 limited accuracy improvement. This suggests that beyond a certain threshold, increases in model size
308 alone may not yield proportionate gains.

309 5 Discussion

310 This study introduces the Omni-Perspective benchmark, a cognitively grounded and scalable frame-
311 work for probing MLLMs along the developmental hierarchy of ToM reasoning. We find that while
312 models perform reliably on Level-1 perspective-taking tasks, they consistently struggle with Level-2
313 visual perspective-taking and intention inference. This pattern generally aligns with developmental
314 theories suggesting that higher-order social reasoning builds upon more basic perceptual capaci-
315 ties, and is thus inherently more demanding. This suggests that MLLMs may be situated within a
316 human-like developmental trajectory for social cognition, albeit currently limited to lower levels
317 of the hierarchy. The observed performance gap reveals a key limitation in current MLLMs: their
318 limited capacity for mental simulation—a mechanism believed to support flexible, context-sensitive
319 social inference. Furthermore, our ablation studies show that model behavior is highly sensitive
320 to distractor configurations and prompt phrasing, indicating a reliance on superficial cues rather
321 than robust mental state representations. Taken together, the Omni-Perspective benchmark offers
322 a controlled and interpretable framework for evaluating social reasoning in MLLMs, while also
323 providing diagnostic insights into their architectural and training limitations.

324 In the meantime, we acknowledge that our benchmark relies on videos with sustained, non-transient
325 task focus as a proxy for intentionality, which may not generalize to brief or socially nuanced
326 intentions. It also assumes access to multiple viewpoints, limiting applicability to monocular settings.

327 **References**

- 328 Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican,
329 K., Reynolds, M., et al. (2022). Flamingo: a visual language model for few-shot learning. *Advances*
330 *in neural information processing systems*, 35:23716–23736.
- 331 Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C. L., and Parikh, D. (2015). Vqa:
332 Visual question answering. In *Proceedings of the IEEE international conference on computer*
333 *vision*, pages 2425–2433.
- 334 Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., et al. (2023).
335 Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- 336 Barlassina, L. and Gordon, R. M. (2017). Folk psychology as mental simulation. In Zalta, E. N.,
337 editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University,
338 summer 2017 edition.
- 339 Barnes-Holmes, Y., McHugh, L., and Barnes-Holmes, D. (2004). Perspective-taking and theory of
340 mind: A relational frame account. *The Behavior Analyst Today*, 5(1):15–25.
- 341 Barresi, J. and Moore, C. (1996). Intentional relations and social understanding. *Behavioral and*
342 *brain sciences*, 19(1):107–122.
- 343 Barsalou, L. W. (2008). Grounded cognition. *Annu. Rev. Psychol.*, 59(1):617–645.
- 344 Cai, W., Ponomarenko, Y., Yuan, J., Li, X., Yang, W., Dong, H., and Zhao, B. (2024). Spatialbot:
345 Precise spatial understanding with vision language models. *arXiv preprint arXiv:2406.13642*.
- 346 Chen, B., Xu, Z., Kirmani, S., Ichter, B., Sadigh, D., Guibas, L., and Xia, F. (2024). Spatialvlm:
347 Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the*
348 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465.
- 349 Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., and Yuille, A. L. (2017). Deeplab: Semantic
350 image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs.
351 *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848.
- 352 Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth,
353 S., and Schiele, B. (2016). The cityscapes dataset for semantic urban scene understanding. In
354 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223.
- 355 Dongxu Li, Junnan Li, H. L. J. C. N. S. C. H. (2022). Align and prompt: Video-and-language
356 pre-training with entity prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
357 *and Pattern Recognition*.
- 358 Fu, C., Chen, P., Shen, Y., Qin, Y., Zhang, M., Lin, X., Yang, J., Zheng, X., Li, K., Sun, X., Wu, Y.,
359 and Ji, R. (2023). Mme: A comprehensive evaluation benchmark for multimodal large language
360 models. *arXiv preprint arXiv: 2306.13394*.
- 361 Gallese, V. (2007). Before and below ‘theory of mind’: embodied simulation and the neural correlates
362 of social cognition. *Philosophical Transactions of the Royal Society B: Biological Sciences*,
363 362(1480):659–669.
- 364 Gallese, V. and Goldman, A. (1998). Mirror neurons and the simulation theory of mind-reading.
365 *Trends in cognitive sciences*, 2(12):493–501.
- 366 Gao, Q., Li, Y., Lyu, H., Sun, H., Luo, D., and Deng, H. (2025). Vision language models see what
367 you want but not what you see.
- 368 Gemini (2023). Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:*
369 *2312.11805*.
- 370 Girdhar, R. and Ramanan, D. (2020). CATER: A diagnostic dataset for Compositional Actions and
371 TEmporal Reasoning. In *ICLR*.

