

ESI-Bench: Towards Embodied Spatial Intelligence that Closes the Perception-Action Loop

Yining Hong¹ Jiageng Liu² Han Yin¹ Manling Li³ Leonidas Guibas¹ Fei-Fei Li¹ Jiajun Wu¹ Yejin Choi¹
¹Stanford University ²UCLA ³Northwestern University

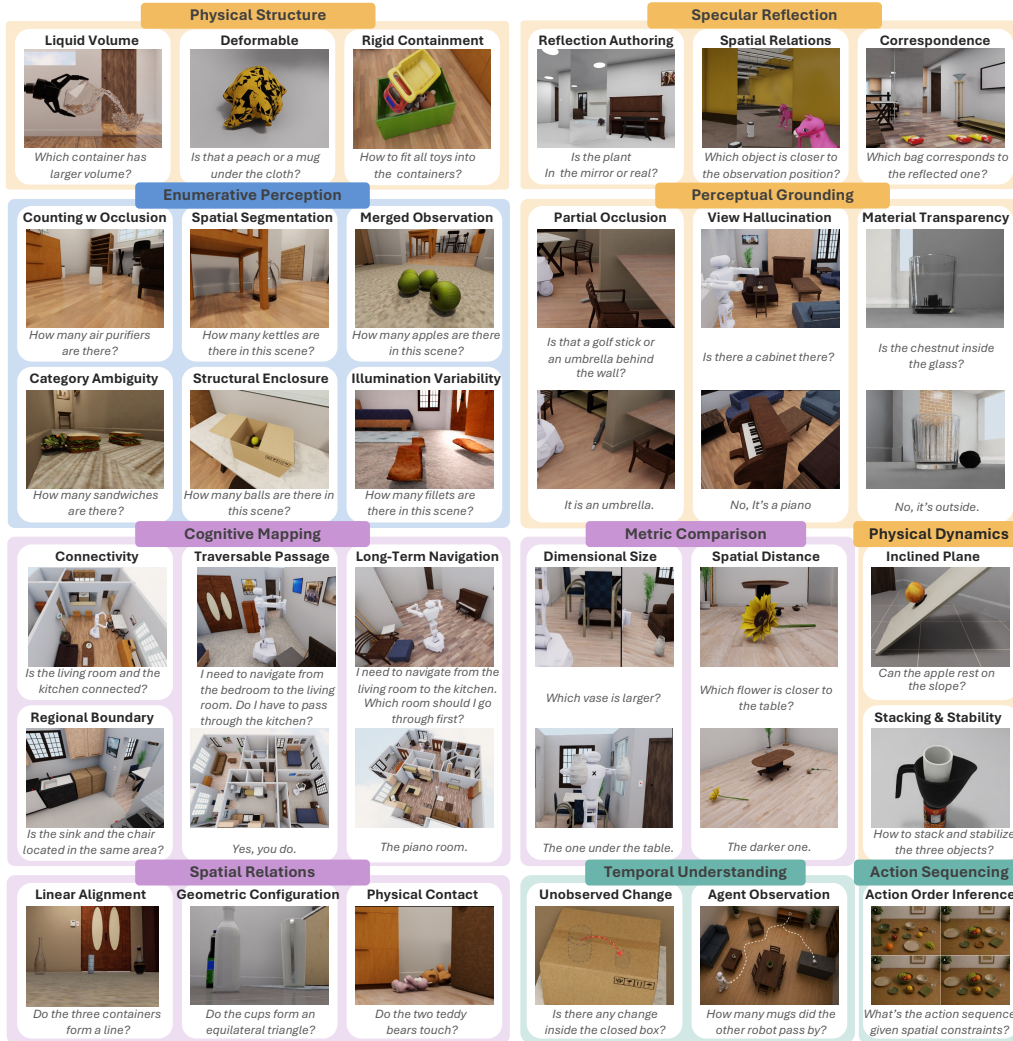


Figure 1: **ESI-Bench** is a comprehensive benchmark for embodied spatial intelligence, spanning 10 task categories and 30 subcategories organized around Spelke’s four core knowledge systems [Spelke and Kinzler, 2007]: **object persistence**, **layout and geometry**, **number representation**, and **agents and goal-directed actions**.

Abstract

1 Spatial intelligence unfolds through a perception–action loop: agents act to acquire
 2 observations, and reason about how observations vary as a function of action.
 3 Rather than passively processing what is seen, they actively uncover what is
 4 unseen—occlusion, dynamics, containment, and functionality—beyond the reach

5 of passive sensing. We take a step beyond prior formulations of spatial intelligence,
6 which often emphasize passive perception or assume access to oracle observations,
7 by recasting the observer as an actor. We introduce ESI-BENCH, a comprehensive
8 benchmark for embodied spatial intelligence spanning 10 task categories and
9 30 subcategories built on OmniGibson, grounded in Spelke’s core knowledge
10 systems. Agents must decide what abilities to deploy — perception, locomotion,
11 and manipulation — and how to act to answer questions that cannot be resolved
12 from passive observation alone. We conduct extensive experiments on state-of-
13 the-art MLLMs and find that active exploration substantially outperforms passive
14 counterparts, with agents spontaneously discovering emergent spatial strategies
15 without explicit instruction, while passive multi-view adds noise rather than signal
16 despite consuming far more images. Most failures stem not from weak perception
17 but from action blindness, and their coupling drives cascading failures where
18 bad actions produce bad views which produce worse actions. While explicit 3D
19 grounding stabilizes reasoning on depth-sensitive tasks, imperfect reconstruction
20 proves more harmful than 2D baselines by actively distorting spatial relations.
21 Human studies further reveal that unlike humans who seek falsifying viewpoints
22 and revise beliefs under contradiction, models commit prematurely with high
23 confidence regardless of evidence quality, exposing a metacognitive gap that neither
24 better perception nor more embodied interaction alone can close. All data, codes
25 are publicly available at our project page <https://esi-bench.github.io/>.

26 1 Introduction

27 Perception is often characterized in cognitive science as *perceptually guided action* [Varela et al.,
28 1991, Gibson, 1979]: a perception–action loop where knowing how observations change as a function
29 of action [O’Regan and Noë, 2001], and which actions elicit effective and informative sensing, are
30 often more challenging than sensing itself. This is especially critical in spatial intelligence, which
31 concerns not only what is *seen*, but also what is *unseen*: latent physical properties such as occlusion,
32 dynamics, containment, and functionality that are inaccessible to passive sensing and must be actively
33 unveiled through interaction.

34 **selective or prioritized sensing? sentences not parsable: perception-action loop where; also which**
35 **predominantly treat make (1) and (2) more different from passive sensing to proactive sensing mirage**
36 **-> hallucination. maybe not incomplete. ambiguity teaser refer to figure epistemic calibration -> spell**
37 **it out. explain more** We take a step beyond prior formulations of spatial intelligence Liu et al. [2023],
38 Yang et al. [2024, 2025d], which predominantly treat perception as passive or assume access to oracle
39 observations, by recasting the observer as an actor. Our work contrasts with prior works in three key
40 ways: (1) *from spatial sensing to spatial competence*, where agents are evaluated not only on what
41 they can perceive, but on whether they know how to act to perceive it; (2) *selective sensing*, where
42 agents must determine which observations are worth acquiring, prioritizing task-relevant information
43 over redundant or uninformative inputs; and (3) *resolving perceptual mirages*, where agents must
44 reason through incomplete or misleading observations to infer the hidden spatial structures and
45 underlying physical constraints beyond what is directly observed.

46 We introduce ESI-BENCH, the first comprehensive benchmark for embodied spatial intelligence
47 spanning 10 task categories and 30 subcategories, to address the critical perception-action gap in
48 existing benchmarks by focusing on questions unanswerable from passive observation. Our category
49 design follows Spelke’s core knowledge systems [Spelke and Kinzler, 2007], which identify four
50 faculties of spatial intelligence: object persistence, layout and geometry, number representation, and
51 agents and intentional action. Building on these theoretical foundations, we conduct human surveys
52 to identify the most challenging spatial tasks that require embodied interaction and manipulation
53 within each faculty, narrowing these into a structured taxonomy spanning diverse forms of spatial
54 reasoning, as illustrated in the teaser across all 10 categories and 30 subcategories. These tasks only
55 become meaningful when an agent has a body, a belief state, and physical stakes in the outcome: the
56 agent must determine what abilities to deploy (perception, locomotion, manipulation), which actions
57 to take (where to move, what to probe, how to manipulate), and how to execute them in the right
58 order.

