



OMNI^{EAR}: BENCHMARKING AGENT REASONING IN EMBODIED TASKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models excel at abstract reasoning but their capacity for embodied agent reasoning remains largely unexplored. We present **OmniEAR**, a comprehensive framework for evaluating how language models reason about physical interactions, tool usage, and multi-agent coordination in embodied tasks. Unlike existing benchmarks that provide predefined tool sets or explicit collaboration directives, OmniEAR requires agents to dynamically acquire capabilities and autonomously determine coordination strategies based on task demands. Through text-based environment representation, we model continuous physical properties and complex spatial relationships across 1,500 scenarios spanning household and industrial domains. Our systematic evaluation reveals severe performance degradation when models must reason from constraints: while achieving 85-96% success with explicit instructions, performance drops to 56-85% for tool reasoning and 63-85% for implicit collaboration, with compound tasks showing over 50% failure rates. Surprisingly, complete environmental information degrades coordination performance, indicating models cannot filter task-relevant constraints. Fine-tuning improves single-agent tasks dramatically (0.6% to 76.3%) but yields minimal multi-agent gains (1.5% to 5.5%), exposing fundamental architectural limitations. These findings demonstrate that embodied reasoning poses fundamentally different challenges than current models can address, establishing OmniEAR as a rigorous benchmark for evaluating and advancing embodied AI systems.

1 INTRODUCTION

Large language models have achieved remarkable success in complex reasoning tasks (Brown et al., 2020; Wei et al., 2022), yet their ability to reason about embodied environments remains poorly understood. In embodied tasks, agents must understand how object properties affect what actions are possible, recognize when their capabilities are insufficient for a task, and determine when collaboration becomes necessary (Ahn et al., 2022; Wu et al., 2023). These reasoning abilities fundamentally differ from abstract problem-solving, as they require understanding the physical principles that govern real-world interactions.

Current evaluation approaches fail to capture this embodied reasoning complexity. Existing benchmarks model environments through discrete states like open/closed doors or picked/placed objects (Shridhar et al., 2020; Puig et al., 2018), overlooking continuous properties such as weight, temperature, or material composition that determine action feasibility. Tool usage evaluations typically provide fixed action sets (Chang et al., 2024; Huang et al., 2022), missing how agents should reason about capability gaps. Multi-agent benchmarks rely on explicit collaboration instructions or efficiency metrics (Kang et al., 2025; Zhang et al., 2024), rather than examining whether agents can recognize when tasks exceed individual abilities. This evaluation paradigm cannot assess understanding of embodied principles.

The core challenge is that real-world embodied reasoning emerges from understanding environmental realities and task requirements. When objects are too heavy for single agents, collaboration naturally becomes necessary. When tasks require manipulating materials beyond native capabilities, tools provide the solution. When spatial layouts limit individual reach, coordinated action enables

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

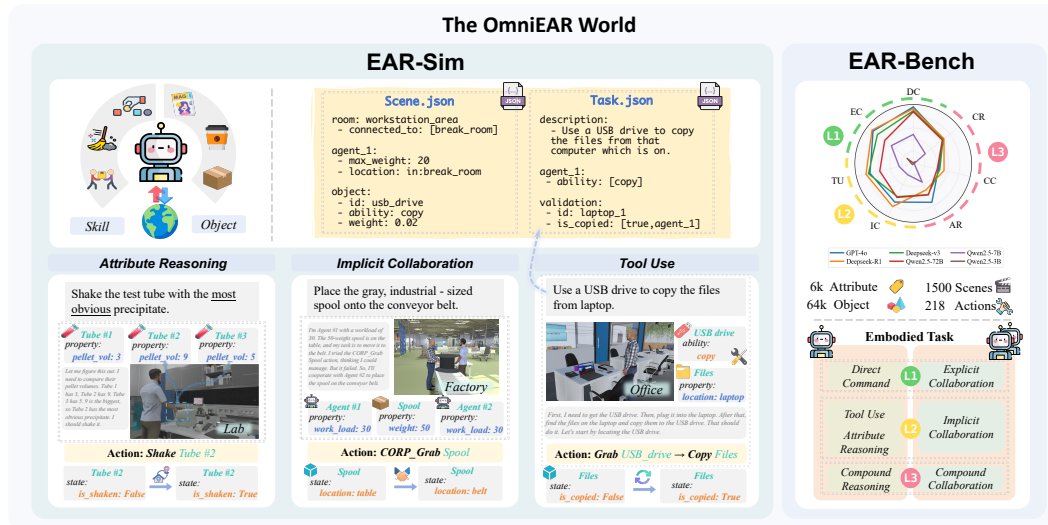


Figure 1: Overview of the OmniEAR framework comprising three integrated components: **EAR-Sim** (left) uses structured text representation to model environments with objects, agents, and spatial relationships, enabling dynamic tool-capability binding and physics-constrained collaboration; **EAR-Bench** (right) presents our comprehensive evaluation matrix spanning single-agent and multi-agent tasks across increasing cognitive complexity levels.

task completion (Zeng et al., 2022; Wang et al., 2023). Current benchmarks rely on static tool sets and explicit collaboration instructions, preventing assessment of how models reason about capability acquisition and coordination needs based on task requirements.

We introduce **OmniEAR**, a comprehensive framework for evaluating agent reasoning in embodied tasks. Our key insight is that embodied reasoning requires understanding how physical properties shape possible actions, how capability limitations necessitate tools, and how task demands drive collaboration.

By designing scenarios where agents must dynamically acquire capabilities and autonomously determine coordination strategies based on task requirements, we can assess whether models genuinely comprehend the principles governing embodied interactions.

OmniEAR employs text-based environment representation to efficiently model rich physical properties while enabling large-scale evaluation. The framework comprises three integrated components: **EAR-Sim** captures detailed object attributes and spatial relationships while supporting dynamic capability evolution through tool acquisition; an automated pipeline generates diverse scenarios where task solutions naturally depend on understanding embodied principles; and **EAR-Bench** provides systematic evaluation through 1,500 scenarios across household and industrial domains.

Our evaluation focuses on three core aspects of embodied reasoning. First, we assess how agents reason about object properties like weight, material, and temperature when determining feasible actions, requiring comparison and inference about continuous attributes. Second, we examine whether agents recognize when tasks demand capabilities beyond their current abilities and plan appropriate tool acquisition. Third, we evaluate autonomous coordination decisions, testing whether agents identify when task requirements exceed individual capacities without explicit collaboration instructions. These capabilities reflect fundamental aspects of embodied intelligence.

Systematic evaluation reveals fundamental gaps in current models’ embodied reasoning abilities. While achieving 85-96% success on explicit instructions, performance degrades sharply when reasoning must emerge from physical constraints. Tool reasoning drops to 56-85% when models must infer capability needs, and implicit collaboration falls to 63-85% compared to 88-92% with explicit coordination. Compound tasks show the steepest decline, with failure rates exceeding 50%. Paradoxically, complete environmental information harms coordination performance, suggesting models cannot filter task-relevant from irrelevant constraints. Even reasoning-specialized models,

108 which excel at logical planning, fail to ground physical constraints effectively, demonstrating that
 109 current architectures lack the mechanisms necessary for autonomous embodied decision-making.

110 Our analysis uncovers important patterns in model capabilities. Smaller models cannot maintain the
 111 planning state necessary for multi-step reasoning about tools and coordination. Reasoning models
 112 excel at logical planning but struggle to ground abstract concepts in concrete physical properties.
 113 While supervised fine-tuning improves single-agent performance, these gains fail to transfer to
 114 multi-agent scenarios, suggesting that coordination reasoning requires architectural capabilities
 115 beyond current training approaches.

116 In summary, our contributions are:

- 117 • We present OmniEAR, a framework that evaluates embodied reasoning through scenarios
 118 requiring agents to understand how physical properties determine actions, capabilities, and
 119 coordination needs, addressing fundamental gaps in current evaluation methods.
- 120 • We develop EAR-Bench, a benchmark of 1,500 scenarios with continuous physical
 121 properties and dynamic capabilities, supported by EAR-Sim and an automated generation
 122 pipeline.
- 123 • We provide empirical evidence that current language models lack core embodied reasoning
 124 capabilities, with performance degrading over 60% when moving from explicit instructions
 125 to embodied reasoning, revealing critical requirements for advancing embodied AI.

126 2 RELATED WORKS

127 Prior embodied benchmarks have made significant contributions to task evaluation but differ
 128 fundamentally in their approach to physical reasoning and collaboration. While ALFRED (Shridhar
 129 et al., 2020) and BEHAVIOR-1K (Li et al., 2024a) provide extensive task coverage, they model
 130 physical states through discrete representations (e.g., binary door states, picked/placed objects)
 131 rather than continuous attributes necessary for reasoning about weight, temperature, or material
 132 properties. Tool usage evaluation spans from low-level manipulation in RoCo (Mandi et al.,
 133 2024) to high-level planning in PARTNR (Chang et al., 2024), yet both maintain static action
 134 spaces determined at initialization, preventing assessment of dynamic capability acquisition.
 135 Recent multi-agent benchmarks including TDW-MAT (Zhang et al., 2024) and EmbodiedBench
 136 (Yang et al., 2025) advance collaboration evaluation through load constraints and task allocation
 137 optimization, but rely on explicit task division instructions or efficiency-driven participation rather
 138 than collaboration that emerges from physical constraints. In contrast, OmniEAR introduces
 139 continuous property reasoning with 6,381 distinct attributes, dynamic tool-capability binding that
 140 expands action spaces during execution, and implicit collaboration where agents must autonomously
 141 recognize when tasks exceed individual capacities based on physical constraints, fundamentally
 142 shifting evaluation from instruction compliance to constraint-based reasoning. A comprehensive
 143 comparison with related work is provided in Appendix 6.1.