372 Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur,
 373 A., Schelten, A., Vaughan, A., Yang, A., Fan, A., Goyal, A., Hartshorn, A., Yang, A., Mitra,
 374 A., Sravankumar, A., Korenev, A., Hinsvark, A., Rao, A., Zhang, A., Rodriguez, A., Gregerson,
 375 A., Spataru, A., Roziere, B., Biron, B., Tang, B., Chern, B., Caucheteux, C., Nayak, C., Bi, C.,
 376 Marra, C., McConnell, C., Keller, C., Touret, C., Wu, C., Wong, C., Ferrer, C. C., Nikolaidis, C.,
 377 Allonsius, D., Song, D., Pintz, D., Livshits, D., Wyatt, D., Esiobu, D., Choudhary, D., Mahajan,
 378 D., Garcia-Olano, D., Perino, D., Hupkes, D., Lakomkin, E., AlBadawy, E., Lobanova, E., Dinan,
 379 E., Smith, E. M., Radenovic, F., Guzmán, F., Zhang, F., Synnaeve, G., Lee, G., Anderson, G. L.,
 380 Thattai, G., Nail, G., Mialon, G., Pang, G., Cucurell, G., Nguyen, H., Korevaar, H., Xu, H.,
 381 Touvron, H., Zarov, I., Ibarra, I. A., Kloumann, I., Misra, I., Evtimov, I., Zhang, J., Copet, J., Lee,
 382 J., Geffert, J., Vranes, J., Park, J., Mahadeokar, J., Shah, J., van der Linde, J., Billock, J., Hong, J.,
 383 Lee, J., Fu, J., Chi, J., Huang, J., Liu, J., Wang, J., Yu, J., Bitton, J., Spisak, J., Park, J., Rocca, J.,
 384 Johnstun, J., Saxe, J., Jia, J., Alwala, K. V., Prasad, K., Upasani, K., Plawiak, K., Li, K., Heafield,
 385 K., Stone, K., El-Arini, K., Iyer, K., Malik, K., Chiu, K., Bhalla, K., Lakhota, K., Rantala-Yearly,
 386 L., van der Maaten, L., Chen, L., Tan, L., Jenkins, L., Martin, L., Madaan, L., Malo, L., Blecher,
 387 L., Landzaat, L., de Oliveira, L., Muzzi, M., Pasupuleti, M., Singh, M., Paluri, M., Kardas, M.,
 388 Tsimpoukelli, M., Oldham, M., Rita, M., Pavlova, M., Kambadur, M., Lewis, M., Si, M., Singh,
 389 M. K., Hassan, M., Goyal, N., Torabi, N., Bashlykov, N., Bogoychev, N., Chatterji, N., Zhang, N.,
 390 Duchenne, O., Çelebi, O., Alrassy, P., Zhang, P., Li, P., Vasic, P., Weng, P., Bhargava, P., Dubal,
 391 P., Krishnan, P., Koura, P. S., Xu, P., He, Q., Dong, Q., Srinivasan, R., Ganapathy, R., Calderer,
 392 R., Cabral, R. S., Stojnic, R., Raileanu, R., Maheswari, R., Girdhar, R., Patel, R., Sauvestre, R.,
 393 Polidoro, R., Sumbaly, R., Taylor, R., Silva, R., Hou, R., Wang, R., Hosseini, S., Chennabasappa,
 394 S., Singh, S., Bell, S., Kim, S. S., Edunov, S., Nie, S., Narang, S., Raparthy, S., Shen, S., Wan, S.,
 395 Bhosale, S., Zhang, S., Vandenhende, S., Batra, S., Whitman, S., Sootla, S., Collot, S., Gururangan,
 396 S., Borodinsky, S., Herman, T., Fowler, T., Sheasha, T., Georgiou, T., Scialom, T., Speckbacher, T.,
 397 Mihaylov, T., Xiao, T., Karn, U., Goswami, V., Gupta, V., Ramanathan, V., Kerkez, V., Gonguet,
 398 V., Do, V., Vogeti, V., Albiero, V., Petrovic, V., Chu, W., Xiong, W., Fu, W., Meers, W., Martinet,
 399 X., Wang, X., Wang, X., Tan, X. E., Xia, X., Xie, X., Jia, X., Wang, X., Goldschlag, Y., Gaur,
 400 Y., Babaei, Y., Wen, Y., Song, Y., Zhang, Y., Li, Y., Mao, Y., Coudert, Z. D., Yan, Z., Chen, Z.,
 401 Papakipos, Z., Singh, A., Srivastava, A., Jain, A., Kelsey, A., Shajnfeld, A., Gangidi, A., Victoria,
 402 A., Goldstand, A., Menon, A., Sharma, A., Boesenberg, A., Baevski, A., Feinstein, A., Kallet, A.,
 403 Sangani, A., Teo, A., Yunus, A., Lupu, A., Alvarado, A., Caples, A., Gu, A., Ho, A., Poulton, A.,
 404 Ryan, A., Ramchandani, A., Dong, A., Franco, A., Goyal, A., Saraf, A., Chowdhury, A., Gabriel,
 405 A., Bharambe, A., Eisenman, A., Yazdan, A., James, B., Maurer, B., Leonhardi, B., Huang, B.,
 406 Loyd, B., Paola, B. D., Paranjape, B., Liu, B., Wu, B., Ni, B., Hancock, B., Wasti, B., Spence, B.,
 407 Stojkovic, B., Gamido, B., Montalvo, B., Parker, C., Burton, C., Mejia, C., Liu, C., Wang, C., Kim,
 408 C., Zhou, C., Hu, C., Chu, C.-H., Cai, C., Tindal, C., Feichtenhofer, C., Gao, C., Civin, D., Beaty,
 409 D., Kreymer, D., Li, D., Adkins, D., Xu, D., Testuggine, D., David, D., Parikh, D., Liskovich, D.,
 410 Foss, D., Wang, D., Le, D., Holland, D., Dowling, E., Jamil, E., Montgomery, E., Presani, E., Hahn,
 411 E., Wood, E., Le, E.-T., Brinkman, E., Arcaute, E., Dunbar, E., Smothers, E., Sun, F., Kreuk, F.,
 412 Tian, F., Kokkinos, F., Ozgenel, F., Caggioni, F., Kanayet, F., Seide, F., Florez, G. M., Schwarz, G.,
 413 Badeer, G., Swee, G., Halpern, G., Herman, G., Sizov, G., Guangyi, Zhang, Lakshminarayanan, G.,
 414 Inan, H., Shojanazeri, H., Zou, H., Wang, H., Zha, H., Habeeb, H., Rudolph, H., Suk, H., Aspegren,
 415 H., Goldman, H., Zhan, H., Damlaj, I., Molybog, I., Tufanov, I., Leontiadis, I., Veliche, I.-E., Gat,
 416 I., Weissman, J., Geboski, J., Kohli, J., Lam, J., Asher, J., Gaya, J.-B., Marcus, J., Tang, J., Chan, J.,
 417 Zhen, J., Reizenstein, J., Teboul, J., Zhong, J., Jin, J., Yang, J., Cummings, J., Carvill, J., Shepard,
 418 J., McPhie, J., Torres, J., Ginsburg, J., Wang, J., Wu, K., U, K. H., Saxena, K., Khandelwal, K.,
 419 Zand, K., Matosich, K., Veeraraghavan, K., Michelena, K., Li, K., Jagadeesh, K., Huang, K.,
 420 Chawla, K., Huang, K., Chen, L., Garg, L., A, L., Silva, L., Bell, L., Zhang, L., Guo, L., Yu, L.,
 421 Moshkovich, L., Wehrstedt, L., Khabsa, M., Avalani, M., Bhatt, M., Mankus, M., Hasson, M.,
 422 Lennie, M., Reso, M., Groshev, M., Naumov, M., Lathi, M., Keneally, M., Liu, M., Seltzer, M. L.,
 423 Valko, M., Restrepo, M., Patel, M., Vyatskov, M., Samvelyan, M., Clark, M., Macey, M., Wang,
 424 M., Hermoso, M. J., Metanat, M., Rastegari, M., Bansal, M., Santhanam, N., Parks, N., White,
 425 N., Bawa, N., Singhal, N., Egebo, N., Usunier, N., Mehta, N., Laptev, N. P., Dong, N., Cheng, N.,
 426 Chernoguz, O., Hart, O., Salpekar, O., Kalinli, O., Kent, P., Parekh, P., Saab, P., Balaji, P., Rittner,
 427 P., Bontrager, P., Roux, P., Dollar, P., Zvyagina, P., Ratanchandani, P., Yuvraj, P., Liang, Q., Alao,
 428 R., Rodriguez, R., Ayub, R., Murthy, R., Nayani, R., Mitra, R., Parthasarathy, R., Li, R., Hogan,
 429 R., Battey, R., Wang, R., Howes, R., Rinott, R., Mehta, S., Siby, S., Bondu, S. J., Datta, S., Chugh,
 430 S., Hunt, S., Dhillon, S., Sidorov, S., Pan, S., Mahajan, S., Verma, S., Yamamoto, S., Ramaswamy,