59 We conduct extensive experiments on state-of-the-art MLLMs across three paradigms: passive
60 single-view, passive multi-view, and active exploration, alongside a ground-truth oracle that separates
61 perception errors from action errors. Our experiments reveal multiple key insights: (1) *active explo-*
62 *ration unlocks emergent spatial strategies*: without explicit instruction, active agents spontaneously
63 discover diverse action compositions, driving substantial gains over passive counterparts while passive
64 multi-view, despite consuming far more images, adds noise rather than signal; (2) *action blindness*
65 *dominates perceptual blindness*: for most tasks the bottleneck is not perception but action selection,
66 as models given oracle viewpoints succeed dramatically, yet certain tasks expose a hard perceptual
67 ceiling where even perfect viewpoints cannot overcome the fundamental limits of 2D recognition;
68 and (3) *failures cascade and compound*: suboptimal actions produce uninformative views, which
69 trigger worse subsequent actions, creating a compounding chain of errors that cannot be recovered
70 within the step budget.

71 We further investigate whether explicit 3D representations can help. We find that while explicit 3D
72 grounding stabilizes reasoning on depth-sensitive tasks by recovering information that 2D projections
73 fundamentally lose, imperfect reconstructions prove more harmful than 2D baselines, as geometric
74 artifacts actively distort fine-grained spatial relations and mislead downstream reasoning. Finally,
75 human studies expose a critical gap in epistemic calibration. We observe that unlike humans
76 who actively seek falsifying viewpoints, explore orthogonal angles, and revise their beliefs when
77 contradicted, models commit prematurely with uniformly high confidence regardless of evidence
78 quality, anchoring to first impressions and ignoring contradictory observations, a metacognitive
79 failure that neither better perception nor more embodied interaction alone can close.

80 2 Related Work

81 **Spatial Reasoning Methods in MLLMs.** Recent advances in multimodal large language models
82 have improved spatial understanding, yet most methods still focus on interpreting a fixed observation
83 set rather than choosing which observations to acquire. One line of work injects geometric priors into
84 otherwise 2D backbones: SpatialVLM Chen et al. [2024] constructs training data with synthesized
85 3D spatial annotations and metric-depth supervision; SpatialBot Cai et al. [2024] leverages RGB-D
86 inputs and depth-oriented question answering; SpatialRGPT Cheng et al. [2024] learns region-level
87 representations from 3D scene graphs with a depth plugin; and newer geometry-grounded models such
88 as Spatial-MLLM Wu et al. [2025a] and VLM-3R Fan et al. [2026] further strengthen 3D reasoning
89 by combining visual geometry priors with monocular video or reconstructive instruction tuning. A
90 second line of work treats spatial inference as an explicit reasoning process: SpatialCoT Liu et al.
91 [2025] aligns vision–language inputs with spatial coordinates for chain-of-thought spatial grounding,
92 while VILASR Wu et al. [2025b] interleaves textual reasoning with visual drawing operations.
93 Together, these methods improve how models reason over observed content, but the observation set
94 remains fixed as input. ESI-Bench instead asks whether a spatial reasoner can actively choose the
95 observations its perception modules need.

96 **Benchmarks for Spatial Reasoning.** Evaluation of spatial intelligence has scaled with model
97 capability, yet most benchmarks still operate under fixed observation settings. Early work such as
98 VSR Liu et al. [2023] formulates spatial reasoning as classification over a predefined set of spatial
99 relations from single images, while BLINK Fu et al. [2024] and 3DSRBench Ma et al. [2025] extend
100 evaluation to core visual perception and fine-grained 3D spatial reasoning. SpatialScore Wu et al.
101 [2026a] further consolidates VGBench and multiple datasets into a unified benchmark with a tool-
102 augmented evaluation agent. More recent benchmarks expand the observation space to multi-view and
103 temporal regimes, including egocentric video understanding (e.g., VSI-Bench Yang et al. [2025a]),
104 multi-image spatial consistency (e.g., MMSI-Bench Yang et al. [2025d]), and spatial mental modeling
105 from partial observations (e.g., MindCube Wang et al. [2026]). Cambrian-S Yang et al. [2025c]
106 pushes toward long-horizon settings with VSI-SUPER, targeting visual spatial recall and continual
107 spatial counting over extended video. In parallel, benchmarks such as PhysBench Chow et al. [2025]
108 and CausalSpatial Ma et al. [2026] probe latent physical structure and object-centric dynamics from
109 passive multimodal inputs. Across this progression, inputs become increasingly rich—multiple
110 images, videos, partial coverage, and temporal dynamics—yet the observation process remains
111 largely fixed, with limited support for agent-driven view selection. As a result, these benchmarks
112 provide limited signals for disentangling perceptual limitations from failures in active information

113 acquisition. ESI-Bench retains the diagnostic spirit of this line of work but treats the observer as an
 114 active agent, making the utility of each observation contingent on the model’s decisions.

115 **Embodied Evaluation and Active Perception.** A separate line of work places models inside
 116 simulators and evaluates them as embodied agents. EmbodiedBench Yang et al. [2025b] and Embod-
 117 iedEval Cheng et al. [2025] measure end-to-end task success across diverse navigation, interaction,
 118 and QA scenarios, while OpenEQA Majumdar et al. [2024]—spanning both episodic-memory (EM-
 119 EQA) and active (A-EQA) settings—and the more recent EXPRESS-Bench Jiang et al. [2025]
 120 formalize embodied question answering with explicit attention to exploration quality, the latter
 121 introducing an Exploration–Answer Consistency (EAC) metric to discourage disembodied reasoning.
 122 Spatial cognition under embodiment has been examined more directly by EmbSpatial-Bench Du
 123 et al. [2024] and ESPIRE Zhao et al. [2026], which provide diagnostic assessments of egocentric
 124 spatial reasoning, with ESPIRE further decomposing each task into localization and execution to
 125 enable fine-grained error analysis. On the side of active perception itself, several recent efforts
 126 develop methods that adaptively select what to observe: Vision in Action Xiong et al. [2025] learns
 127 active observation strategies from human demonstrations; Thinking in 360° Yu et al. [2025] studies
 128 humanoid head-rotation visual search over panoramic environments; and SaPaVe Liu et al. [2026a]
 129 and ActiveVLA Liu et al. [2026b] integrate viewpoint selection into vision–language–action policies
 130 for manipulation. The closest concurrent benchmark is CHAIN Wu et al. [2026b], which evaluates
 131 closed-loop physical reasoning over interlocking mechanical puzzles and 3D stacking and packing in
 132 a physics engine; CHAIN and ESI-Bench are complementary in scope—CHAIN focuses on tabletop
 133 physical reasoning, while ESI-Bench targets the broader range of spatial faculties an embodied agent
 134 must exercise across object, geometry, number, and agent reasoning, including latent properties such
 135 as containment, occlusion, transparency, reflection, and unobserved scene change that remain only
 136 partially covered by current embodied evaluation. Table 1 situates ESI-Bench within this landscape
 137 by combining active embodiment with hidden-state probing for spatial reasoning.

Table 1: Comparison of spatial reasoning and embodied benchmarks. **Action:** L = locomotion, P = active perception, M = manipulation. **Core Knowledge:** Obj = object, Geom = geometry, Num = number, Phys = physics, Agt = agent. Hidden-State Probing denotes whether the benchmark targets latent properties (e.g., occlusion, containment, dynamics) inaccessible to passive observation.