144 3 FRAMEWORK

145 We present OmniEAR, a comprehensive framework for evaluating agent reasoning in embodied
 146 tasks. Our framework addresses the fundamental challenge of assessing whether language models
 147 understand embodied principles. We achieve this through three key design principles: (1) tasks
 148 must require reasoning about physical properties and constraints rather than following explicit
 149 instructions, (2) agent capabilities should dynamically evolve based on tool acquisition rather than
 150 remaining static, and (3) collaboration needs should emerge from task requirements rather than
 151 predetermined protocols.

152 3.1 TASK DESIGN AND FORMALIZATION

153 **Environment Representation.** We formalize embodied environments as directed graphs $G_t =$
 154 (V_t, E_t, A_t) that capture the essential structure of physical spaces. The node set V_t encompasses
 155 three entity types: spatial nodes representing rooms and areas, object nodes for interactive items,
 156 and agent nodes for autonomous entities. Each node maintains an attribute dictionary A_t storing

continuous physical properties such as weight, temperature, material composition, and geometric dimensions. The edge set E_t encodes spatial relationships through static containment relations (e.g., “in”, “on”) and dynamic proximity relations E_{near} that track which objects fall within an agent’s interaction range. This graph representation enables efficient reasoning about spatial constraints while avoiding the computational overhead of continuous 3D simulation.

Task Formalization. Each evaluation task is defined as a tuple $\mathcal{T} = (S_{\text{init}}, I, G_{\text{goal}}, \mathcal{A}_{\text{task}})$, where S_{init} specifies the initial environment state, I provides the natural language instruction, G_{goal} defines success conditions through logical predicates, and $\mathcal{A}_{\text{task}}$ identifies participating agents. The evaluation objective is to assess whether agents can generate an action sequence $\Pi = (\pi_1, \dots, \pi_T)$ that transforms the environment from S_{init} to a terminal state S_{final} satisfying all predicates in G_{goal} . This formalization captures both the planning and execution aspects of embodied reasoning.

3.2 HIERARCHICAL TASK TAXONOMY

Our evaluation framework organizes tasks along two orthogonal dimensions: agent configuration (single vs. multi-agent) and cognitive complexity (L1: basic, L2: intermediate, L3: advanced). This structure enables systematic assessment of how reasoning capabilities scale with task demands.

Single-Agent Tasks. Single-agent scenarios ($|\mathcal{A}_{\text{task}}| = 1$) isolate individual reasoning capabilities across three complexity levels. At the basic level, **Direct Command** tasks require straightforward instruction following, such as “place cup#1 on table#1,” establishing baseline comprehension abilities. Intermediate complexity introduces two parallel challenges: **Attribute Reasoning** tasks require comparing continuous properties to identify targets (e.g., “move the heaviest cup” requires solving $v^* = \arg \max_{v \in V_{\text{cups}}} A_t(v, \text{weight})$), while **Tool Use** tasks demand recognizing capability gaps and acquiring right tools. For instance, “clean the table” requires agents to identify that cleaning actions are unavailable in their base action set \mathcal{A}_i , locate cleaning tools, and execute $\text{grasp}(v_{\text{tool}})$ to dynamically expand their capabilities. Advanced **Compound Reasoning** tasks integrate multiple challenges, such as “clean the heaviest table,” requiring simultaneous attribute comparison, tool acquisition, and multi-step planning.

Multi-Agent Tasks. Multi-agent scenarios ($|\mathcal{A}_{\text{task}}| > 1$) evaluate coordination capabilities through parallel complexity progression. Basic **Explicit Collaboration** tasks provide clear coordination directives, such as “Agent A and Agent B cooperate to open the heavy cabinet,” testing fundamental synchronization abilities. Intermediate **Implicit Collaboration** removes explicit instructions, requiring agents to autonomously recognize when tasks exceed individual capabilities. For example, “move the dining table to the storage room” requires agents to infer that $A_t(v_{\text{table}}, \text{weight}) > C_{\text{max}}(i)$ for any individual agent i , necessitating collaborative effort. Advanced **Compound Collaboration** combines all elements, such as “cooperatively repair the malfunctioning television,” demanding tool acquisition, capability assessment, and coordinated execution.

3.3 EAR-SIM: EFFICIENT ENVIRONMENT SIMULATION

State Representation and Updates. EAR-Sim employs text-based environment modeling to achieve efficient simulation at scale. The graph structure G_t maintains spatial relationships through topological connections rather than continuous coordinates, eliminating expensive collision detection while preserving essential spatial constraints. State updates follow an incremental approach where actions modify only directly affected nodes and edges. For instance, when an agent executes $\text{GOTO}(\text{table})$, the system updates only the relevant proximity relations in E_{near} rather than recomputing global spatial relationships.

Dynamic Capability Management. A key innovation in EAR-Sim is the dynamic tool-capability binding system. Agent actions are partitioned into basic actions (movement, grasping, opening) available to all agents, and tool-dependent actions (cleaning, heating, repairing) that require specific tools. Each tool object maintains a `capability` attribute specifying which actions it enables. When an agent grasps a tool, the system dynamically binds the associated capabilities to the agent’s action set. Upon releasing the tool, these capabilities are automatically unbound. This mechanism

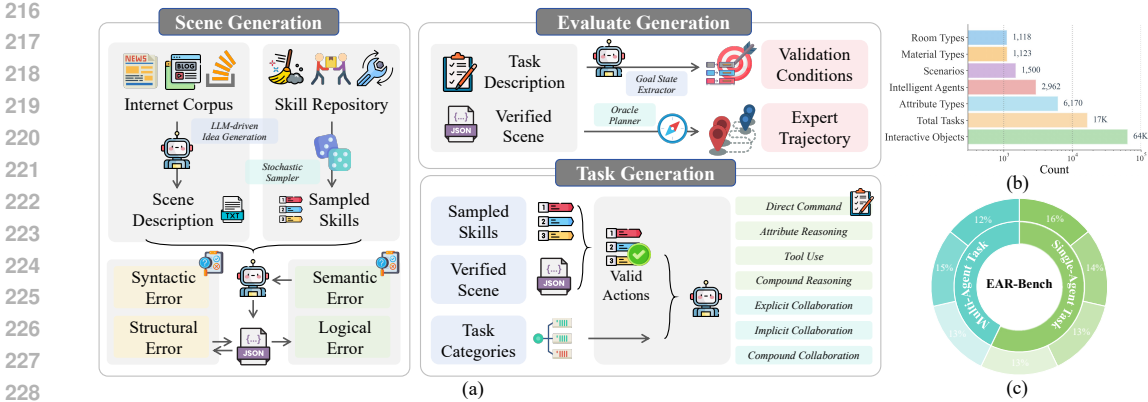


Figure 2: OmniEAR automated benchmark generation and evaluation framework. (a) Four-stage generation pipeline combining LLMs with rule-based validation: scene generation from internet corpus, task generation with skill sampling, evaluation logic extraction, and expert trajectory generation with human validation. (b) EAR-Bench statistics: 1,500 scenarios, 64K objects, 6K attribute types, spanning diverse domains and material compositions. (c) Balanced task distribution across seven categories spanning single-agent (Direct Command, Tool Use, Attribute Reasoning, Compound Reasoning) and multi-agent tasks (Explicit/Implicit/Compound Collaboration).

enables realistic modeling of how agents extend their abilities through tool use, moving beyond the static action spaces of existing benchmarks.

Emergent Collaboration. EAR-Sim supports collaboration that emerges from physical constraints rather than explicit programming. When agents attempt actions on objects whose properties exceed individual capabilities, the system enables collaboration request mechanisms. For instance, if an agent attempts to move an object where $A_t(v, \text{weight}) > C_{\max}(\text{agent})$, it can initiate collaboration by identifying suitable partners and coordinating joint actions. The system validates preconditions for all participating agents and maintains consistency throughout collaborative execution, ensuring realistic multi-agent interactions.

3.4 AUTOMATED BENCHMARK GENERATION

Generation Pipeline. Creating diverse, physically consistent scenarios at scale requires careful orchestration of neural generation and symbolic validation. As shown in 2, our pipeline operates in four stages, each combining the creative capabilities of large language models with rule-based consistency checking. This hybrid approach enables generating thousands of unique scenarios while maintaining physical realism and task solvability.