- 431 S., Lindsay, S., Lindsay, S., Feng, S., Lin, S., Zha, S. C., Patil, S., Shankar, S., Zhang, S., Zhang,
432 S., Wang, S., Agarwal, S., Sajuyigbe, S., Chintala, S., Max, S., Chen, S., Kehoe, S., Satterfield,
433 S., Govindaprasad, S., Gupta, S., Deng, S., Cho, S., Virk, S., Subramanian, S., Choudhury, S.,
434 Goldman, S., Remez, T., Glaser, T., Best, T., Koehler, T., Robinson, T., Li, T., Zhang, T., Matthews,
435 T., Chou, T., Shaked, T., Vontimitta, V., Ajayi, V., Montanez, V., Mohan, V., Kumar, V. S., Mangla,
436 V., Ionescu, V., Poenaru, V., Mihailescu, V. T., Ivanov, V., Li, W., Wang, W., Jiang, W., Bouaziz,
437 W., Constable, W., Tang, X., Wu, X., Wang, X., Wu, X., Gao, X., Kleinman, Y., Chen, Y., Hu, Y.,
438 Jia, Y., Qi, Y., Li, Y., Zhang, Y., Zhang, Y., Adi, Y., Nam, Y., Yu, Wang, Zhao, Y., Hao, Y., Qian, Y.,
439 Li, Y., He, Y., Rait, Z., DeVito, Z., Rosnbrick, Z., Wen, Z., Yang, Z., Zhao, Z., and Ma, Z. (2024).
440 The llama 3 herd of models.
- 441 Grauman, K., Westbury, A., Torresani, L., Kitani, K., Malik, J., Afouras, T., Ashutosh, K., Baiyya, V.,
442 Bansal, S., Boote, B., Byrne, E., Chavis, Z., Chen, J., Cheng, F., Chu, F.-J., Crane, S., Dasgupta, A.,
443 Dong, J., Escobar, M., Forigua, C., Gebreselasie, A., Hareh, S., Huang, J., Islam, M. M., Jain, S.,
444 Khirodkar, R., Kukreja, D., Liang, K. J., Liu, J.-W., Majumder, S., Mao, Y., Martin, M., Mavroudi,
445 E., Nagarajan, T., Ragusa, F., Ramakrishnan, S. K., Seminara, L., Somayazulu, A., Song, Y., Su,
446 S., Xue, Z., Zhang, E., Zhang, J., Castillo, A., Chen, C., Fu, X., Furuta, R., Gonzalez, C., Gupta,
447 P., Hu, J., Huang, Y., Huang, Y., Khoo, W., Kumar, A., Kuo, R., Lakhavani, S., Liu, M., Luo, M.,
448 Luo, Z., Meredith, B., Miller, A., Oguntola, O., Pan, X., Peng, P., Pramanick, S., Ramazanov, M.,
449 Ryan, F., Shan, W., Somasundaram, K., Song, C., Southerland, A., Tateno, M., Wang, H., Wang,
450 Y., Yagi, T., Yan, M., Yang, X., Yu, Z., Zha, S. C., Zhao, C., Zhao, Z., Zhu, Z., Zhuo, J., Arbelaez,
451 P., Bertasius, G., Crandall, D., Damen, D., Engel, J., Farinella, G. M., Furnari, A., Ghanem, B.,
452 Hoffman, J., Jawahar, C. V., Newcombe, R., Park, H. S., Rehg, J. M., Sato, Y., Savva, M., Shi, J.,
453 Shou, M. Z., and Wray, M. (2024). Ego-exo4d: Understanding skilled human activity from first-
454 and third-person perspectives.
- 455 Hamilton, A. F. d. C., Brindley, R., and Frith, U. (2009). Visual perspective taking impairment in
456 children with autistic spectrum disorder. *Cognition*, 113(1):37–44.
- 457 Herath, S., Harandi, M., and Porikli, F. (2017). Going deeper into action recognition: A survey.
458 *Image and vision computing*, 60:4–21.
- 459 Hudson, D. A. and Manning, C. D. (2019). Gqa: A new dataset for real-world visual reasoning and
460 compositional question answering.
- 461 Jhuang, H., Gall, J., Zuffi, S., Schmid, C., and Black, M. J. (2013). Towards understanding action
462 recognition. In *Proceedings of the IEEE international conference on computer vision*, pages
463 3192–3199.
- 464 Jia, B., Chen, Y., Huang, S., Zhu, Y., and chun Zhu, S. (2020). Lemma: A multi-view dataset for
465 learning multi-agent multi-task activities.
- 466 Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand,
467 F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L.,
468 Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. (2023). Mistral 7b. *arXiv preprint arXiv:
469 2310.06825*.
- 470 Johnson, J., Hariharan, B., van der Maaten, L., Fei-Fei, L., Zitnick, C. L., and Girshick, R. (2017).
471 Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*.
- 472 Kessler, K. and Rutherford, H. (2010). The two forms of visuo-spatial perspective taking are
473 differently embodied and subserve different spatial prepositions. *Frontiers in Psychology*, Volume
474 1 - 2010.
- 475 Kong, Y. and Fu, Y. (2022). Human action recognition and prediction: A survey. *International
476 Journal of Computer Vision*, 130(5):1366–1401.
- 477 Lan, T., Chen, T.-C., and Savarese, S. (2014). A hierarchical representation for future action prediction.
478 In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September
479 6-12, 2014, Proceedings, Part III 13*, pages 689–704. Springer.
- 480 Li, B., Ge, Y., Ge, Y., Wang, G., Wang, R., Zhang, R., and Shan, Y. (2024). Seed-bench: Benchmark-
481 ing multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer
482 Vision and Pattern Recognition*, pages 13299–13308.