Benchmark	Active Obs.	Action	Hidden-State	Core Knowledge
VSR Liu et al. [2023]	✗	–	✗	Geom
BLINK Fu et al. [2024]	✗	–	✗	Geom
3DSRBench Ma et al. [2025]	✗	–	✗	Obj, Geom
VSI-Bench Yang et al. [2025a]	✗	–	✗	Obj, Geom, Num
MMSI-Bench Yang et al. [2025d]	✗	–	✗	Obj, Geom
MindCube Wang et al. [2026]	✗	–	✓	Geom, Phys
PhysBench Chow et al. [2025]	✗	–	✓	Obj, Geom, Phys
CausalSpatial Ma et al. [2026]	✗	–	✓	Obj, Phys
EmbSpatial-Bench Du et al. [2024]	✗	–	✗	Geom
OpenEQA Majumdar et al. [2024]	✓	L	✗	Obj, Geom
EmbodiedBench Yang et al. [2025b]	✓	L, M	✗	Obj, Geom, Agt
EmbodiedEval Cheng et al. [2025]	✓	L, M	✗	Obj, Geom, Agt
EXPRESS-Bench Jiang et al. [2025]	✓	L	✗	Geom, Agt
ESPIRE Zhao et al. [2026]	✓	L, M	✗	Obj, Geom, Agt
CHAIN Wu et al. [2026b]	✓	M	✓	Obj, Geom, Phys
ESI-Bench (ours)	✓	L, P, M	✓	Obj, Geom, Num, Phys, Agt

138 3 ESI-BENCH

139 3.1 Benchmark Setup

140 **Task Definition** Each task in ESI-Bench is defined by a tuple $(\mathcal{S}, p_0, q, y^*)$, where \mathcal{S} is a 3D scene
 141 instantiated from the BEHAVIOR-1K scene pool with a fixed set of pre-loaded objects, p_0 is the

142 initial pose of the agent within the scene, q is a natural language question about a spatial property of
143 the scene, and y^* is the ground-truth answer. The environment is formalized as $\mathcal{E} = \langle \mathcal{S}, \mathcal{A}, \mathcal{O}, T \rangle$,
144 where \mathcal{A} is the action space, \mathcal{O} is the egocentric visual observation space, and $T : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ is
145 the transition function governing scene evolution. Given (\mathcal{S}, p_0, q) , the agent receives observation
146 $o_t \in \mathcal{O}$ at each timestep t , issues action $a_t \in \mathcal{A}$, and transitions to a new scene state via T , producing
147 a trajectory $\tau = (o_0, a_0, o_1, a_1, \dots)$ until it commits to a final answer \hat{y} . The agent produces \hat{y} as
148 free-form text; while the answer space is open, the phrasing of q implicitly suggests the expected
149 format: binary yes/no for relational tasks, a category name for comparative tasks, an integer count for
150 enumerative tasks, an ordering for procedural tasks, and so on. A response is correct if $\hat{y} = y^*$.

151 **Simulation Environment.** We build ESI-BENCH on BEHAVIOR-1K within the OmniGibson
152 simulator. BEHAVIOR-1K provides 51 fully interactive 3D scenes spanning residential, commercial,
153 and institutional environments (houses, gardens, offices, restaurants, schools, and stores), totaling
154 over 300 rooms across all scenes, with over 9,000 object instances across 1,829 categories, each
155 annotated with rich physical properties including friction, mass, and articulation. OmniGibson
156 is built on top of NVIDIA’s Isaac Sim and PhysX 5, providing the broad simulation capabilities
157 essential for embodied spatial evaluation: rigid-body contact physics, particle-based fluid simulation,
158 transparency rendering, realistic lighting and reflections, and extended object states including fill
159 levels, and toggled states. For each task instance (described in the next paragraph), we randomly
160 sample a scene from the BEHAVIOR-1K scene pool and select rooms based on room type and task
161 category requirements, determined from a combined room-object list. We load the selected room
162 into OmniGibson, allow the physics simulation to settle, and query the simulator state to extract a
163 structured scene graph containing per-object bounding boxes, categories, spatial relationships, room
164 assignments, and object states such as fillable capacity, toggled state, and contact flags. This scene
165 graph serves as the foundation for scenario construction: it provides the object inventory from which
166 task-relevant objects are selected, the spatial layout from which camera positions are computed, and
167 the ground-truth geometric state and object states from which labels are automatically derived.

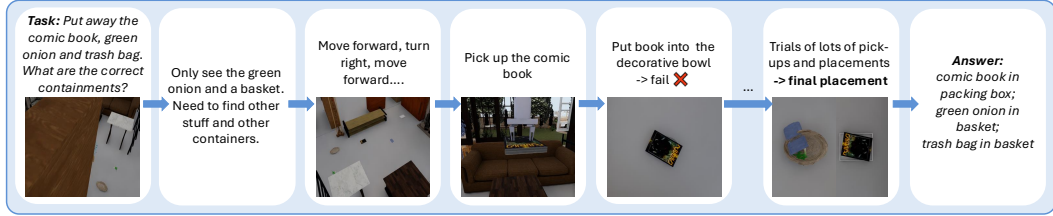
168 **Action Space.** Active agents operate over a unified action vocabulary spanning perception, locomotion
169 and manipulation, summarized in Table ???. We impose a maximum step budget of $T_{\max} = 30$.

170 3.2 Task Construction

171 **Task Proposal.** GPT-4o is prompted with the scene graph alongside task category requirements to
172 select task-relevant objects from a random sample of 200 candidate categories drawn from the full
173 BEHAVIOR-1K inventory, applying task-specific physical criteria to choose the most appropriate
174 categories and resolve a specific model instance per category. Beyond object selection, GPT-4o also
175 determines the initial positions of both the objects and the agent within the scene, and generates a
176 ground-truth action trajectory providing the optimal sequence of actions needed to resolve the task.
177 The selected objects and their spatial configuration implicitly define the task, with the ground-truth
178 answer y^* derived directly from the resulting scene state.

179 **Scene Instantiation.** Taking the GPT-4o-proposed object selections and initial positions as input,
180 selected objects are loaded into the scene at data generation time. Before placement, a conflict check
181 is performed against existing scene objects via bounding box intersection tests to ensure no overlap
182 with pre-existing scene content. Objects are then placed on supporting surfaces via physics-based
183 kinematic sampling and allowed to settle under simulation for a fixed number of steps. After settling,
184 the configuration is checked through a battery of tests: placement stability via bounding box re-
185 querying, per-view object existence verification via segmentation masks, and contact flag validation
186 where applicable. Configurations that fail any check are rejected.

187 **Agent Trajectory Collection.** The agent is initialized at the GPT-4o-proposed pose, determined
188 by a combination of randomization and pre-defined placement principles designed to intentionally
189 withhold the full scene configurations and properties from the agent. At initialization, the same battery
190 of checks applied during scene instantiation is performed: per-view object existence verification
191 via segmentation masks, bounding box re-querying, and contact flag validation where applicable.
192 The proposed actions are then executed step by step in the environment, with each step rendered to
193 produce egocentric visual observations and verified through the same battery of checks. Trajectories
194 where any step fails verification are discarded.



Action	Description
<i>Locomotion</i>	
move_forward / move_backward	Translate agent along viewing axis
move_left / move_right	Translate agent laterally
move_up / move_down	Translate agent vertically
<i>Perception</i>	
turn_left / turn_right	Rotate agent horizontally
turn_up / turn_down	Rotate agent vertically
<i>Manipulation</i>	
pick_up obj	Pick up object
put obj inside obj	Put object inside object
put obj on obj	Put object on object
fill obj with water	Fill object with water
pour obj from obj to obj	Pour object from object to object
<i>Terminal</i>	
answer(\hat{y} , c)	Commit to final answer with confidence c

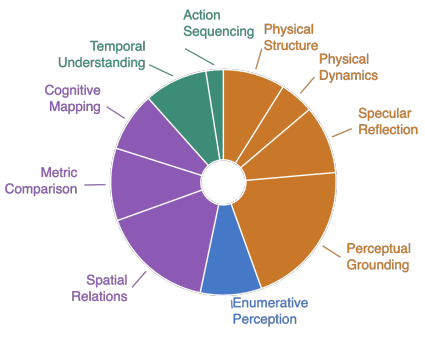


Figure 2: **Top:** Dataset example. **Bottom-left:** ESI-BENCH agent action space. **Bottom-right:** Distribution of task categories.

195 **Metadata Saving.** Upon completion of agent trajectory collection, all task instance information
 196 is saved to a structured JSON file. This includes scene name, room name, floor name, per-object
 197 category, model instance, verified initial position and quaternion, agent initial pose and quat p_0 ,
 198 per-view object existence flags derived from segmentation masks, question, and ground-truth answer
 199 and action trajectory. This metadata file forms a self-contained and fully reproducible record of the
 200 task instance.

201 **Human Verification.** All generated task instances are examined by human annotators using the
 202 rendered per-step observations alongside the metadata. Annotators verify along three axes: correct-
 203 ness, ensuring the initial position and whole trajectory is physically valid; answerability, ensuring the
 204 task is solvable through interaction and the spatial configuration is unambiguous; and non-triviality,
 205 ensuring the task cannot be resolved through visual bias or prior knowledge without interaction and
 206 presents genuine spatial uncertainty. Task instances that fail any criterion are discarded.