Scene and Task Generation. Scene generation begins with semantic seeds extracted from diverse text sources(Li et al., 2024b), which guide a neural generator g_{scene} in creating structured environment descriptions. The generator, implemented using high-temperature language models for diversity, produces initial scenes S_0 containing objects, spatial layouts, and agent configurations. Task generation follows a two-stage process: first, an environment analyzer C_{env} extracts feasible actions based on the scene structure, then a task generator g_{task} creates instructions anchored in physical possibilities. This grounding prevents generation of impossible tasks while maintaining creative diversity.

Evaluation Logic and Trajectories. For each generated task, we automatically derive evaluation criteria by parsing the instruction and scene to extract minimal state changes required for success. This produces a goal predicate set G_{goal} that serves as an objective success measure. Expert trajectories are generated using oracle agents with complete environmental knowledge, creating high-quality demonstrations for each task. These trajectories undergo filtering to remove suboptimal sequences, providing ideal solutions for comparison and learning.

Quality Assurance. As shown in Figure 3, generated scenarios undergo a three-stage validation pipeline. Automated validators check structural consistency and physical feasibility, correcting 87.6% of detected errors. Trained human annotators then verify scenarios against explicit criteria for instruction clarity, feasibility, and solvability (Cohen’s $\kappa = 0.84$). Finally, oracle agents validate expert trajectories through simulation replay. This rigorous process yielded 1,500 validated scenarios from 2,100 initial candidates (71.4% acceptance rate). Detailed procedures, annotation interface, and quality threshold examples are provided in Appendix 6.3.

3.5 BENCHMARK STATISTICS AND COVERAGE

EAR-Bench encompasses 1,500 scenarios across 11 domains including laboratory (39%), office (19%), industrial (12%), and medical environments, containing 64,057 interactive objects with rich physical properties. The dataset maintains careful balance across our task taxonomy: 65% single-agent tasks spanning all complexity levels, and 35% multi-agent tasks with emphasis on implicit collaboration scenarios that require genuine reasoning about coordination needs. With 6,381 distinct property types and 214 action types, EAR-Bench provides comprehensive coverage of embodied reasoning challenges while maintaining tractable evaluation scope. Detailed statistics are provided in Appendix 6.4.

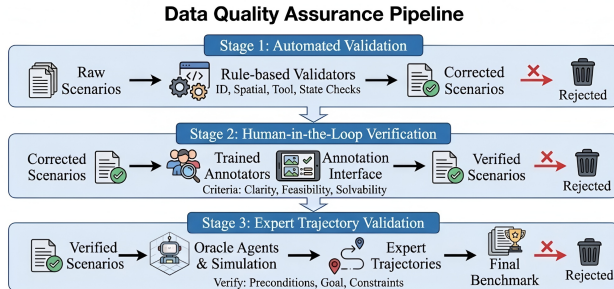


Figure 3: Overview of the Data Quality Assurance Pipeline. The pipeline consists of three sequential stages: (1) Automated Validation for rule-based error correction, (2) Human-in-the-Loop Verification for semantic and physical feasibility checks, and (3) Expert Trajectory Validation to ensure genuine solvability via oracle agents.

4 EXPERIMENTS

We systematically evaluate current LLMs on EAR-Bench to assess their physical reasoning capabilities in embodied tasks. Our experiments examine: (1) How performance degrades when models must dynamically acquire tools and determine coordination requirements from task contexts, (2) Whether model scale and architectural choices affect constraint-based reasoning capabilities, and (3) How environmental information presentation and training approaches impact autonomous decision-making in embodied scenarios.

4.1 EXPERIMENTAL SETUP

Model Selection. We evaluate nine representative models spanning three architectural paradigms. Closed-source models include GPT-4o (Hurst et al., 2024) and Gemini-2.5-Flash (Comanici et al., 2025), representing current commercial state-of-the-art. Open-source foundation models cover a wide parameter range: Deepseek-V3 (Liu et al., 2024) at 671B parameters, the Qwen2.5 series (Team, 2024) at 3B, 7B, and 72B parameters, and Llama3.1-8B (Touvron et al., 2023). This selection enables analysis of how model scale affects embodied reasoning. We also include reasoning-specialized models: Deepseek-R1 (Guo et al., 2025) and QwQ-32B (Li et al., 2024b), which employ explicit chain-of-thought reasoning during inference.

Evaluation Protocol. All models undergo identical evaluation to ensure fair comparison. We implement partial observability where agents must explore environments to discover object locations and properties, reflecting realistic deployment conditions. Each model completes 2,800 test scenarios across seven task categories with three independent runs for statistical reliability. We standardize prompts, environment descriptions, and action vocabularies across all models, with tool-dependent actions dynamically enabled based on context. This design ensures performance differences reflect reasoning capabilities rather than implementation artifacts. Detailed experimental configurations are provided in Appendix 6.8.

Model	Single-Agent Tasks								Multi-Agent Tasks					
	Direct Command		Tool Use		Attribute Reasoning		Compound Reasoning		Explicit Collab.		Implicit Collab.		Compound Collab.	
	SR	Step	SR	Step	SR	Step	SR	Step	SR	Step	SR	Step	SR	Step
<i>Closed-source Models</i>														
GPT-4o	96.6	12.9	80.0	13.6	77.8	12.3	69.2	14.5	90.0	13.9	<u>77.5</u>	14.4	32.0	22.9
Gemini-2.5-Flash	90.5	11.0	82.3	16.5	56.3	17.5	59.4	20.0	88.5	8.4	85.5	7.1	40.5	16.2
<i>Reasoning-specialized Models</i>														
Deepseek-R1	94.1	10.3	85.8	14.1	41.9	12.2	70.6	16.2	92.0	7.4	84.5	9.6	48.5	12.5
QwQ-32B	85.2	10.3	73.4	13.0	44.9	11.0	54.1	13.6	88.0	8.5	84.0	8.3	36.5	19.0
<i>Open-source Foundation Models</i>														
Deepseek-V3	91.1	11.2	82.3	15.1	56.3	10.3	67.1	16.0	82.0	9.4	63.0	9.7	36.0	20.2
Qwen2.5-72B	89.7	14.7	56.4	21.7	57.4	17.2	66.7	21.1	56.0	24.1	65.4	15.6	28.6	29.5
Llama3.1-8B	24.9	34.4	8.3	34.6	9.9	34.8	12.4	34.3	4.0	3.5	1.5	2.1	0.0	3.4
Qwen2.5-7B	40.2	24.1	15.4	31.7	22.2	26.6	16.5	30.5	38.5	25.0	13.5	24.1	1.0	27.2
Qwen2.5-3B	0.6	30.5	1.8	31.3	0.6	34.0	2.9	32.9	8.5	20.4	1.5	16.3	0.5	16.8
+ SFT	76.3	15.4	45.0	24.7	33.5	22.8	36.5	24.7	22.5	29.2	5.5	28.3	1.0	27.1

Table 1: Performance across task categories. Success Rate (SR) measures task completion percentage, Step Count indicates average actions for successful completion. Bold indicates best in category, underline shows overall best.

Fine-tuning Configuration. To assess whether supervised learning can address reasoning limitations, we fine-tune Qwen2.5-3B on expert trajectories. We collect 1,942 successful demonstrations from Qwen2.5-72B with complete environmental access, filtering for optimal action sequences. The resulting 20,346 instruction-action pairs train the model using standard causal language modeling objectives, testing whether smaller models can learn embodied reasoning patterns from larger models. Complete hyperparameters are listed in Appendix 6.6.

Deployment Configurations. We evaluate models in two configurations. Single-agent scenarios test individual reasoning capabilities without collaborative complexity. Multi-agent scenarios employ centralized coordination where one model controls all agents with complete state visibility, isolating collaborative reasoning from communication challenges. This design choice allows us to assess pure multi-agent reasoning capabilities without confounding factors from limited observability or communication protocols.

4.2 MAIN RESULTS

Table 1 presents comprehensive evaluation results across our task hierarchy. The results reveal systematic performance patterns that validate our framework design and expose fundamental limitations in current models.

Task Complexity Hierarchy. Figure 4 reveals systematic performance degradation across our task hierarchy, with success rates declining from 85.2-96.6% on Direct Commands to 32.0-48.5% on Compound Collaboration tasks. This consistent pattern confirms that performance differences reflect reasoning complexity rather than task difficulty alone. Tool Use (73.4-85.8%) requires recognizing capability gaps from context, while Attribute Reasoning (41.9-77.8%) demands grounding language in physical properties. Both involve inferring requirements from environmental constraints rather than following explicit instructions. Notably, Explicit Collaboration outperforms several single-agent tasks, indicating that reasoning about physical constraints poses greater challenges than multi-agent coordination when guidance is provided. The severe performance drop in compound tasks demonstrates that current models cannot integrate multiple constraints simultaneously, supporting our framework’s focus on autonomous inference from physical context as the key determinant of embodied reasoning difficulty.