- 483 Li, C., Xu, H., Tian, J., Wang, W., Yan, M., Bi, B., Ye, J., Chen, H., Xu, G., Cao, Z., Zhang, J., Huang,
484 S., Huang, F., Zhou, J., and Si, L. (2022). mplug: Effective and efficient vision-language learning
485 by cross-modal skip-connections. *arXiv preprint arXiv: 2205.12005*.
- 486 Li, J., Li, D., Savarese, S., and Hoi, S. (2023a). Blip-2: Bootstrapping language-image pre-training
487 with frozen image encoders and large language models. *CONFERENCE*.
- 488 Li, J., Li, D., Savarese, S., and Hoi, S. (2023b). Blip-2: Bootstrapping language-image pre-training
489 with frozen image encoders and large language models. In *International conference on machine*
490 *learning*, pages 19730–19742. PMLR.
- 491 Li, L. H., Yatskar, M., Yin, D., Hsieh, C.-J., and Chang, K.-W. (2019). Visualbert: A simple and
492 performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*.
- 493 Li, Y., Gao, Q., Zhao, T., Wang, B., Sun, H., Lyu, H., Luo, D., and Deng, H. (2025). Core knowledge
494 deficits in multi-modal language models.
- 495 Lin, T.-Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., Perona, P., Ramanan, D.,
496 Zitnick, C. L., and Dollár, P. (2015). Microsoft coco: Common objects in context.
- 497 Liu, H., Li, C., Wu, Q., and Lee, Y. J. (2024a). Visual instruction tuning. *Advances in neural*
498 *information processing systems*, 36.
- 499 Liu, Y., Duan, H., Zhang, Y., Li, B., Zhang, S., Zhao, W., Yuan, Y., Wang, J., He, C., Liu, Z., Chen,
500 K., and Lin, D. (2024b). Mmbench: Is your multi-modal model an all-around player?
- 501 Liu, Y., Duan, H., Zhang, Y., Li, B., Zhang, S., Zhao, W., Yuan, Y., Wang, J., He, C., Liu, Z., et al.
502 (2024c). Mmbench: Is your multi-modal model an all-around player? In *European conference on*
503 *computer vision*, pages 216–233. Springer.
- 504 Liu, Y., Li, Z., Yang, B., Li, C., Yin, X., Liu, C.-l., Jin, L., and Bai, X. (2023). On the hidden mystery
505 of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*.
- 506 Marino, K., Rastegari, M., Farhadi, A., and Mottaghi, R. (2019). Ok-vqa: A visual question answering
507 benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer*
508 *vision and pattern recognition*, pages 3195–3204.
- 509 Mitchell, M. and Krakauer, D. C. (2023). The debate over understanding in ai’s large language
510 models. *Proceedings of the National Academy of Sciences*, 120(13):e2215907120.
- 511 Mori, S., Nishida, H., and Yamada, H. (1999). *Optical character recognition*. John Wiley & Sons,
512 Inc.
- 513 OpenAI (2023). Gpt-4 technical report. *arXiv preprint arXiv: 2303.08774*.
- 514 Park, D. H., Hendricks, L. A., Akata, Z., Rohrbach, A., Schiele, B., Darrell, T., and Rohrbach, M.
515 (2018). Multimodal explanations: Justifying decisions and pointing to the evidence. In *2018*
516 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8779–8788.
- 517 Peng, B., Li, C., He, P., Galley, M., and Gao, J. (2023). Instruction tuning with gpt-4. *arXiv preprint*
518 *arXiv:2304.03277*.
- 519 Piaget, J. and Inhelder, B. (1969). *The Psychology of the Child*. Basic Books, New York.
- 520 Premack, D. and Woodruff, G. (1978). Does the chimpanzee have a theory of mind? *Behavioral and*
521 *Brain Sciences*, 1(4):515–526.
- 522 Quesque, F. and Rossetti, Y. (2020). What do theory-of-mind tasks actually measure? theory and
523 practice. *Perspectives on Psychological Science*, 15(2):384–396.
- 524 Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A.,
525 Mishkin, P., Clark, J., Krueger, G., and Sutskever, I. (2021). Learning transferable visual models
526 from natural language supervision. *arXiv preprint arXiv: 2103.00020*.