207 **3.3 Task Categories and Statistics**

208 ESI-Bench comprises N total scenarios across 10 categories and 30 subcategories, summarized in
 209 Table 2. Figure ?? reports the task contribution. Each category targets a distinct spatial faculty
 210 structurally inaccessible to passive sensing. Physical Capacity and Physical Dynamics tasks require
 211 manipulation to reveal hidden physical properties such as containment capacity and gravitational
 212 stability. Reflective Reasoning and Perceptual Grounding tasks require active repositioning to
 213 disambiguate viewpoint-dependent phenomena where the correct answer changes critically with the
 214 observer’s position. Metric Comparison and Enumerative Perception tasks require locomotion and
 215 manipulation to overcome forced-perspective distortions where objects appear merged or equal in
 216 size from any fixed view. Spatial Relations tasks require navigation to vantage points that break
 217 projective symmetry, since configurations such as collinearity and contact are indistinguishable
 218 from most angles. Cognitive Mapping tasks require multi-step locomotion to construct topological
 219 representations of space that no single observation can provide. Temporal Scene Understanding
 220 tasks require manipulation and interaction to trigger or observe state changes, while Procedural
 221 Sequencing tasks require executing and reasoning over ordered action sequences to determine causal
 222 dependencies. Across all categories, the correct answer emerges not from any single image but from
 223 the agent’s capacity to act selectively and reason over the result.

224 potentially appendix if boring.... or should be interesting. select more interesting ones....

Table 2: **ESI-Bench task taxonomy**. Ten categories spanning 30 subcategories, each targeting a distinct spatial faculty requiring embodied action to resolve.

Category	Subcategory	Description	Example
Physics Capacity	Rigid Containment	Plan the placement of multiple objects across multiple containers	How can all the toys be fit into the boxes?
	Liquid Volume	Compare liquid-holding capacity of containers	Which container has larger volume?
	Deformable Fitting	Whether a deformable container can conform to an object	How to fold the sweater to make it fit?
Physics Dynamics	Inclined Plane	Predict object motion and stability on slopes	Can the apple sit stable on the slope?
	Stacking Stability	Whether objects stack or balance given shape, mass and geometry	How to stack and stabilize these three objects?
Specular Reflection	Reflection Authorization	Distinguish real objects from mirror reflections	Is the object real or reflected?
	Spatial Relationship	Infer relations among objects across mirror and real-world views	What relations does the mirror reveal?
	Scene Correspondence	Identify which objects appear in mirror given the real-world scene	Which of the 3 snacks appears in the mirror?
Perceptual Grounding	Partial Occlusion	Reason about objects hidden behind other scene elements	Golf Stick or Umbrella behind the Wall?
	Viewpoint Hallucination	Detect objects whose visibility changes critically with viewing angle	Cabinet or piano from this view?
	Material Transparency	Reason about objects seen through transparent surfaces	Is the object inside the glass or not?
Metric Comparison	Dimensional Size	Compare relative sizes of objects	Which vase is larger?
	Spatial Distance	Compare relative distances with respect to a reference object	Which flower is closer to the table?
Enumerative Perception	Occluded Counting	Count objects partially obscured by other scene elements	How many balls are under the blanket?
	Spatial Segmentation	Count objects separated across distinct spatial regions	One or two cylinders separated by the post?
	Category Ambiguity	Count visually similar objects requiring fine-grained distinction	How many apples among the balls?
	Merged Observation	Count groups that appear visually merged from a single viewpoint	How many separate stacks of books?
	Illumination Variability	Count objects under challenging or non-uniform lighting	Objects in the dim scene?
	Structural Enclosure	Count objects hidden within enclosed or covered spaces	Objects in the microwave?
Spatial Relations	Linear Alignment	Whether objects are arranged along a common axis	Do the glasses form a line?
	Geometric Configuration	Identify the shape formed by a set of objects	Do the cups form an equilateral triangle?
	Physical Contact	Detect whether two or more objects are in direct contact	Are the teddy bears touching each other?
Cognitive Mapping	Topology & Connectivity	Whether two locations or regions are mutually reachable	Is region A connected to C?
	Traversable Passage	Identify navigable corridors or passageways between regions	Passage between the rooms?
	Regional Boundary	Identify and delineate distinct functional spatial regions	Boundaries of this region?
Temporal Scene	Long-Horizon Navigation	Plan multi-step navigation trajectories toward a distant goal	What room is behind the wall?
	Unobserved State Change	Infer scene changes that occurred during an unobserved interval	What changed when I looked away?
Procedural Sequencing	Multi-Agent Interaction	Reason about scene dynamics from other agents	Other robot's world model?
	Action Ordering	Determine the correct procedural ordering of a sequence of actions	What is the action sequence for assembly?

225 4 Experiments

226 4.1 Models and Evaluation Setup

227 We evaluate models across four paradigms, organized by the degree of action and perceptual access
 228 granted to the agent: (1) **Passive Single-View** provides a single fixed observation from the initial
 229 pose, establishing a baseline under conditions identical to existing spatial benchmarks; (2) **Passive**
 230 **Multi-View** provides a set of pre-defined randomly sampled views designed to broadly cover the full
 231 environment, simulating exhaustive passive scene coverage without any agent action or viewpoint
 232 selection; (3) **Active Exploration** places the agent the initial pose with full access to the action space,
 233 requiring it to gather evidence through deliberate movement and interaction before committing to a
 234 final answer; Ground-Truth Passive provides the sequence of views rendered along the ground-truth
 235 action trajectory, serving as an oracle ablation that *separates perception errors from action errors*.
 236 Comparing (1) vs. (2) reveals whether exhaustive passive coverage helps; comparing (2) vs. (4)
 237 isolates the benefit of action-guided over passive coverage; comparing (4) vs. (3) isolates whether
 238 failures stem from the agent's inability to select informative actions or from perceptual limitations
 239 given the perfect views themselves.

240 We evaluate two families of models. 2D vision-language models include GPT-5 and Gemini 3.1, each
 241 taking egocentric visual observations as input, evaluated across all four paradigms. 3D-augmented
 242 models include VGGT+Gemini, where explicit 3D scene representations are reconstructed from
 243 multi-view observations via VGGT, from which scene graphs are constructed and provided to the
 244 language model; and Ground-Truth 3D+Gemini, where perfect point clouds derived directly from
 245 simulator state are used to construct scene graphs provided to the language model, serving as an
 246 oracle ablation for 3D grounding. Human performance is additionally collected on a subset of tasks to
 247 establish a human upper bound. All models are evaluated zero-shot with no task-specific fine-tuning,
 248 and results are reported as accuracy on matched subsets to ensure fair cross-paradigm comparison.

249 4.2 What Really Holds Spatial Intelligence Back: Seeing or Acting?

250 Table 3 reports accuracy across all models, paradigms, and task categories. Without any explicit
 251 instruction, active agents *spontaneously discover emergent spatial strategies*: to determine whether

Table 3: **ESI-Bench results (accuracy %)**. *2D + VLM*: GPT-5, Gemini 3.1. *3D + LLM*: VGGT + Gemini, Ground-truth 3D+ Gemini.