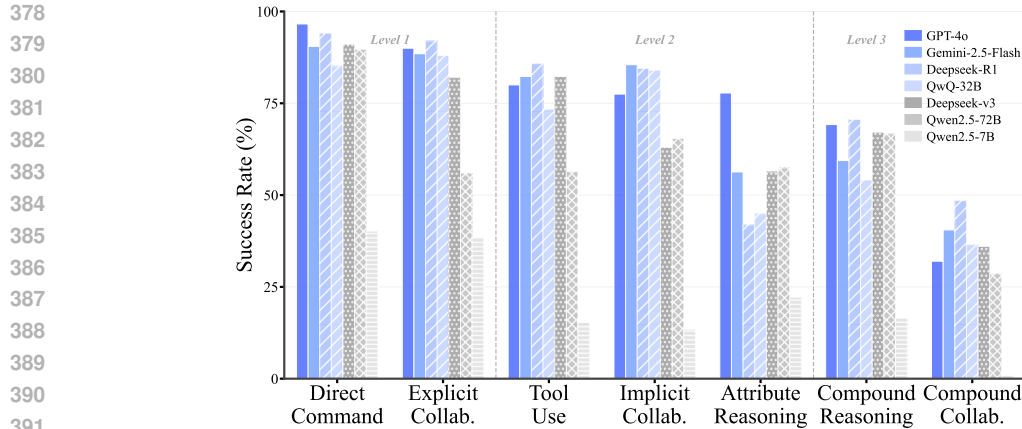


Figure 4: Performance comparison across task categories demonstrating systematic difficulty hierarchy and distinct model performance tiers.

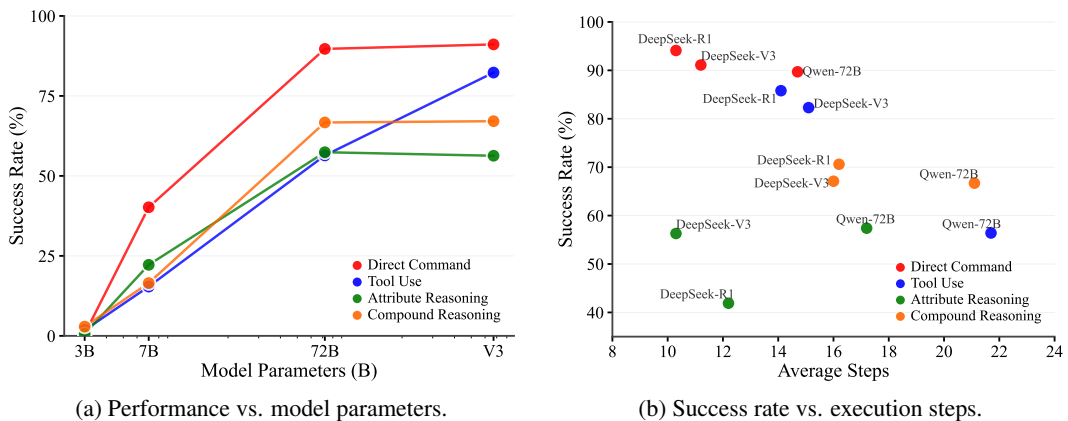


Figure 5: Scaling patterns reveal distinct thresholds for embodied reasoning capabilities. (a) Direct Command and Tool Use scale sharply with parameters while Attribute/Compound Reasoning plateau early. (b) Reasoning-specialized models achieve higher success through longer execution paths.

Model Scale and Reasoning Capabilities. Figure 5a reveals distinct scaling patterns across task types. While Direct Command performance improves sharply with model size (from near-zero at 3B to over 90% at 72B), tasks requiring physical constraint reasoning show more complex relationships. Tool Use exhibits similar steep scaling, suggesting that maintaining multi-step plans for capability acquisition correlates strongly with model capacity. However, Attribute Reasoning and Compound Reasoning plateau earlier, with diminishing returns beyond 72B parameters. This differential scaling indicates that raw parameter count enables better execution and planning but does not necessarily improve understanding of physical properties.

Table 1 provides further evidence distinguishing execution capability from genuine reasoning. Reasoning-specialized models like Deepseek-R1 achieve the highest performance on Compound Collaboration (48.5%) despite lower scores on Attribute Reasoning (41.9%) compared to GPT-4o (77.8%). This performance inversion suggests these models excel at explicit logical planning but struggle with grounding abstract properties in physical contexts. The success rate versus step count trade-off in Figure 5b reinforces this interpretation: reasoning models achieve higher success through longer, more deliberate execution paths rather than efficient understanding of constraints. Fine-tuning results provide the clearest evidence that current models lack true embodied reasoning: while Qwen2.5-3B improves dramatically on single-agent tasks through imitation (0.6% to 76.3%), multi-agent performance remains negligible (1.5% to 5.5%), indicating that learned

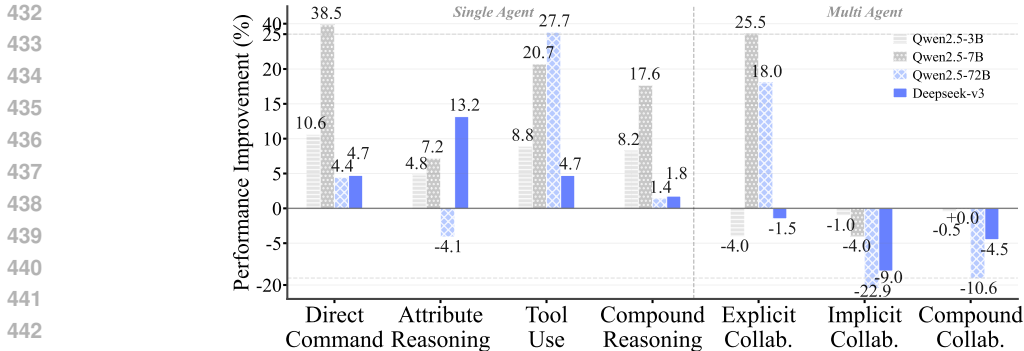


Figure 6: Performance changes with World Graph enhancement. Tool Use and Attribute Reasoning benefit substantially, while Implicit Collaboration shows degradation, suggesting information overload effects.

behaviors cannot generalize to scenarios requiring autonomous assessment of physical constraints and coordination needs.

4.3 DETAILED ANALYSIS

We conduct analyses to understand the factors driving model performance and identify specific capability bottlenecks.

Environmental Representation Impact.

Table 2 and Figure 6 reveal task-specific effects of structured environmental knowledge. Tool Use benefits most significantly (up to 27.7% improvement), as World Graph transforms spatial search into direct tool selection. Smaller models gain more than larger ones, suggesting that full environmental knowledge compensates for limited working memory. Conversely, Implicit Collaboration consistently drops with World Graph across all model scales. This counterintuitive pattern indicates that exploration-based discovery helps models focus on task-relevant constraints, while complete information introduces distraction. The divergent effects across task types demonstrate that optimal information presentation depends on reasoning requirements, not information quantity.

Task	3B		7B		72B		671B	
	w/o	w/	w/o	w/	w/o	w/	w/o	w/
Direct Cmd	0.6	11.2	40.2	78.7	89.7	94.1	91.1	95.9
Tool Use	1.8	10.7	15.4	36.1	56.4	84.0	82.3	87.0
Attr. Reas.	0.6	5.4	22.2	29.3	57.4	53.3	56.3	69.5
Comp. Reas.	2.9	11.2	16.5	34.1	64.5	65.9	67.1	68.8
Expl. Coll.	8.5	4.5	38.5	64.0	62.5	80.5	82.0	80.5
Impl. Coll.	1.5	0.5	13.5	9.5	65.4	42.5	63.0	54.0
Comp. Coll.	0.5	0.0	1.0	1.0	28.6	18.0	36.0	31.5

Table 2: Success rates (%) with and without World Graph enhancement across model scales, revealing task-specific gains and unexpected drop in implicit collaboration.

Computational Efficiency Trade-offs. Figure 7a identifies three efficiency regimes with distinct cost-performance profiles. Foundation models achieve moderate performance with minimal tokens (456-1400), while commercial models trade higher token usage (1817-2457) for improved success rates. Reasoning models consume up to 12,000 tokens but excel on complex tasks. The efficiency frontier shifts dramatically between single and multi-agent scenarios: Gemini-2.5-Flash optimizes single-agent efficiency, but Deepseek-R1 becomes necessary for multi-agent tasks despite 75% higher costs. This shift reflects the irreducible computational complexity of modeling multiple agent states and coordination protocols, suggesting no universal optimization exists across task types.