- 527 Schaafsma, S. M., Pfaff, D. W., Spunt, R. P., and Adolphs, R. (2015). Deconstructing and reconstruct-
528 ing theory of mind. *Trends in cognitive sciences*, 19(2):65–72.
- 529 Schurz, M., Radua, J., Tholen, M. G., Maliske, L., Margulies, D. S., Mars, R. B., Sallet, J., and
530 Kanske, P. (2021). Toward a hierarchical model of social cognition: A neuroimaging meta-analysis
531 and integrative review of empathy and theory of mind. *Psychological bulletin*, 147(3):293.
- 532 Sucholutsky, I., Muttenthaler, L., Weller, A., Peng, A., Bobu, A., Kim, B., Love, B. C., Grant, E.,
533 Groen, I., Achterberg, J., et al. (2023). Getting aligned on representational alignment. *arXiv*
534 *preprint arXiv:2310.13018*.
- 535 Sun, Q., Cui, Y., Zhang, X., Zhang, F., Yu, Q., Luo, Z., Wang, Y., Rao, Y., Liu, J., Huang, T., and
536 Wang, X. (2023). Generative multimodal models are in-context learners. *Computer Vision and*
537 *Pattern Recognition*.
- 538 Tan, H. and Bansal, M. (2019). Lxmert: Learning cross-modality encoder representations from
539 transformers. *arXiv preprint arXiv:1908.07490*.
- 540 Team, C. (2024). Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:*
541 *2405.09818*.
- 542 Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal,
543 N., Hambro, E., Azhar, F., et al. (2023). Llama: Open and efficient foundation language models.
544 *arXiv preprint arXiv:2302.13971*.
- 545 Wang, W., Lv, Q., Yu, W., Hong, W., Qi, J., Wang, Y., Ji, J., Yang, Z., Zhao, L., Song, X., et al. (2023).
546 Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*.
- 547 Xu, X., Wu, C., Rosenman, S., Lal, V., Che, W., and Duan, N. (2023). Bridgetower: Building
548 bridges between encoders in vision-language representation learning. In *Proceedings of the AAAI*
549 *Conference on Artificial Intelligence*, volume 37, pages 10637–10647.
- 550 Yin, Z., Wang, J., Cao, J., Shi, Z., Liu, D., Li, M., Huang, X., Wang, Z., Sheng, L., Bai, L.,
551 et al. (2024). Lamm: Language-assisted multi-modal instruction-tuning dataset, framework, and
552 benchmark. *Advances in Neural Information Processing Systems*, 36.
- 553 Yu, W., Yang, Z., Li, L., Wang, J., Lin, K., Liu, Z., Wang, X., and Wang, L. (2023). Mm-vet:
554 Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*.

555 **Appendices**

556 **A Dataset Details**

557 **A.1 Ground-Truth Image-Intention Pair Generation**

558 The section contains the essential information used to scale the ground-truth image-intention pair
559 generation process. Below, we detail key design choices and procedures.

560 **Scenario and Task Selection** - Scenarios and tasks with repetitive behaviors (e.g., dancing, instru-
561 ments playing) are excluded. The Table 3 lists all scenarios and tasks considered.

Table 3: Scenario and Applicable Tasks

Scenario	Applicable_task_name
Bike Repair	Install a Wheel, Remove a Wheel, Fix a Flat Tire - Replace a Bike Tube, Clean and Lubricate the Chain
CPR	First Aid - CPR
Covid Test	Covid-19 Rapid Antigen Test
Cooking	Making Cucumber & Tomato Salad, Making Greek Salad, Making Sesame-Ginger Asian Salad, Making Chai Tea, Making a Milk Tea, Cooking Noodles, Cooking an Omelet, Cooking Scrambled Eggs, Cooking Tomato & Eggs, Cooking Dumplings, Cooking Pasta, Cooking Sushi Rolls, Cooking Samosas, Making Greek Salad, Making White Radish & Lettuce & Tomato & Cucumber Salad

562 **Intention Definition and Keywords Mapping** - For each selected scenario, we define a set of
563 high-level intentions (Table 4). We apply a two-stage matching process:

- 564 1. For each take, we extract all action-level narrations and compute cosine similarity between
565 narration sentences and the keyword list associated with each intention (Table 5).
- 566 2. From each take, we select up to three frames (from the annotated *best_exo* camera) with the
567 highest similarity scores for each intention, ensuring a minimum 10-second separation to
568 avoid look-alike images. These are used as first-pass image-intention candidates.

Table 4: Scenarios and Associated Intentions

Scenario	Intention
Bike Repair	Install a wheel
	Replace the tire tube on the wheel
	Clean and lubricate the chain
	Remove a wheel
CPR	Confirm patient consciousness
	Call for help
	Press for heart rate
Covid Test	Set up for test
	Understand instruction
	Perform test
Cooking	Prepare ingredient
	Preheat pan for cooking
	Add flavor to dish
	Clean up work station

Table 5: Intention to Keywords Mapping

Intention	Keywords
Install a wheel	install, attach, bike fork
Replace the tire tube on the wheel	tire level, tire valve, inflate/deflate, tire tube, bike inner tube, fit the bike tire
Clean and lubricate the chain	chain lube, degreaser spray, lubricant bottle, hold the towel, clean the chain, pick up a brush, spray water
Remove a wheel	removes the bicycle wheel, removes the wheel, take off wheel
Confirm patient consciousness	pat, check for breathing, observe, tap
Call for help	wave her hands, extend right hand, extend left hand, call for help
Press for heart rate	interlace the fingers of this hands, compress, interlock, press
Set up for test	put on desk, place on desk, pick out from box, set up, open the box
Understand instruction	test manual, test instruction, read, understand, flip
Perform test	insert test swab, pick up the collection swab, dip the swab, nostril, nose
Prepare ingredient	chopping board, tomato, onion, scallion, knife, cut, carrot, potato, banana
Heat pan for cooking	press a switch, take the skillet, turn on heat, adjust the heat, turn on gas stove, picks the frying pan
Add flavor to dish	pick up black pepper, pick up the salt, soy sauce, sauce, sugar
Clean up work station	wash, turns on the tap, opens the tap, waste bins, push dirt into sink hole, picks the dirt, trash can

569 **Confounding Distractors** - As shown in Table 6, for some intentions, we define the confounding
570 distractors that are either visually similar with or sequentially entailed to each other, and avoid
571 presenting them within the same question.

Table 6: Intention and Confounding Distractor Pairs

Intention	Confounding Distractor
Install a wheel	Remove a wheel
Remove a wheel	Install a wheel
Confirm patient consciousness	Press for heart rate
Press for heart rate	Confirm patient consciousness
Set up for test	Perform test
Understand instruction	Set up for test
Perform test	Set up for test
Prepare ingredient	Clean up work station
Clean up work station	Prepare ingredient

572 **LLM Validation** - We then use GPT-4o to validate each image-intention pair.