Category / Subcategory	2D + VLM								3D + LLM								Human Performance			
	GPT-5				Gemini 3.1				VGGT + Gemini				GT 3D + Gemini							
	Passive Single	Passive Multiple	Active Multiple	GT	Passive Single	Passive Multiple	Active Multiple	GT	Passive Single	Passive Multiple	Active Multiple	GT	Passive Single	Passive Multiple	Active Multiple	GT	Passive Single	Passive Multiple	Active Multiple	GT
Perceptual Grounding																				
Partial Occlusion	30.5	32.9	62.4	91.5	14.6	20.7	78.5	95.1	45.4	40.2	65.9	80.4	52.6	53.6	79.4	93.8	38.1	37.1	79.4	86.6
View Hallucination	11.7	20.2	60.1	87.8	39.9	32.9	68.1	91.1	56.3	59.0	76.1	87.6	61.6	60.6	74.1	89.0	52.2	45.8	80.8	83.8
Material Transparency	30.3	36.7	66.1	96.3	44.0	45.0	52.3	88.0	37.4	29.4	31.8	90.9	27.8	31.2	60.4	100	41.3	44.0	93.6	97.2
Physical Structure																				
Rigid Containment	45.0	42.5	42.5	95.0	47.5	40.0	67.5	97.5	27.5	37.5	57.5	72.5	45.0	42.5	65.0	95.0	47.5	42.5	92.5	100
Liquid Volume	66.2	66.2	81.6	86.0	69.9	67.6	80.9	86.8	65.4	56.7	71.3	77.2	74.5	74.5	83.1	86.8	79.4	66.9	89.7	86.8
Deformable	42.9	41.8	55.1	75.5	34.7	41.8	43.9	78.6	43.9	42.9	49.0	81.6	98.0	98.0	99.0	99.0	57.1	60.2	76.5	82.7
Physical Dynamics																				
Inclined Plane	57.4	60.7	77.0	86.9	65.6	62.3	83.6	88.5	67.2	62.3	80.3	86.9	63.9	63.9	83.6	91.8	60.7	62.3	83.6	83.6
Stacking & Stability	34.8	37.1	62.9	86.5	38.2	36.0	78.7	84.3	34.8	39.3	55.1	62.9	27.2	33.7	68.5	86.5	36.0	39.3	84.3	86.5
Specular Reflection																				
Reflection Authoring	68.7	70.7	70.3	73.7	60.6	60.6	64.9	67.0	52.5	54.6	53.5	55.6	55.6	58.6	55.6	60.6	94.9	88.9	96.0	96.0
Spatial Relations	50.4	38.9	54.8	58.4	43.4	42.5	44.2	46.9	41.6	39.8	41.6	43.4	38.9	35.4	54.2	41.6	78.8	78.8	84.1	87.6
Correspondence	39.8	51.1	52.3	56.8	37.5	40.9	42.0	42.0	37.5	43.2	48.9	48.9	31.8	39.8	42.3	48.9	85.2	80.7	89.8	92.0
Enumerative Perception																				
Counting w Occlusion	3.3	3.3	13.3	56.7	3.3	3.3	10.0	63.3	0.0	13.3	53.3	33.3	80.0	83.3	100.0	6.7	10.0	53.3	76.7	
Spatial Segmentation	3.3	6.7	26.7	63.3	3.3	3.3	16.7	70.0	3.3	3.3	23.3	56.7	43.3	80.0	66.7	86.7	13.3	23.3	63.3	70.0
Merged Observation	38.3	35.0	51.7	50.0	43.3	51.7	61.7	75.0	63.3	51.7	55.0	65.0	100.0	100.0	93.3	100.0	60.0	60.0	65.0	73.3
Category Ambiguity	8.3	10.0	8.3	41.7	13.3	13.3	15.0	48.3	21.7	26.7	26.7	46.7	18.3	76.7	40.0	93.3	23.3	25.0	51.7	61.7
Structural Enclosure	5.0	10.0	22.5	67.5	2.5	0.0	10.0	52.5	0.0	0.0	12.5	52.5	40.0	60.0	50.0	100.0	10.0	22.5	42.5	77.5
Illumination Variability	6.0	22.0	22.0	46.0	12.0	16.0	22.0	58.0	10.0	28.0	30.0	40.0	20.0	84.0	42.0	96.0	20.0	34.0	62.0	66.0
Spatial Relations																				
Linear Alignment	27.7	31.9	42.6	60.6	47.9	44.6	67.0	77.7	45.7	38.3	53.2	59.6	73.4	73.4	84.0	89.4	51.1	39.4	73.4	79.8
Geometric Configuration	25.3	20.4	26.0	26.0	27.5	22.3	27.5	44.6	9.9	11.3	18.6	32.6	70.8	73.6	88.0	100	32.4	33.5	74.3	86.3
Physical Contact	40.0	41.7	64.2	90.0	60.8	55.8	70.0	74.2	59.2	41.7	59.2	78.2	65.8	67.5	70.8	72.5	35.8	40.8	88.3	90.8
Metric Comparison																				
Dimensional Size	42.5	44.9	67.7	80.3	44.3	41.3	68.3	82.6	50.4	40.1	59.3	80.8	58.6	60.5	69.9	85.9	48.5	56.9	82.6	91.9
Spatial Distance	53.9	49.1	58.6	73.7	52.6	50.5	59.9	80.3	50.7	47.3	55.0	71.4	59.6	61.2	64.5	96.7	57.9	59.9	69.1	78.9
Cognitive Mapping																				
Connectivity	68.3	70.0	68.3	78.3	51.7	48.3	60.0	85.0	66.7	66.7	73.3	83.3	65.0	65.0	71.7	86.7	68.3	68.3	81.7	91.7
Traversable Passage	68.3	66.7	71.7	73.3	66.7	61.7	73.3	78.3	68.3	63.3	66.7	78.3	71.7	73.3	71.7	80.0	70.0	71.7	78.3	85.0
Regional Boundary	65.0	63.8	65.0	67.5	63.8	62.5	65.0	65.0	65.0	65.0	62.5	67.5	60.0	68.8	62.5	68.8	61.3	62.5	67.5	70.0
Long-Term Navigation	40.0	38.3	41.7	50.0	33.3	35.0	33.3	41.7	36.7	40.0	43.3	51.6	36.7	41.7	43.3	40.0	35.0	36.7	51.7	61.7
Temporal Understanding																				
Unobserved Change	40.5	41.2	51.4	77.0	37.2	37.8	47.3	74.5	37.8	41.9	45.8	76.6	39.9	42.6	46.5	74.5	39.4	39.4	70.8	81.0
Agent Observation	40.6	30.1	51.0	72.7	37.6	36.1	58.6	65.4	32.7	27.9	34.2	56.5	36.3	33.3	76.7	87.8	38.3	42.1	83.5	90.2
Action Sequencing																				
Action Order Inference	36.4	37.7	41.6	67.5	44.2	46.8	54.5	72.7	35.1	31.1	50.6	58.4	44.2	46.8	51.9	75.3	40.3	41.6	74.0	81.8

252 a chestnut is inside a glass, agents independently develop four distinct approaches, moving behind
 253 the object, repositioning top-down, picking it up, and pouring it out (Figure 3a), none of which were
 254 prescribed. These emergent abilities drive consistent and substantial gains across tasks requiring
 255 deliberate repositioning: on Dimensional Size GPT-5 improves from 42.5% to 67.7%, on Physical
 256 Contact from 40.0% to 64.2%, and on Viewpoint Hallucination Gemini 3.1 jumps from 39.9% to
 257 68.1%. In contrast, passive multi-view provides negligible or negative gains despite consuming far
 258 more images (*e.g.*, PT-5 drops from 53.9% to 49.1% on Spatial Distance compared with single-image
 259 passive baseline), confirming that observation quantity without selective action adds noise rather than
 260 signal.

261 Ground-truth passive trajectories reveal that the bottleneck is almost entirely in action selection, not
 262 perception: GPT-5 reaches 90.0% on Physical Contact and 95.0% on Rigid Containment under oracle
 263 views, and Gemini 3.1 jumps from 14.6% to 95.1% on Partial Occlusion. Geometric Configuration
 264 and Specular Reflection stand as exceptions where oracle views provide minimal benefit: GPT-5
 265 reaches only 26.0% on Geometric Configuration under ground-truth passive, and models consistently
 266 fail to determine whether three objects form an equilateral triangle even from the perfect viewpoint;
 267 on Specular Reflection, models hallucinate objects in mirrors that do not exist, or fail to identify the

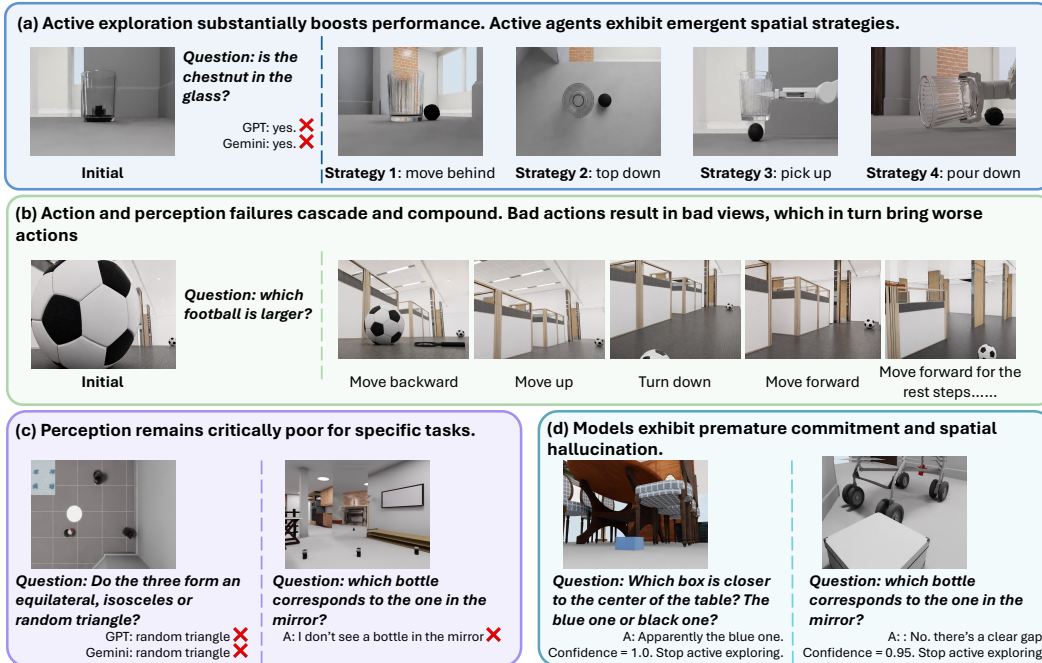


Figure 3: Qualitative results of ESI-Bench.