Execution Efficiency Analysis. Figure 8 compares model solutions to expert demonstrations via Relative Step Ratios (RSR = L_{expert}/L_{model}). Single-agent tasks show consistent moderate efficiency (median RSR 0.40-0.55), while multi-agent tasks exhibit both lower efficiency and higher variance, reflecting uncertainty in coordination timing and strategy selection. Compound Collaboration

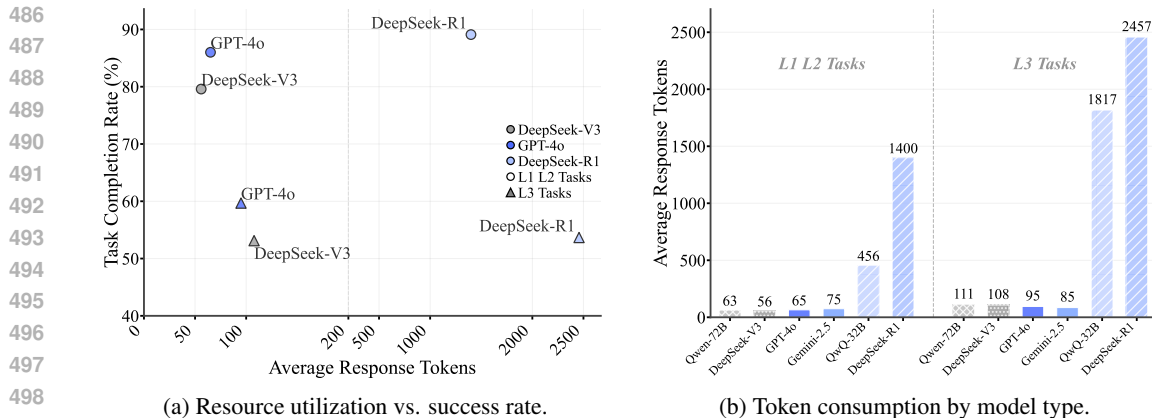


Figure 7: Reasoning-specialized models achieve higher performance through increased computational overhead. (a) Efficiency-performance trade-offs across model architectures. (b) Token consumption patterns revealing computational costs of reasoning approaches.

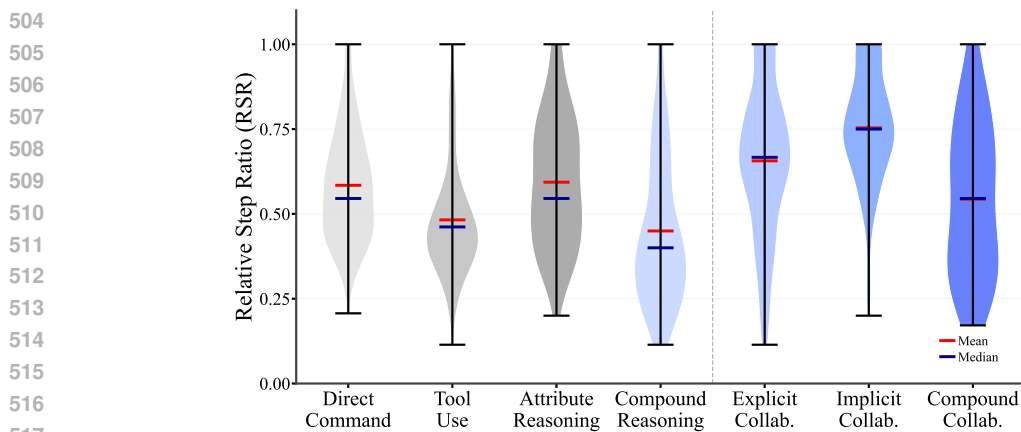


Figure 8: Relative Step Ratio distributions showing execution efficiency compared to expert trajectories. Multi-agent tasks show both lower efficiency and higher variance than single-agent tasks.

reveals a striking bimodal distribution: models either adopt simple sequential execution or attempt complex parallel coordination, with no successful middle strategies. This polarization suggests current models lack adaptive coordination mechanisms, defaulting to extreme approaches rather than selecting strategies based on task constraints.

5 CONCLUSION

We presented OmniEAR, a benchmark for evaluating embodied agent reasoning through 1,500 scenarios requiring inference from physical constraints. Our evaluation reveals that current models show severe performance degradation when moving from explicit instructions to constraint-based reasoning, with performance dropping from over 85% to below 65% across tool usage and coordination tasks. We identify critical parameter thresholds for maintaining multi-step plans, paradoxical effects of environmental information on coordination, and the inability of fine-tuning to address multi-agent reasoning gaps. Results demonstrate that embodied reasoning requires fundamentally different computational mechanisms than those underlying current language models. OmniEAR provides systematic diagnostics of these limitations and a rigorous platform for developing next-generation embodied AI systems. We discuss broader implications and future research directions in Appendix 6.7.

REFERENCES

- 540
541
542 Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea
543 Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine
544 Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally
545 Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee,
546 Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka
547 Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander
548 Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy
549 Zeng. Do as i can, not as i say: Grounding language in robotic affordances, 2022. URL <https://arxiv.org/abs/2204.01691>.
- 550
551 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
552 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
553 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 554
555 Matthew Chang, Gunjan Chhablani, Alexander Clegg, Mikael Dallaire Cote, Ruta Desai, Michal
556 Hlavac, Vladimir Karashchuk, Jacob Krantz, Roozbeh Mottaghi, Priyam Parashar, Siddharth
557 Patki, Ishita Prasad, Xavier Puig, Akshara Rai, Ram Ramrakhya, Daniel Tran, Joanne Truong,
558 John M. Turner, Eric Undersander, and Tsung-Yen Yang. Partnr: A benchmark for planning and
559 reasoning in embodied multi-agent tasks, 2024. URL <https://arxiv.org/abs/2411.00081>.
- 560
561 Sadia Sultana Chowa, Riasad Alvi, Subhey Sadi Rahman, Md Abdur Rahman, Mohaimenul
562 Azam Khan Raiaan, Md Rafiqul Islam, Mukhtar Hussain, and Sami Azam. From language to
563 action: A review of large language models as autonomous agents and tool users, 2025. URL
<https://arxiv.org/abs/2508.17281>.
- 564
565 Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit
566 Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing
567 the frontier with advanced reasoning, multimodality, long context, and next generation agentic
568 capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- 569
570 Pascale Fung, Yoram Bachrach, Asli Celikyilmaz, Kamalika Chaudhuri, DeLong Chen, Willy Chung,
571 Emmanuel Dupoux, Hongyu Gong, Hervé Jégou, Alessandro Lazaric, Arjun Majumdar, Andrea
572 Madotto, Franziska Meier, Florian Metze, Louis-Philippe Morency, Théo Moutakanni, Juan Pino,
573 Basile Terver, Joseph Tighe, Paden Tomasello, and Jitendra Malik. Embodied ai agents: Modeling
574 the world, 2025. URL <https://arxiv.org/abs/2506.22355>.
- 575
576 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
577 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
578 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 579
580 Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan
581 Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through
582 planning with language models. *arXiv preprint arXiv:2207.05608*, 2022.
- 583
584 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark,
585 AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint*
586 *arXiv:2410.21276*, 2024.
- 587
588 Li Kang, Xiufeng Song, Heng Zhou, Yiran Qin, Jie Yang, Xiaohong Liu, Philip Torr, Lei Bai, and
589 Zhenfei Yin. Viki-r: Coordinating embodied multi-agent cooperation via reinforcement learning,
590 2025. URL <https://arxiv.org/abs/2506.09049>.
- 591
592 Chengshu Li, Fei Xia, Roberto Martín-Martín, Michael Lingelbach, Sanjana Srivastava, Bokui
593 Shen, Kent Vainio, Cem Gokmen, Gokul Dharan, Tanish Jain, et al. igibson 2.0: Object-centric
simulation for robot learning of everyday household tasks. *arXiv preprint arXiv:2108.03272*,
2021.
- Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín-Martín, Chen Wang, Gabriel Levine, Wensi Ai, Benjamin Martinez, et al. Behavior-1k: A human-centered, embodied ai benchmark with 1,000 everyday activities and realistic simulation. *arXiv preprint arXiv:2403.09227*, 2024a.

- 594 Dingming Li, Hongxing Li, Zixuan Wang, Yuchen Yan, Hang Zhang, Siqi Chen, Guiyang Hou,
595 Shengpei Jiang, Wenqi Zhang, Yongliang Shen, Weiming Lu, and Yueting Zhuang. Viewspatial-
596 bench: Evaluating multi-perspective spatial localization in vision-language models, 2025. URL
597 <https://arxiv.org/abs/2505.21500>.
- 598
599 Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik
600 Bansal, Etash Guha, Sedrick Scott Keh, Kushal Arora, et al. Datacomp-1m: In search of the
601 next generation of training sets for language models. *Advances in Neural Information Processing
602 Systems*, 37:14200–14282, 2024b.
- 603 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
604 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint
605 arXiv:2412.19437*, 2024.
- 606
607 Zhao Mandi, Shreya Jain, and Shuran Song. Roco: Dialectic multi-robot collaboration with large
608 language models. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*,
609 pp. 286–299. IEEE, 2024.
- 610 Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba.
611 Virtualhome: Simulating household activities via programs, 2018. URL [https://arxiv.
612 org/abs/1806.07011](https://arxiv.org/abs/1806.07011).
- 613
614 Xavier Puig, Eric Undersander, Andrew Szot, Mikael Dallaire Cote, Tsung-Yen Yang, Ruslan
615 Partsey, Ruta Desai, Alexander William Clegg, Michal Hlavac, So Yeon Min, et al. Habitat
616 3.0: A co-habitat for humans, avatars and robots. *arXiv preprint arXiv:2310.13724*, 2023.
- 617 Neil Rabinowitz, Frank Perbet, Francis Song, Chiyuan Zhang, SM Ali Eslami, and Matthew
618 Botvinick. Machine theory of mind. In *International conference on machine learning*, pp. 4218–
619 4227. PMLR, 2018.
- 620
621 Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng
622 Li, and Yueting Zhuang. Taskbench: Benchmarking large language models for task automation.
623 In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang
624 (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 4540–4574. Curran
625 Associates, Inc., 2024. URL [https://proceedings.neurips.cc/paper_files/
626 paper/2024/file/085185ea97db31ae6dcac7497616fd3e-Paper-Datasets_
627 and_Benchmarks_Track.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/085185ea97db31ae6dcac7497616fd3e-Paper-Datasets_ and_Benchmarks_Track.pdf).
- 628
629 Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi,
630 Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions
631 for everyday tasks, 2020. URL <https://arxiv.org/abs/1912.01734>.
- 632
633 Nan Sun, Chengming Shi, et al. [aml] interactgen: Enhancing human-involved embodied task
634 reasoning through llm-based multi-agent collaboration. 2024.
- 635
636 Qwen Team. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.
- 637
638 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
639 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
640 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 641
642 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan,
643 and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models,
644 2023. URL <https://arxiv.org/abs/2305.16291>.
- 645
646 Yuntao Wang, Yanghe Pan, Quan Zhao, Yi Deng, Zhou Su, Linkang Du, and Tom H Luan. Large
647 model agents: State-of-the-art, cooperation paradigms, security and privacy, and future trends.
arXiv e-prints, pp. arXiv–2409, 2024.
- 648
649 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
650 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
651 neural information processing systems*, 35:24824–24837, 2022.