573 Sample Prompt:



Figure 6: Sample Image Input for LLM Qtype4 Distractor Generation - Cooking

- 574 - I will provide an image of a person performing an action related to *Cooking* (*note:*
 575 *Scenario*), and a phrase that tries to describe the intention of the person: "*Add*
 576 *flavor to dish*" (*note: Intention*). Return only the required strings in a list format
 577 based on the following instructions, without additional explanations.
- 578 - Return 'great' if you are confident that the phrase accurately describes the intention
 579 of the person in the image.
 - 580 - Return 'good' if you think the phrase describes the intention, but not as confidently.
 - 581 - Return 'wrong' if the phrase is unrelated to the image, is not the intention that a
 582 normal non-technical human viewer could infer from the image, or has a better
 583 alternative from the following list: [*Prepare ingredients, Clean up work station,*
 584 *Add flavor to dish, Preheat pan for cooking*] (*note: All intentions in the scenario*).
 - 585 - If you choose 'wrong', also return the best alternative option from the list. If none
 586 of the alternatives work, return 'None'.

587 **A.2 Qtype 3 Question Generation**

588 We utilize GPT-3o to scale the question generation process for Qtype3. Below documents the detailed
 589 prompt we provide to the LLM.

590 **Context**

591 You will receive one or more third-person photos of everyday scenes. Each image contains:

- 592 1. a **red gaze line** that starts at the eyes of the **primary person** (the "subject"), and
- 593 2. several clearly identifiable objects.

594 Your task is to write **perspective-based multiple-choice questions (MCQs)** that test spatial reasoning
 595 **from the subject's viewpoint** (not the camera's).

596 **MCQ Templates**

- 597 • Type: Visibility - From the perspective of SUBJECT, which of the following items in the
 598 image are visible?
- 599 • Type: Direction - From the perspective of SUBJECT, in which direction is TARGET-
 600 OBJECT?
- 601 • Type: Leftmost/Rightmost - From the perspective of SUBJECT, which of the following
 602 items appears leftmost / rightmost?

603 Note on choices: All options must be generic and unambiguous (e.g., "a red box on the counter"
 604 rather than "a toolbox"). Label the correct answer A–D.

605 **Workflow**

606
607
608
609
610
611
612
613
614
615
616
617

1. **Load the image**

- (a) Note the general setting (kitchen, bike workshop, etc.).
- (b) Locate the subject (person with the red line).
- (c) **Determine subject orientation** — choose exactly one:
 - facing-camera
 - back-to-camera
 - profile-left (subject looking toward **camera-left**)
 - profile-right (subject looking toward **camera-right**)

If the body is roughly 45°, combine them, such as facing-camera & profile-right

- (d) **Build a subject-centric frame**
 - **Forward** = the red gaze line.
 - **Left / Right** = rotate the frame $\pm 90^\circ$ around the subject.

Subject Orientation	Subject-Left	Subject-Right	Quick Visual Cue
facing-camera	camera- right	camera- left	(mirror rule)
back-to-camera	camera- left	camera- right	(mirror rule)
profile-left	down in photo	up in photo	
profile-right	up in photo	down in photo	

618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640

- **Behind** = opposite of forward.
- If subject orientation is combined (e.g., facing-camera & profile-right), the projection should also be combined.

2. **Parse objects**

List every salient object as *minimal-adjective + generic noun* (e.g., “blue mug,” “metal faucet”). Re-use these exact names in the MCQs.

3. **Generate three MCQs (one of each type) per image**

- Describe the subject succinctly (e.g., “the woman in a blue apron”).
- **Direction**: pick a clear {TARGET-OBJECT}; options = front / behind / left / right.
- **Visibility & Leftmost/Rightmost**: provide four distinct objects.
- Mark the correct answer.

4. **Quality check (mandatory)**

- Verify every spatial relation in the subject-centric frame.
- Ensure wording is concise, bias-free, and each referenced object is clearly visible.

5. **Output** — one JSON record per question. {

```
"image_id": "<image filename or UID>",
"subject_direction": "facing-camera | back-to-camera | profile-left | profile-right | <combined>",
"question_type": "visibility | direction | leftmost | rightmost",
"question": "<full question text>",
"options": "A": "...", "B": "...", "C": "...", "D": "...",
"answer_key": "A/B/C/D"
}
```

641 **A.3 Qtype 4 Distractor Generation**

642 The distractor generation process for Qtype 4 requires special attention due to its textual nature.

643 For **Wrong Intention** distractor type, we randomly sample other intentions from the same scenario, while explicitly avoiding confounding distractors (Table 6). When the number of suitable alternatives is insufficient, we supplement the set with manually created pseudo-intentions that are plausible yet not part of our dataset (e.g. *Taste the food, Throw away food waste*).

647 For **Wrong Action** and **Correct Action** distractor types, we leverage a LLM (GPT-4o) to scale generation and validation.

649 Sample Prompt:



Figure 7: Sample Image Input for LLM Ground-Truth Validation - Cooking

650 You are an expert in linguistics and are good at coming up natural alternative
 651 expression if given a sentence in English.

652 Give the sentence '*C takes the dark soy sauce with his right hand.*', please come
 653 up with the following, without including any explanations.

654 1. Type 3: 5 concise phrases that describe the action (atomic description) in the
 655 sentence. If the sentence doesn't have 'C' (a human) as the subject, make sure to
 656 phrase the action such that it sounds reasonable if the subject is a human.

657 2. Type 2: 5 concise phrases that describe different but similar actions. For example,
 658 these alternate phrases can EITHER a) describe the same action on a different
 659 object, OR b) describe different action on the same object. Do not replace both
 660 action and object at the same time. It is preferred that if a human is to perform these
 661 phrases, their body gestures and/or scenario will look like the original sentence.

662 General requests:

663 1. return phrases without explicit subject. For example, 'C does something' should
 664 be shortened to 'do something'.