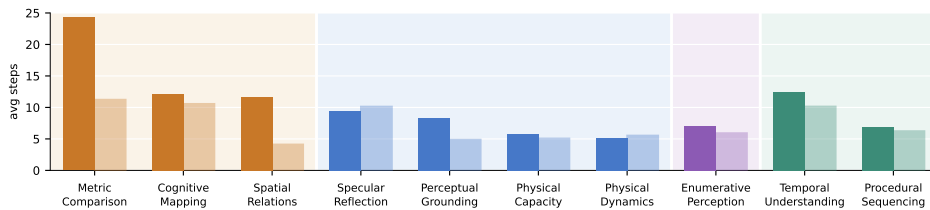


Figure 4: Average number of active exploration steps to reach a correct answer for GPT-5 (solid) and Gemini 3.1 (outlined), grouped by Spekke's four core knowledge systems.

268 correct real-world correspondence altogether (Figure 3c). These indicate hard perceptual limits that
 269 no action strategy can overcome.

270 The active-to-oracle gap confirms that actions and perception failures cascade and compound: on
 271 Occluded Counting the gap reaches 33.6%, and on Structural Enclosure 49.7%, confirming that bad
 272 actions result in bad views, which in turn bring worse actions that cannot be recovered within the
 273 step budget (Figure 3b).

💡 Action Blindness Dominates Perceptual Blindness, while Their Coupling Drives Failure Cascades

- Without explicit instruction, agents spontaneously acquire spatial strategies, yielding substantial gains over passive multi-view, which often adds noise rather than signal.
- For most tasks, perception is not the bottleneck: with the right viewpoint, the agent succeeds; yet some tasks hit a hard perceptual ceiling that no action can overcome.
- Suboptimal actions produce bad views, which in turn produce worse actions, cascading into reasoning failures unrecoverable within the step budget.

274

275 4.3 When Does 3D Help, and When Does It Hurt?

276 Ground-Truth 3D+Gemini reaches 60.4% on Material Transparency versus 44.0% for Gemini 3.1, a
 277 16.4 point improvement from 3D grounding alone, most pronounced on tasks where 2D projections

278 fundamentally lose depth information. However, VGGT+Gemini drops to 9.9% on Geometric
279 Configuration versus Gemini 3.1’s 27.5%, a 17.6 point degradation from reconstruction artifacts.
280 Imperfect geometry actively misleads rather than merely failing to help: VGGT reconstructions
281 distort fine-grained spatial relations, causing the model to reason over a corrupted scene graph rather
282 than the true geometry, making 3D augmentation a strategy that amplifies failures.

💡 3D Helps When Geometry Is Perfect, But Imperfect Reconstruction Actively Misleads

- Explicit 3D grounding improves performance on tasks where depth and occlusion render 2D projections fundamentally ambiguous.
- Imperfect 3D reconstruction actively degrades performance, proving more harmful than 2D baselines.

283

284 4.4 How Far Are Models From Human-Level Spatial Reasoning?

285 Human annotators consistently gather more observations before committing, actively seek viewpoints
286 that falsify their current hypothesis, and reduce confidence under ambiguity. Models commit after
287 fewer steps with uniformly high confidence regardless of evidence quality, manifesting as spatial
288 hallucination: asserting object properties that directly contradict the scene state (Figure 3d). The
289 gap is not perceptual but metacognitive: models lack the awareness to recognize when current
290 observations are insufficient, a failure that neither better perception nor more embodied interaction
291 alone can fully close.

💡 Models Can See But Do Not Know When They Have Seen Enough

- Models commit prematurely with high confidence; humans treat uncertainty as a signal to keep looking, not a reason to answer.
- Humans seek viewpoints that falsify their hypothesis; models seek confirmation and move in the same direction repeatedly.
- Humans revise beliefs when contradicted; models anchor to their first impression, revealing a fundamental absence of belief updating.

292

293 5 Conclusion

294 We introduced ESI-Bench, a benchmark spanning 10 task categories and 30 subcategories that
295 requires agents to close the perception-action loop rather than reason from pre-given observations.
296 Our experiments show that active exploration substantially outperforms passive counterparts, while
297 passive multi-view coverage adds noise rather than signal. The primary bottleneck across nearly
298 all tasks is not perception but action selection: given oracle viewpoints, models succeed; given the
299 freedom to act, they struggle to acquire them. Beyond accuracy, we identify a persistent gap in
300 epistemic calibration: unlike humans, models commit prematurely with high confidence regardless
301 of evidence quality, a failure that neither better perception nor more exploration alone can resolve.
302 ESI-Bench establishes a clear frontier for embodied spatial intelligence, where closing the perception-
303 action loop remains an open and fundamental challenge.

304 References

305 References

- 306 Wenxiao Cai, Yaroslav Ponomarenko, Jianhao Yuan, Xiaoqi Li, Wankou Yang, Hao Dong, and
307 Bo Zhao. Spatialbot: Precise spatial understanding with vision language models. *arXiv preprint*
308 *arXiv:2406.13642*, 2024.
- 309 Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia.
310 Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings*
311 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14455–
312 14465, June 2024.