648 Zhenyu Wu, Ziwei Wang, Xiuwei Xu, Jiwen Lu, and Haibin Yan. Embodied task planning with
649 large language models, 2023. URL <https://arxiv.org/abs/2307.01848>.
650

651 Rui Yang, Hanyang Chen, Junyu Zhang, Mark Zhao, Cheng Qian, Kangrui Wang, Qineng Wang,
652 Teja Venkat Koripella, Marziyeh Movahedi, Manling Li, et al. Embodiedbench: Comprehensive
653 benchmarking multi-modal large language models for vision-driven embodied agents. *arXiv*
654 *preprint arXiv:2502.09560*, 2025.

655 Andy Zeng, Maria Attarian, Brian Ichter, Krzysztof Choromanski, Adrian Wong, Stefan Welker,
656 Federico Tombari, Aveek Purohit, Michael Ryoo, Vikas Sindhwani, Johnny Lee, Vincent
657 Vanhoucke, and Pete Florence. Socratic models: Composing zero-shot multimodal reasoning
658 with language, 2022. URL <https://arxiv.org/abs/2204.00598>.

659 Hongxin Zhang, Weihua Du, Jiaming Shan, Qinhong Zhou, Yilun Du, Joshua B. Tenenbaum,
660 Tianmin Shu, and Chuang Gan. Building cooperative embodied agents modularly with large
661 language models, 2024. URL <https://arxiv.org/abs/2307.02485>.
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

6 APPENDIX

Dataset	Scenes	Domain	Task Types	Actions	Action Space	Collab.	Auto Gen.
ALFRED	120	House	D	7	Static	—	×
PARTNR	60	House	D	11	Static	Effic.	✓
BEHAVIOR-1K	50	Diverse	D,T	6	Static	—	×
WAH	7	House	D	10	Static	Effic.	×
TDW-MAT	6	House	D,E	7	Static	Effic.	×
C-WAH	6	House	D,E	7	Static	Effic.	×
Overcooked	5	Kitchen	E,I,C	6	Static	Effic.	×
OmniEAR	1.5K	Diverse	D,A,T,R,E,I,C	218	Dynamic	Phys.	✓

Table 3: Comparison of embodied AI datasets and benchmarks. Task types: D (Direct Command), A (Attribute Reasoning), T (Tool Use), R (Compound Reasoning), E (Explicit Collaboration), I (Implicit Collaboration), C (Compound Collaboration). Actions: number of available action types. Collab.: collaboration mechanism (Effic. = efficiency-based, Phys. = physical necessity-driven). Auto Gen.: automated task generation capability. Our framework uniquely combines comprehensive task coverage, dynamic action spaces, physical necessity-driven collaboration, and scalable automated generation.

6.1 RELATED WORK

Embodied Intelligence Benchmarks The embodied intelligence evaluation landscape has established diverse benchmark frameworks spanning navigation to complex manipulation tasks (Puig et al., 2023; Li et al., 2021). ALFRED (Shridhar et al., 2020) provides foundational standards for instruction-following task evaluation, while BEHAVIOR-1K (Li et al., 2024a) extends coverage to 1,000 daily activity scenarios. These benchmarks effectively assess task execution capabilities, yet physical property modeling predominantly employs discrete state representations, such as binary door operations and object pickup/placement, with limited requirements for reasoning about continuous attributes including weight, hardness, and temperature. Our framework addresses this limitation by introducing continuous physical property reasoning tasks that require agents to compare object attributes and make decisions based on physical constraints.

Embodied Tool Use Tool usage evaluation in embodied AI exhibits stratified characteristics across different complexity levels. RoCo (Mandi et al., 2024) focuses on low-level manipulation skills such as grasping precision, while high-level benchmarks like PARTNR (Chang et al., 2024) adopt predefined tool configurations with agent action spaces fixed at task initialization. This design effectively simplifies evaluation complexity but presents limitations in assessing dynamic tool reasoning capabilities based on task requirements. Current approaches typically provide static tool sets (Chowa et al., 2025; Fung et al., 2025), preventing evaluation of how agents should reason about capability gaps and tool acquisition needs. Our framework introduces dynamic tool acquisition mechanisms, requiring agents to autonomously infer tool requirements and expand their action spaces based on task demands, thereby supplementing existing evaluation dimensions.

Multi-Agent Collaboration Multi-agent embodied intelligence evaluation has emerged as a significant research direction, with related work achieving valuable progress in collaboration modeling (Sun et al., 2024; Wang et al., 2024). PARTNR evaluates multi-agent planning capabilities through heterogeneous task design, TDW-MAT (Zhang et al., 2024) creates collaborative scenarios using load capacity constraints, and EmbodiedBench (Yang et al., 2025) focuses on task allocation and execution optimization. Existing approaches primarily model collaboration requirements through two pathways: explicit collaboration instructions that clearly specify inter-agent task division, and efficiency optimization that drives multi-agent participation to enhance task completion speed. However, real-world collaboration decisions often stem from physical constraints rather than external instructions or efficiency considerations. Our framework employs implicit collaboration design requiring agents to autonomously assess whether tasks exceed single-agent capability ranges based on physical constraints and determine collaboration strategies accordingly, transforming collaboration judgment from external instructions to constraint-driven internal reasoning processes.

6.2 TASK CATEGORY EXAMPLES

Figure 9 illustrates our hierarchical task taxonomy with representative examples, organized by agent configuration (single-agent vs. multi-agent) and cognitive complexity level (L1-L3). The figure demonstrates the progression from basic instruction following to advanced reasoning requiring dynamic capability acquisition and autonomous coordination. Table 4 provides detailed descriptions of the key reasoning requirements for each task category.

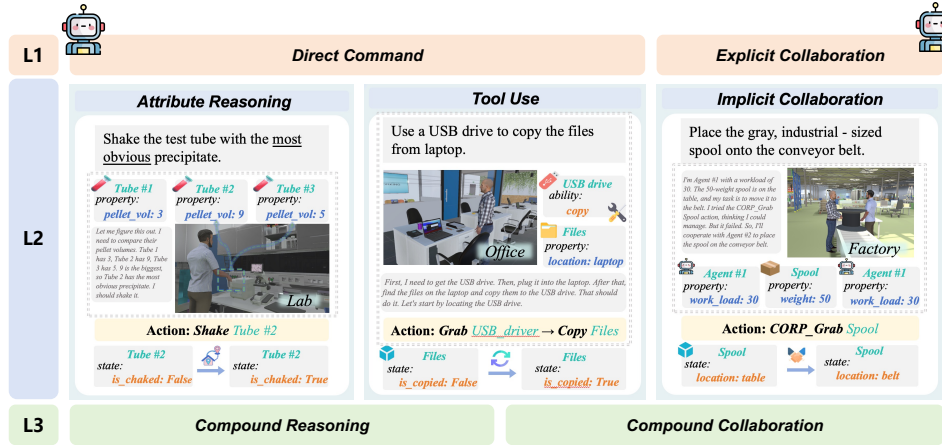


Figure 9: Representative examples for each task category in EAR-Bench. Left column shows single-agent tasks progressing from Direct Command (L1) through Attribute Reasoning and Tool Use (L2) to Compound Reasoning (L3). Right column shows multi-agent tasks from Explicit Collaboration (L1) through Implicit Collaboration (L2) to Compound Collaboration (L3). Each example includes the task instruction, relevant object properties, agent reasoning process, and resulting state transitions.