665 2. the phrases should use verbs and nouns that are natural and colloquial.

666 3. the phrases should make sense with human as the subject, even if the subject in
 667 original sentence may not be a human. Rephrase the original sentence to human-
 668 subject first, then generate alternatives.

669 The output format should follow: {'type_3': [phrases1, phrases2, ...], 'type_2':
 670 [phrases1, phrases2, ...]}

671 Sample Output:

672 {'type_3': ['grab soy sauce', 'hold dark soy', 'pick up sauce', 'lift dark soy', 'take
 673 soy bottle'], 'type_2': ['grab light soy sauce', 'hold ketchup bottle', 'pick up olive
 674 oil', 'lift sesame oil', 'take vinegar bottle']}

675 **B NeurIPS Paper Checklist**

676 The checklist is designed to encourage best practices for responsible machine learning research,
677 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
678 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
679 follow the references and follow the (optional) supplemental material. The checklist does NOT count
680 towards the page limit.

681 Please read the checklist guidelines carefully for information on how to answer these questions. For
682 each question in the checklist:

- 683 • You should answer [\[Yes\]](#) , [\[No\]](#) , or [\[NA\]](#) .
- 684 • [\[NA\]](#) means either that the question is Not Applicable for that particular paper or the
685 relevant information is Not Available.
- 686 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

687 **The checklist answers are an integral part of your paper submission.** They are visible to the
688 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it
689 (after eventual revisions) with the final version of your paper, and its final version will be published
690 with the paper.

691 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
692 While "[\[Yes\]](#)" is generally preferable to "[\[No\]](#)", it is perfectly acceptable to answer "[\[No\]](#)" provided a
693 proper justification is given (e.g., "error bars are not reported because it would be too computationally
694 expensive" or "we were unable to find the license for the dataset we used"). In general, answering
695 "[\[No\]](#)" or "[\[NA\]](#)" is not grounds for rejection. While the questions are phrased in a binary way, we
696 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
697 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
698 supplemental material, provided in appendix. If you answer [\[Yes\]](#) to a question, in the justification
699 please point to the section(s) where related material for the question can be found.

700 IMPORTANT, please:

- 701 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 702 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 703 • **Do not modify the questions and only use the provided macros for your answers.**

704 1. **Claims**

705 Question: Do the main claims made in the abstract and introduction accurately reflect the
706 paper’s contributions and scope?

707 Answer: [\[Yes\]](#)

708 Justification: Both the abstract and introduction crystallize the curation structure, use case,
709 and scaling goal of our benchmark with literature, statistics, and procedural visualizations.

710 Guidelines:

- 711 • The answer NA means that the abstract and introduction do not include the claims
712 made in the paper.
- 713 • The abstract and/or introduction should clearly state the claims made, including the
714 contributions made in the paper and important assumptions and limitations. A No or
715 NA answer to this question will not be perceived well by the reviewers.
- 716 • The claims made should match theoretical and experimental results, and reflect how
717 much the results can be expected to generalize to other settings.
- 718 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
719 are not attained by the paper.

720 2. **Limitations**

721 Question: Does the paper discuss the limitations of the work performed by the authors?

722 Answer: [\[Yes\]](#) .

723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775

Justification: The paper discusses several limitations, including reliance on synthetic prompts that may not generalize across all domains, limited access to closed-source models for full comparison, and the challenge of verifying whether LLMs perform true mental simulation or merely pattern match.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [No]

Justification: The paper does not present theoretical results or formal proofs. It is primarily an empirical benchmark contribution.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

776 Justification: The paper includes detailed descriptions of the dataset, the benchmark structure,
777 prompt generation pipelines, and evaluation procedures. All steps required to reproduce the
778 experimental setup are clearly documented in both the main paper and appendix

779 Guidelines:

- 780 • The answer NA means that the paper does not include experiments.
- 781 • If the paper includes experiments, a No answer to this question will not be perceived
782 well by the reviewers: Making the paper reproducible is important, regardless of
783 whether the code and data are provided or not.
- 784 • If the contribution is a dataset and/or model, the authors should describe the steps taken
785 to make their results reproducible or verifiable.
- 786 • Depending on the contribution, reproducibility can be accomplished in various ways.
787 For example, if the contribution is a novel architecture, describing the architecture fully
788 might suffice, or if the contribution is a specific model and empirical evaluation, it may
789 be necessary to either make it possible for others to replicate the model with the same
790 dataset, or provide access to the model. In general, releasing code and data is often
791 one good way to accomplish this, but reproducibility can also be provided via detailed
792 instructions for how to replicate the results, access to a hosted model (e.g., in the case
793 of a large language model), releasing of a model checkpoint, or other means that are
794 appropriate to the research performed.
- 795 • While NeurIPS does not require releasing code, the conference does require all submis-
796 sions to provide some reasonable avenue for reproducibility, which may depend on the
797 nature of the contribution. For example
 - 798 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
799 to reproduce that algorithm.
 - 800 (b) If the contribution is primarily a new model architecture, the paper should describe
801 the architecture clearly and fully.
 - 802 (c) If the contribution is a new model (e.g., a large language model), then there should
803 either be a way to access this model for reproducing the results or a way to reproduce
804 the model (e.g., with an open-source dataset or instructions for how to construct
805 the dataset).
 - 806 (d) We recognize that reproducibility may be tricky in some cases, in which case
807 authors are welcome to describe the particular way they provide for reproducibility.
808 In the case of closed-source models, it may be that access to the model is limited in
809 some way (e.g., to registered users), but it should be possible for other researchers
810 to have some path to reproducing or verifying the results.

811 5. Open access to data and code

812 Question: Does the paper provide open access to the data and code, with sufficient instruc-
813 tions to faithfully reproduce the main experimental results, as described in supplemental
814 material?

815 Answer: [Yes]

816 Justification: The dataset and code will be made available upon the sharing of private
817 link. Instructions for dataset usage and evaluation are provided in paper body and in the
818 supplemental material. The paper also includes a structured overview of asset preparation.
819 Further annotation and usage instructions will be added upon official release to the public.