- 313 An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang,
314 and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision language models, 2024. URL
315 <https://arxiv.org/abs/2406.01584>.
- 316 Zhili Cheng, Yuge Tu, Ran Li, Shiqi Dai, Jinyi Hu, Shengding Hu, Jiahao Li, Yang Shi, Tianyu Yu,
317 Weize Chen, Lei Shi, and Maosong Sun. Embodiedeval: Evaluate multimodal llms as embodied
318 agents, 2025. URL <https://arxiv.org/abs/2501.11858>.
- 319 Wei Chow, Jiageng Mao, Boyi Li, Daniel Seita, Vitor Guizilini, and Yue Wang. Physbench: Bench-
320 marking and enhancing vision-language models for physical world understanding, 2025. URL
321 <https://arxiv.org/abs/2501.16411>.
- 322 Mengfei Du, Binhao Wu, Zejun Li, Xuanjing Huang, and Zhongyu Wei. Embspatial-bench: Bench-
323 marking spatial understanding for embodied tasks with large vision-language models, 2024. URL
324 <https://arxiv.org/abs/2406.05756>.
- 325 Zhiwen Fan, Jian Zhang, Renjie Li, Junge Zhang, Runjin Chen, Hezhen Hu, Kevin Wang, Huaizhi
326 Qu, Shijie Zhou, Dilin Wang, Zhicheng Yan, Hongyu Xu, Justin Theiss, Tianlong Chen, Jiachen
327 Li, Zhengzhong Tu, Zhangyang Wang, and Rakesh Ranjan. Vlm-3r: Vision-language models
328 augmented with instruction-aligned 3d reconstruction, 2026. URL <https://arxiv.org/abs/2505.20279>.
- 330 Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A. Smith,
331 Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not
332 perceive, 2024. URL <https://arxiv.org/abs/2404.12390>.
- 333 James J. Gibson. *The Ecological Approach to Visual Perception*. Houghton Mifflin, 1979.
- 334 Kaixuan Jiang, Yang Liu, Weixing Chen, Jingzhou Luo, Ziliang Chen, Ling Pan, Guanbin Li, and
335 Liang Lin. Beyond the destination: A novel benchmark for exploration-aware embodied question
336 answering, 2025. URL <https://arxiv.org/abs/2503.11117>.
- 337 Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial reasoning, 2023. URL <https://arxiv.org/abs/2205.00363>.
- 339 Mengzhen Liu, Enshen Zhou, Cheng Chi, Yi Han, Shanyu Rong, Liming Chen, Pengwei Wang,
340 Zhongyuan Wang, and Shanghang Zhang. Sapave: Towards active perception and manipulation
341 in vision-language-action models for robotics, 2026a. URL <https://arxiv.org/abs/2603.12193>.
- 343 Yuecheng Liu, Dafeng Chi, Shiguang Wu, Zhanguang Zhang, Yaochen Hu, Lingfeng Zhang, Yingxue
344 Zhang, Shuang Wu, Tongtong Cao, Guowei Huang, Helong Huang, Guangjian Tian, Weichao
345 Qiu, Xingyue Quan, Jianye Hao, and Yuzheng Zhuang. Spatialcot: Advancing spatial reasoning
346 through coordinate alignment and chain-of-thought for embodied task planning, 2025. URL
347 <https://arxiv.org/abs/2501.10074>.
- 348 Zhenyang Liu, Yongchong Gu, Yikai Wang, Xiangyang Xue, and Yanwei Fu. Activevla: Injecting
349 active perception into vision-language-action models for precise 3d robotic manipulation, 2026b.
350 URL <https://arxiv.org/abs/2601.08325>.
- 351 Wenxin Ma, Chenlong Wang, Ruisheng Yuan, Hao Chen, Nanru Dai, S. Kevin Zhou, Yijun Yang,
352 Alan Yuille, and Jieneng Chen. Causalspatial: A benchmark for object-centric causal spatial
353 reasoning, 2026. URL <https://arxiv.org/abs/2601.13304>.
- 354 Wufei Ma, Haoyu Chen, Guofeng Zhang, Yu-Cheng Chou, Celso M de Melo, and Alan Yuille.
355 3dsrbench: A comprehensive 3d spatial reasoning benchmark, 2025. URL <https://arxiv.org/abs/2412.07825>.
- 357 Arjun Majumdar, Anurag Ajay, Xiaohan Zhang, Pranav Putta, Sriram Yenamandra, Mikael Henaff,
358 Sneha Silwal, Paul Mcvay, Oleksandr Maksymets, Sergio Arnaud, Karmesh Yadav, Qiyang Li,
359 Ben Newman, Mohit Sharma, Vincent Berges, Shiqi Zhang, Pulkit Agrawal, Yonatan Bisk, Dhruv
360 Batra, Mrinal Kalakrishnan, Franziska Meier, Chris Paxton, Sasha Sax, and Aravind Rajeswaran.
361 Openeqa: Embodied question answering in the era of foundation models. In *Conference on*
362 *Computer Vision and Pattern Recognition (CVPR)*, 2024.

- 363 J. Kevin O'Regan and Alva Noë. A sensorimotor account of vision and visual consciousness.
364 *Behavioral and Brain Sciences*, 24(5):939–973, 2001.
- 365 Elizabeth S. Spelke and Katherine D. Kinzler. Core knowledge. *Developmental Science*, 10(1):89–96,
366 2007.
- 367 Francisco J. Varela, Evan Thompson, and Eleanor Rosch. *The Embodied Mind: Cognitive Science*
368 *and Human Experience*. MIT Press, 1991.
- 369 Qineng Wang, Baiqiao Yin, Pingyue Zhang, Jianshu Zhang, Kangrui Wang, Zihan Wang, Jieyu
370 Zhang, Keshigeyan Chandrasegaran, Han Liu, Ranjay Krishna, Saining Xie, Jiajun Wu, Li Fei-
371 Fei, and Manling Li. Mindcube: Spatial mental modeling from limited views, 2026. URL <https://arxiv.org/abs/2506.21458>.
- 373 Diankun Wu, Fangfu Liu, Yi-Hsin Hung, and Yueqi Duan. Spatial-mlm: Boosting mllm capabilities
374 in visual-based spatial intelligence, 2025a. URL <https://arxiv.org/abs/2505.23747>.
- 375 Haoning Wu, Xiao Huang, Yaohui Chen, Ya Zhang, Yanfeng Wang, and Weidi Xie. Spatialscore:
376 Towards comprehensive evaluation for spatial intelligence, 2026a. URL <https://arxiv.org/abs/2505.17012>.
- 378 Junfei Wu, Jian Guan, Kaituo Feng, Qiang Liu, Shu Wu, Liang Wang, Wei Wu, and Tieniu Tan.
379 Reinforcing spatial reasoning in vision-language models with interwoven thinking and visual
380 drawing, 2025b. URL <https://arxiv.org/abs/2506.09965>.
- 381 Yuhao Wu, Maojia Song, Yihuai Lan, Lei Wang, Zhiqiang Hu, Yao Xiao, Heng Zhou, Weihua Zheng,
382 Dylan Raharja, Soujanya Poria, and Roy Ka-Wei Lee. From perception to action: An interactive
383 benchmark for vision reasoning, 2026b. URL <https://arxiv.org/abs/2602.21015>.
- 384 Haoyu Xiong, Xiaomeng Xu, Jimmy Wu, Yifan Hou, Jeannette Bohg, and Shuran Song. Vision in
385 action: Learning active perception from human demonstrations, 2025. URL <https://arxiv.org/abs/2506.15666>.
- 387 Jihan Yang, Shusheng Yang, Anjali W. Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking
388 in space: How multimodal large language models see, remember, and recall spaces, 2024. URL
389 <https://arxiv.org/abs/2412.14171>.
- 390 Jihan Yang, Shusheng Yang, Anjali W. Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking in
391 space: How multimodal large language models see, remember, and recall spaces, 2025a. URL
392 <https://arxiv.org/abs/2412.14171>.
- 393 Rui Yang, Hanyang Chen, Junyu Zhang, Mark Zhao, Cheng Qian, Kangrui Wang, Qineng Wang,
394 Teja Venkat Koripella, Marziyeh Movahedi, Manling Li, Heng Ji, Huan Zhang, and Tong Zhang.
395 Embodiedbench: Comprehensive benchmarking multi-modal large language models for vision-
396 driven embodied agents, 2025b. URL <https://arxiv.org/abs/2502.09560>.
- 397 Shusheng Yang, Jihan Yang, Pinzhi Huang, Ellis Brown, Zihao Yang, Yue Yu, Shengbang Tong,
398 Zihan Zheng, Yifan Xu, Muhan Wang, Daohan Lu, Rob Fergus, Yann LeCun, Li Fei-Fei, and
399 Saining Xie. Cambrian-s: Towards spatial supersensing in video, 2025c. URL <https://arxiv.org/abs/2511.04670>.
- 401 Sihan Yang, Runsen Xu, Yiman Xie, Sizhe Yang, Mo Li, Jingli Lin, Chenming Zhu, Xiaochen
402 Chen, Haodong Duan, Xiangyu Yue, Dahua Lin, Tai Wang, and Jiangmiao Pang. Mmsi-bench: A
403 benchmark for multi-image spatial intelligence, 2025d. URL <https://arxiv.org/abs/2505.23764>.
- 405 Heyang Yu, Yinan Han, Xiangyu Zhang, Baiqiao Yin, Bowen Chang, Xiangyu Han, Xinhao Liu, Jing
406 Zhang, Marco Pavone, Chen Feng, Saining Xie, and Yiming Li. Thinking in 360°: Humanoid
407 visual search in the wild, 2025. URL <https://arxiv.org/abs/2511.20351>.
- 408 Yanpeng Zhao, Wentao Ding, Hongtao Li, Baoxiong Jia, and Zilong Zheng. Espire: A diagnostic
409 benchmark for embodied spatial reasoning of vision-language models, 2026. URL <https://arxiv.org/abs/2603.13033>.

411 **NeurIPS Paper Checklist**

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

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

- 419 • You should answer [Yes], [No], or [N/A].
- 420 • [N/A] means either that the question is Not Applicable for that particular paper or the
421 relevant information is Not Available.
- 422 • Please provide a short (1–2 sentence) justification right after your answer (even for [N/A]).

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

427 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
428 While [Yes] is generally preferable to [No], it is perfectly acceptable to answer [No] provided a
429 proper justification is given (e.g., error bars are not reported because it would be too computationally
430 expensive” or “we were unable to find the license for the dataset we used”). In general, answering
431 [No] or [N/A] is not grounds for rejection. While the questions are phrased in a binary way, we
432 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
433 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
434 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
435 please point to the section(s) where related material for the question can be found.