Level	Task Type	Example Task	Key Reasoning Requirements
<i>Single-Agent Tasks</i>			
L1	Direct Command	Place the red cup on the kitchen table.	Basic object manipulation and spatial understanding through straightforward instruction following.
L2	Tool Use	Clean the dirty table in the living room.	Recognizing capability gaps, locating appropriate tools, and dynamically expanding action capabilities through tool acquisition.
L2	Attribute Reasoning	Move the heaviest box to the storage room.	Comparing continuous physical properties across multiple objects to identify correct manipulation targets.
L3	Compound Reasoning	Clean the heaviest table in the room.	Integrating attribute comparison, tool acquisition, and multi-step planning simultaneously.
<i>Multi-Agent Tasks</i>			
L1	Explicit Collab.	Agent A and Agent B cooperate to move the heavy dining table.	Fundamental multi-agent synchronization with explicit coordination directives provided.
L2	Implicit Collab.	Move the piano to the music room.	Autonomously recognizing when tasks exceed individual capabilities without explicit coordination instructions.
L3	Compound Collab.	Cooperatively repair the malfunctioning television.	Combining tool acquisition, capability assessment, and coordinated execution with autonomous collaboration recognition.

Table 4: Representative examples for each task category in EAR-Bench. Tasks span single-agent scenarios (Direct Command, Tool Use, Attribute Reasoning, Compound Reasoning) and multi-agent scenarios (Explicit, Implicit, and Compound Collaboration) across three complexity levels.

6.3 DATA QUALITY ASSURANCE AND VALIDATION

A key challenge in LLM-based benchmark generation is ensuring that generated scenarios are physically consistent, unambiguous, and genuinely solvable (Li et al., 2025; Shen et al., 2024). To address this challenge and ensure benchmark reliability at scale, we implement a rigorous three-stage validation pipeline that combines automated error correction, systematic human-in-the-loop verification, and simulation-based trajectory validation, as illustrated in Figure 3. This hybrid approach balances scalability with quality control, achieving validation standards comparable to manually curated benchmarks while enabling generation at scale.

6.3.1 STAGE 1: AUTOMATED VALIDATION

The first stage addresses systematic errors that can be detected and corrected programmatically, reducing the burden on human annotators while ensuring basic structural integrity.

Error Analysis and Correction. Systematic analysis of 2,100 initially generated scenarios identified five primary error categories, summarized in Table 5. We developed specialized rule-based validators to automatically detect and correct each error type: ID mismatches are resolved via semantic similarity matching (cosine threshold 0.7), spatial violations are corrected via dimension checking and constraint propagation, tool-capability conflicts are resolved via attribute cross-referencing against our capability ontology, and state inconsistencies are handled via priority-based resolution that preserves task-critical states. This automated pipeline successfully corrected 87.6% of detected errors. The remaining 12.4% of scenarios—primarily those with complex logical inconsistencies requiring semantic understanding—were rejected, yielding 1,804 candidates for human review.

Error Type	Frequency	Auto-Correctable
ID Mismatch	32.4%	91%
Invalid Spatial Relations	28.7%	85%
Tool-Capability Conflict	23.1%	89%
State Inconsistency	10.3%	94%
Other Logical Errors	5.5%	62%
Overall	100%	87.6%

Table 5: Distribution of LLM generation errors and automatic correction rates by category.

6.3.2 STAGE 2: HUMAN-IN-THE-LOOP VERIFICATION

While automated validation effectively handles systematic errors, subtle semantic ambiguities, non-obvious unsolvability, and nuanced physical inconsistencies require human expertise. We therefore implement rigorous human-in-the-loop verification with carefully selected and trained annotators, standardized evaluation criteria, and quantitative quality control measures.

Annotator Selection and Training. Five annotators were recruited from graduate programs in robotics and embodied AI, each with at least two years of research experience in related domains. All annotators completed a standardized 4-hour training program comprising three components: conceptual training on embodied reasoning principles and physical constraints (1.5 hours), hands-on system familiarization with environment representation and the validation interface (1 hour), and calibration exercises on 20 expert-validated scenarios with group discussion and feedback (1.5 hours). The calibration scenarios were independently validated by two senior researchers to establish reliable ground truth labels. Annotators were required to achieve Cohen’s $\kappa \geq 0.75$ against these gold labels before beginning independent annotation; those falling below this threshold received additional training until reaching the required level. During the main annotation phase, each annotator reviewed approximately 720 scenarios (with systematic overlap for agreement computation), averaging 3-4 minutes per scenario to ensure careful evaluation.

Annotation Interface. To ensure consistent and reliable human evaluation, we developed a standardized annotation interface illustrated in Figure 10. The interface presents three integrated components: (1) a structured scene representation panel displaying room layouts, object properties (weight, material, dimensions), spatial relationships, and agent configurations in a clear hierarchical format; (2) a task instruction panel with automatically highlighted key elements including target objects, required actions, and relevant constraints; and (3) an evaluation panel presenting the three assessment criteria as binary selections along with a mandatory text field for rejection justification. This standardized interface ensures consistent evaluation conditions across all annotators and scenarios, minimizing variance arising from interface interpretation differences and enabling systematic collection of rejection rationales for error analysis.

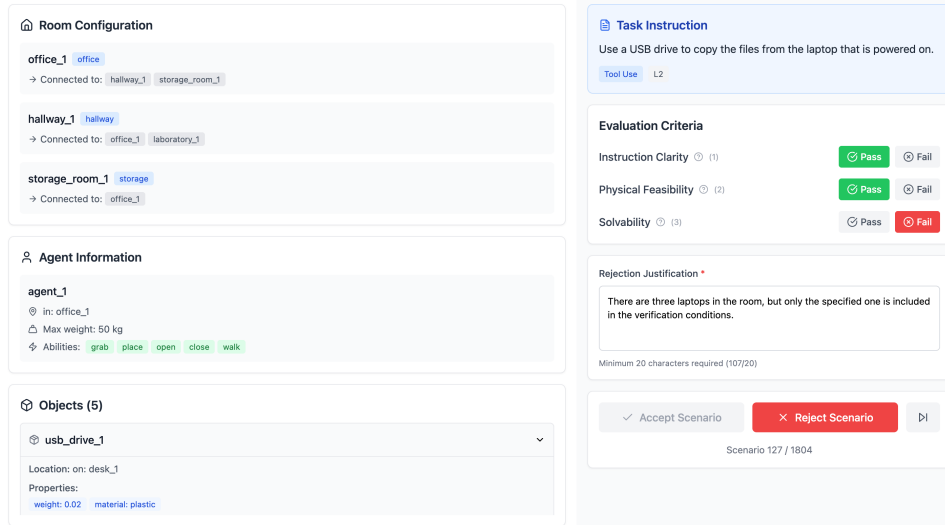


Figure 10: Human annotation interface for scenario verification. The left panel displays structured scene representation including room configuration, agent properties, and interactive objects with their attributes. The right panel shows the task instruction with highlighted key elements (top), three evaluation criteria with binary pass/fail selections (middle), and a rejection justification field requiring written rationale for any failed criterion (bottom).

Annotation Guidelines. Annotators evaluated each scenario against three binary criteria with explicit operational definitions designed to ensure consistent judgment:

- **Instruction Clarity:** The task instruction unambiguously specifies the goal state, target objects, and required actions without relying on implicit assumptions or domain-specific knowledge not provided in the scene description.
- **Physical Feasibility:** All actions required to complete the task respect the physical constraints encoded in the environment, including weight limits, spatial accessibility, tool requirements, and material properties.
- **Solvability:** A valid action sequence exists that transforms the initial environment state to a goal state satisfying all success criteria within the maximum step limit of 25 steps.

Each scenario received independent review from two randomly assigned annotators to eliminate individual bias. Acceptance required unanimous approval on all three criteria; any rejection required written justification specifying the failed criterion and the specific issue identified, enabling systematic failure analysis and iterative improvement of the generation pipeline.

Quality Threshold Examples. To provide concrete grounding for our annotation guidelines and demonstrate the practical application of our quality criteria, Table 6 presents representative examples of accepted and rejected scenarios across different task categories. These contrastive examples illustrate the key distinctions between scenarios meeting our quality standards and those rejected for ambiguity, physical infeasibility, or unclear collaboration requirements. The examples were selected to highlight common failure modes and the reasoning behind acceptance decisions.

Task Type	Instruction	Decision	Rationale
Tool Use	Clean the dusty shelf using supplies from the cabinet which is located in the corner of the room next to the window.	✓	Clear goal specification, explicit tool source, verifiably solvable path exists, and the operation steps are clear and easy to execute.
Tool Use	Clean everything in the room including the floor, the ceiling, the furniture and all the small ornaments on the shelf.	✗	Ambiguous task scope; “everything” lacks clear termination criteria and target enumeration, which may lead to endless work and cannot be completed effectively.
Attribute Reasoning	Move the heaviest box which is marked with red label to storage room B on the second floor of the warehouse.	✓	Unambiguous quantitative comparison criterion (weight), clearly specified destination, and the marking of the box avoids confusion.
Attribute Reasoning	Move the large box to storage which is not clearly designated with a specific room number or location.	✗	“Large” is subjectively interpretable; multiple boxes may satisfy this criterion, and the storage location is not clear, leading to operational ambiguity.
Implicit Collab.	Move the 150kg industrial cabinet which is made of solid steel to warehouse located 500 meters away from the current location.	✓	Object weight (150kg) clearly exceeds single-agent capacity (50kg); collaboration requirement is inferable from physical constraints, and the material and location information are clear.
Implicit Collab.	Move the cabinet which is of unknown material and weight to warehouse without any additional information.	✗	Cabinet weight (30kg) is within single-agent capacity; collaboration necessity unclear and potentially unnecessary, and the lack of detailed information increases operational uncertainty.