820 Guidelines:

- 821 • The answer NA means that paper does not include experiments requiring code.
- 822 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
823 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 824 • While we encourage the release of code and data, we understand that this might not be
825 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
826 including code, unless this is central to the contribution (e.g., for a new open-source
827 benchmark).
- 828 • The instructions should contain the exact command and environment needed to run to
829 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
830 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.

- 831
- 832
- 833
- 834
- 835
- 836
- 837
- 838
- 839
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
 - The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
 - At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
 - Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

840 6. Experimental setting/details

841 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
842 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
843 results?

844 Answer: [Yes]

845 Justification: The paper specifies all relevant experimental configurations, including prompt
846 types, number of frames per input, prompting strategies, and accuracy calculation. Full
847 prompting templates are included in the appendix.

848 Guidelines:

- 849
- 850
- 851
- 852
- 853
- The answer NA means that the paper does not include experiments.
 - The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
 - The full details can be provided either with the code, in appendix, or as supplemental material.

854 7. Experiment statistical significance

855 Question: Does the paper report error bars suitably and correctly defined or other appropriate
856 information about the statistical significance of the experiments?

857 Answer: [Yes]

858 Justification: We performed statistical significance testing and showed the results in 4

859 Guidelines:

- 860
- 861
- 862
- 863
- 864
- 865
- 866
- 867
- 868
- 869
- 870
- 871
- 872
- 873
- 874
- 875
- 876
- 877
- 878
- 879
- The answer NA means that the paper does not include experiments.
 - The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
 - The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
 - The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
 - The assumptions made should be given (e.g., Normally distributed errors).
 - It should be clear whether the error bar is the standard deviation or the standard error of the mean.
 - It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

880 8. Experiments compute resources

881 Question: For each experiment, does the paper provide sufficient information on the com-
882 puter resources (type of compute workers, memory, time of execution) needed to reproduce
883 the experiments?

884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933

Answer: [Yes]

Justification: We explained our experiment setup and inference details in Section 4.1

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: We made sure to conform with the Code of Ethics. Please see the below bullet points for more explanation.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: This paper presents a foundational benchmark for evaluating multimodal reasoning in machine learning models. It does not propose a deployed system or application, nor does it yet involve real-world users or decision-making contexts. As such, the work does not present direct societal impacts—positive or negative—in its current form.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.

- 934
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
- 935
936
937

938 11. Safeguards

939 Question: Does the paper describe safeguards that have been put in place for responsible
940 release of data or models that have a high risk for misuse (e.g., pretrained language models,
941 image generators, or scraped datasets)?

942 Answer: [NA]

943 Justification: The constructed dataset contains only predefined and reviewed output, and
944 thus poses no such risks,

945 Guidelines:

- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.
- 946
947
948
949
950
951
952
953
954
955

956 12. Licenses for existing assets

957 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
958 the paper, properly credited and are the license and terms of use explicitly mentioned and
959 properly respected?

960 Answer: [Yes]

961 Justification: We credited and introduced the dataset we use in the paper in the Section 3.1
962 and in References. All authors have been granted licenses to download and use the dataset
963 from the official website.

964 Guidelines:

- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.
- 965
966
967
968
969
970
971
972
973
974
975
976
977
978
979

980 13. New assets

981 Question: Are new assets introduced in the paper well documented and is the documentation
982 provided alongside the assets?

983 Answer: [Yes]

984 Justification: The benchmark dataset and code will be made available upon the sharing
985 of private link. Instructions for dataset usage and evaluation are provided in paper body
986 and in the supplemental material. The paper also includes a structured overview of asset
987 preparation. Further annotation and usage instructions will be added upon official release to
988 the public.

989 Guidelines:

- 990 • The answer NA means that the paper does not release new assets.
- 991 • Researchers should communicate the details of the dataset/code/model as part of their
992 submissions via structured templates. This includes details about training, license,
993 limitations, etc.
- 994 • The paper should discuss whether and how consent was obtained from people whose
995 asset is used.
- 996 • At submission time, remember to anonymize your assets (if applicable). You can either
997 create an anonymized URL or include an anonymized zip file.

998 14. Crowdsourcing and research with human subjects

999 Question: For crowdsourcing experiments and research with human subjects, does the paper
1000 include the full text of instructions given to participants and screenshots, if applicable, as
1001 well as details about compensation (if any)?

1002 Answer: [NA]

1003 Justification: The paper does not involve crowdsourcing nor research with human subjects

1004 Guidelines:

- 1005 • The answer NA means that the paper does not involve crowdsourcing nor research with
1006 human subjects.
- 1007 • Including this information in the supplemental material is fine, but if the main contribu-
1008 tion of the paper involves human subjects, then as much detail as possible should be
1009 included in the main paper.
- 1010 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1011 or other labor should be paid at least the minimum wage in the country of the data
1012 collector.

1013 15. Institutional review board (IRB) approvals or equivalent for research with human 1014 subjects

1015 Question: Does the paper describe potential risks incurred by study participants, whether
1016 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1017 approvals (or an equivalent approval/review based on the requirements of your country or
1018 institution) were obtained?

1019 Answer: [NA]

1020 Justification: This paper does not involve crowdsourcing nor research with human subjects.

1021 Guidelines:

- 1022 • The answer NA means that the paper does not involve crowdsourcing nor research with
1023 human subjects.
- 1024 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1025 may be required for any human subjects research. If you obtained IRB approval, you
1026 should clearly state this in the paper.
- 1027 • We recognize that the procedures for this may vary significantly between institutions
1028 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1029 guidelines for their institution.
- 1030 • For initial submissions, do not include any information that would break anonymity (if
1031 applicable), such as the institution conducting the review.

1032 16. Declaration of LLM usage

1033 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1034 non-standard component of the core methods in this research? Note that if the LLM is used
1035 only for writing, editing, or formatting purposes and does not impact the core methodology,
1036 scientific rigor, or originality of the research, declaration is not required.

1037

Answer: [Yes]

1038

Justification: We use LLMs for generating our VQA questions. We documented the detailed procedures and prompts in Sections A.1, A.2, and A.3

1039

1040

Guidelines:

1041

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.

1042

1043

- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.

1044