436 **IMPORTANT, please:**

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

440 **1. Claims**

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

443 Answer: **[TODO]**

444 Justification: **[TODO]**

445 Guidelines:

- 446 • The answer [N/A] means that the abstract and introduction do not include the claims
447 made in the paper.
- 448 • The abstract and/or introduction should clearly state the claims made, including the
449 contributions made in the paper and important assumptions and limitations. A [No] or
450 [N/A] answer to this question will not be perceived well by the reviewers.
- 451 • The claims made should match theoretical and experimental results, and reflect how
452 much the results can be expected to generalize to other settings.
- 453 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
454 are not attained by the paper.

455 **2. Limitations**

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

457 Answer: **[TODO]**

458 Justification: **[TODO]**

459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509

Guidelines:

- The answer [N/A] 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: [TODO]

Justification: [TODO]

Guidelines:

- The answer [N/A] 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: [TODO]

Justification: [TODO]

Guidelines:

- The answer [N/A] means that the paper does not include experiments.

- 510 • If the paper includes experiments, a [No] answer to this question will not be perceived
511 well by the reviewers: Making the paper reproducible is important, regardless of
512 whether the code and data are provided or not.
- 513 • If the contribution is a dataset and/or model, the authors should describe the steps taken
514 to make their results reproducible or verifiable.
- 515 • Depending on the contribution, reproducibility can be accomplished in various ways.
516 For example, if the contribution is a novel architecture, describing the architecture fully
517 might suffice, or if the contribution is a specific model and empirical evaluation, it may
518 be necessary to either make it possible for others to replicate the model with the same
519 dataset, or provide access to the model. In general, releasing code and data is often
520 one good way to accomplish this, but reproducibility can also be provided via detailed
521 instructions for how to replicate the results, access to a hosted model (e.g., in the case
522 of a large language model), releasing of a model checkpoint, or other means that are
523 appropriate to the research performed.
- 524 • While NeurIPS does not require releasing code, the conference does require all submis-
525 sions to provide some reasonable avenue for reproducibility, which may depend on the
526 nature of the contribution. For example
 - 527 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
528 to reproduce that algorithm.
 - 529 (b) If the contribution is primarily a new model architecture, the paper should describe
530 the architecture clearly and fully.
 - 531 (c) If the contribution is a new model (e.g., a large language model), then there should
532 either be a way to access this model for reproducing the results or a way to reproduce
533 the model (e.g., with an open-source dataset or instructions for how to construct
534 the dataset).
 - 535 (d) We recognize that reproducibility may be tricky in some cases, in which case
536 authors are welcome to describe the particular way they provide for reproducibility.
537 In the case of closed-source models, it may be that access to the model is limited in
538 some way (e.g., to registered users), but it should be possible for other researchers
539 to have some path to reproducing or verifying the results.

540 5. Open access to data and code

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

544 Answer: [TODO]

545 Justification: [TODO]

546 Guidelines:

- 547 • The answer [N/A] means that paper does not include experiments requiring code.
- 548 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
549 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 550 • While we encourage the release of code and data, we understand that this might not
551 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
552 including code, unless this is central to the contribution (e.g., for a new open-source
553 benchmark).
- 554 • The instructions should contain the exact command and environment needed to run to
555 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
556 neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 557 • The authors should provide instructions on data access and preparation, including how
558 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 559 • The authors should provide scripts to reproduce all experimental results for the new
560 proposed method and baselines. If only a subset of experiments are reproducible, they
561 should state which ones are omitted from the script and why.
- 562 • At submission time, to preserve anonymity, the authors should release anonymized
563 versions (if applicable).

- 564 • Providing as much information as possible in supplemental material (appended to the
565 paper) is recommended, but including URLs to data and code is permitted.

566 6. Experimental setting/details

567 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
568 rameters, how they were chosen, type of optimizer) necessary to understand the results?

569 Answer: **[TODO]**

570 Justification: **[TODO]**

571 Guidelines:

- 572 • The answer [N/A] means that the paper does not include experiments.
- 573 • The experimental setting should be presented in the core of the paper to a level of detail
574 that is necessary to appreciate the results and make sense of them.
- 575 • The full details can be provided either with the code, in appendix, or as supplemental
576 material.

577 7. Experiment statistical significance

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

580 Answer: **[TODO]**

581 Justification: **[TODO]**

582 Guidelines:

- 583 • The answer [N/A] means that the paper does not include experiments.
- 584 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
585 intervals, or statistical significance tests, at least for the experiments that support the
586 main claims of the paper.
- 587 • The factors of variability that the error bars are capturing should be clearly stated (for
588 example, train/test split, initialization, random drawing of some parameter, or overall
589 run with given experimental conditions).
- 590 • The method for calculating the error bars should be explained (closed form formula,
591 call to a library function, bootstrap, etc.)
- 592 • The assumptions made should be given (e.g., Normally distributed errors).
- 593 • It should be clear whether the error bar is the standard deviation or the standard error
594 of the mean.
- 595 • It is OK to report 1-sigma error bars, but one should state it. The authors should
596 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
597 of Normality of errors is not verified.
- 598 • For asymmetric distributions, the authors should be careful not to show in tables or
599 figures symmetric error bars that would yield results that are out of range (e.g., negative
600 error rates).
- 601 • If error bars are reported in tables or plots, the authors should explain in the text how
602 they were calculated and reference the corresponding figures or tables in the text.

603 8. Experiments compute resources

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

607 Answer: **[TODO]**

608 Justification: **[TODO]**

609 Guidelines:

- 610 • The answer [N/A] means that the paper does not include experiments.
- 611 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
612 or cloud provider, including relevant memory and storage.
- 613 • The paper should provide the amount of compute required for each of the individual
614 experimental runs as well as estimate the total compute.

615 • The paper should disclose whether the full research project required more compute
616 than the experiments reported in the paper (e.g., preliminary or failed experiments that
617 didn't make it into the paper).

618 9. Code of ethics

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

621 Answer: **[TODO]**

622 Justification: **[TODO]**

623 Guidelines:

- 624 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
625 Ethics.
- 626 • If the authors answer [No], they should explain the special circumstances that require a
627 deviation from the Code of Ethics.
- 628 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
629 eration due to laws or regulations in their jurisdiction).

630 10. Broader impacts

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

633 Answer: **[TODO]**

634 Justification: **[TODO]**

635 Guidelines:

- 636 • The answer [N/A] means that there is no societal impact of the work performed.
- 637 • If the authors answer [N/A] or [No], they should explain why their work has no societal
638 impact or why the paper does not address societal impact.
- 639 • Examples of negative societal impacts include potential malicious or unintended uses
640 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
641 (e.g., deployment of technologies that could make decisions that unfairly impact specific
642 groups), privacy considerations, and security considerations.
- 643 • The conference expects that many papers will be foundational research and not tied
644 to particular applications, let alone deployments. However, if there is a direct path to
645 any negative applications, the authors should point it out. For example, it is legitimate
646 to point out that an improvement in the quality of generative models could be used to
647 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
648 that a generic algorithm for optimizing neural networks could enable people to train
649 models that generate Deepfakes faster.
- 650 • The authors should consider possible harms that could arise when the technology is
651 being used as intended and functioning correctly, harms that could arise when the
652 technology is being used as intended but gives incorrect results, and harms following
653 from (intentional or unintentional) misuse of the technology.
- 654 • If there are negative societal impacts, the authors could also discuss possible mitigation
655 strategies (e.g., gated release of models, providing defenses in addition to attacks,
656 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
657 feedback over time, improving the efficiency and accessibility of ML).

658 11. Safeguards

659 Question: Does the paper describe safeguards that have been put in place for responsible
660 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
661 image generators, or scraped datasets)?

662 Answer: **[TODO]**

663 Justification: **[TODO]**

664 Guidelines:

- 665 • The answer [N/A] means that the paper poses no such risks.

- 666
- 667
- 668
- 669
- 670
- 671
- 672
- 673
- 674
- 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.

675 12. Licenses for existing assets

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

679 Answer: **[TODO]**

680 Justification: **[TODO]**

681 Guidelines:

- 682
- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- The answer [N/A] 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.

697 13. New assets

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

700 Answer: **[TODO]**

701 Justification: **[TODO]**

702 Guidelines:

- 703
- 704
- 705
- 706
- 707
- 708
- 709
- 710
- The answer [N/A] means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

711 14. Crowdsourcing and research with human subjects

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

715 Answer: **[TODO]**

716 Justification: **[TODO]**

717
718
719
720
721
722
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

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

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

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

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

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

16. Declaration of LLM usage

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

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.