Table 6: Representative examples of accepted (✓) and rejected (✗) scenarios with detailed rationale illustrating our quality threshold criteria.

Inter-Annotator Agreement. We assessed annotation reliability by measuring pairwise agreement using Cohen’s κ across all 1,804 scenarios evaluated in Stage 2. The results demonstrate strong overall agreement: $\kappa = 0.84$ (95% CI: [0.82, 0.86], computed via bootstrap resampling with 10,000 iterations). Agreement varied by criterion: Instruction Clarity achieved the highest agreement ($\kappa = 0.88$), followed by Physical Feasibility ($\kappa = 0.82$) and Solvability ($\kappa = 0.79$). The relatively lower agreement on Solvability reflects the inherent cognitive difficulty of determining solution existence through mental simulation alone, which validates our design decision to include simulation-based verification in Stage 3. Disagreements between annotators (204 cases, 11.3% of scenarios) were systematically adjudicated by a senior researcher through interactive simulation-based verification, with final decisions documented for quality tracking.

Acceptance Results. Human verification accepted 1,500 of 1,804 candidate scenarios (83.1%), demonstrating that automated validation alone is insufficient for ensuring benchmark quality—16.9% of automatically validated scenarios contained subtle issues detectable only through human expertise. Analysis of rejection rationales revealed the following distribution: ambiguous or underspecified instructions (37.2%), unsolvable tasks due to missing preconditions or unreachable goals (31.6%), physical inconsistencies between object properties and task requirements (22.4%), and other issues including edge cases and annotation uncertainties (8.8%).

Acceptance rates varied systematically by task complexity, as shown in Table 7. The observed inverse relationship between task complexity and acceptance rate provides indirect validation of our hierarchical task taxonomy: more complex tasks involving multiple reasoning dimensions naturally present greater opportunities for subtle generation errors. This pattern also confirms the importance of human verification, particularly for advanced task categories where automated checks are less effective.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Task Category	Acceptance Rate
Direct Command (L1)	91.2%
Tool Use (L2)	85.7%
Attribute Reasoning (L2)	82.3%
Compound Reasoning (L3)	78.9%
Explicit Collaboration (L1)	84.1%
Implicit Collaboration (L2)	79.6%
Compound Collaboration (L3)	76.4%

Table 7: Human verification acceptance rates by task category, showing systematic variation with task complexity.

6.3.3 STAGE 3: EXPERT TRAJECTORY VALIDATION

The final validation stage ensures that all accepted scenarios are not only well-formed but genuinely solvable through executable action sequences. Oracle agents with complete environmental knowledge generated expert demonstration trajectories for all 1,500 human-verified scenarios.

Trajectory Generation and Verification. Each generated trajectory underwent rigorous simulation replay to verify three critical properties: (1) precondition satisfaction at each action step, ensuring no invalid state transitions; (2) goal state achievement upon trajectory completion, confirming task success; and (3) physical constraint compliance throughout execution, validating that all intermediate states respect environmental limitations. Trajectories failing any verification criterion were regenerated using alternative planning strategies, with a maximum of three attempts per scenario. Scenarios with persistent trajectory failures were removed from the final dataset (18 scenarios, 1.2%), as such failures indicate potential issues not detected in earlier stages.

Oracle Quality Validation. To verify that oracle trajectories represent high-quality reference solutions rather than merely valid but suboptimal ones, we conducted an additional validation study on a randomly sampled subset of 100 scenarios stratified across task categories. Two human experts with extensive experience in embodied AI independently solved each scenario without access to oracle solutions, and we compared their solutions with the corresponding oracle trajectories using step count as the primary efficiency metric. Oracle solutions matched or outperformed human expert solutions in 94% of evaluated cases, with the remaining 6% representing alternative but equally valid strategies (typically involving different but equivalent tool choices or spatial paths). This validation confirms that oracle trajectories provide appropriate reference solutions for both evaluation and potential use in supervised learning approaches.

Final Dataset Statistics. The validated dataset comprises 1,500 scenarios instantiated into 16,592 task instances (multiple task variations per scenario). Expert trajectory statistics: mean length 9.1 steps (SD 4.2), median 9.0 steps, range [3, 25] steps. The high trajectory generation success rate (98.8% of human-verified scenarios, averaging 1.3 generation attempts per scenario) provides empirical validation that our human verification process effectively identifies genuinely solvable scenarios.

Ethics Statement. All human annotators provided written informed consent prior to participation and were compensated at standard institutional rates (\$20 USD per hour). The annotation task involved evaluation of synthetic embodied AI scenarios and did not involve deceptive practices, collection of personal information, or exposure to sensitive or harmful content. This study protocol was reviewed and approved by our institutional review board (IRB) prior to data collection.

Summary. Our three-stage validation pipeline systematically addresses the quality challenges inherent in LLM-based benchmark generation. Automated validation handles 87.6% of systematic errors efficiently, human verification with trained annotators (Cohen’s $\kappa = 0.84$) catches subtle semantic and physical issues, and simulation-based trajectory validation confirms genuine solvability. The resulting benchmark of 1,500 validated scenarios provides a reliable foundation for evaluating embodied reasoning capabilities.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Metric	Count
Total Scenarios	1,500
Total Task Files	1,481
Total Task Instances	16,592
Interactive Objects	64,057
Spatial Nodes (Rooms)	6,634
Average Objects per Scene	42.7
Average Rooms per Scene	4.4
Collaborative Agent Pairs	1,481

Table 8: Dataset Overview and Scale

Task Category	Count	Percentage
<i>Single-Agent Tasks (65%)</i>		
Direct Command	2,684	16.2%
Attribute Reasoning	2,669	16.1%
Tool Use	2,190	13.2%
Compound Reasoning	2,214	13.3%
<i>Multi-Agent Tasks (35%)</i>		
Explicit Collaboration	2,160	13.0%
Implicit Collaboration	2,582	15.6%
Compound Collaboration	2,093	12.6%
Total	16,592	100%

Table 9: Hierarchical Task Distribution

6.4 BENCHMARK STATISTICS AND COVERAGE

EAR-Bench encompasses 1,500 scenarios with 64,057 interactive objects, providing comprehensive coverage across diverse domains and task complexities. Tables 8 through 11 present detailed statistics demonstrating the scale and diversity of our benchmark.

Physical Property Modeling. The benchmark features exceptional attribute diversity with 6,381 distinct property types. Core physical properties are comprehensively modeled: weight (64,047 objects), material composition (35,411 objects), size dimensions (22,820 objects), color (28,034 objects), and dynamic states (17,547 objects). This rich attribute space enables sophisticated reasoning about physical constraints and object affordances.

Action Space and Tool Ecosystem. The framework supports 214 distinct action types, partitioned into basic actions (60%) available to all agents and tool-dependent actions (40%) requiring specific capabilities. Among the 64,057 objects, 15,134 are classified as tools (23.6%), with 13,482 objects possessing the `provides_abilities` attribute that enables dynamic capability extension. This design enables realistic modeling of how agents acquire new abilities through tool use.

Cross-Domain Coverage. The benchmark spans diverse application domains, with laboratory environments comprising 39.0% of scenarios, followed by office (18.8%), industrial (11.5%), and medical (6.2%) settings. This distribution reflects our emphasis on professional environments where embodied reasoning is particularly critical. Each domain presents unique challenges: laboratory settings require precise tool usage and material handling, office environments emphasize multi-agent coordination, and industrial scenarios demand reasoning about heavy equipment and safety constraints.

Quality Assurance and Expert Trajectories. All 16,592 task instances include expert demonstration trajectories averaging 8.7 steps, providing optimal solutions for comparison and learning. Each trajectory undergoes validation to ensure physical feasibility and task completion. The evaluation framework supports multi-level verification including spatial relationships (1,300 location

Category/Material	Count	Percentage
<i>Object Categories</i>		
Container	17,632	27.5%
Tool	15,134	23.6%
Appliance	8,963	14.0%
Furniture	6,234	9.7%
Consumable	4,890	7.6%
Others	11,204	17.6%
<i>Material Types (Top 10 of 1,123)</i>		
Plastic	13,767	21.5%
Metal	11,274	17.6%
Wood	8,263	12.9%
Glass	6,277	9.8%
Fabric	5,060	7.9%
Ceramic	3,843	6.0%
Silicon	1,794	2.8%
Aluminum	1,601	2.5%
Steel	1,153	1.8%
Others	11,025	17.2%

Table 10: Object Categories and Material Distribution

Domain/Room Type	Count	Percentage
<i>Application Domains</i>		
Laboratory	585	39.0%
Office	282	18.8%
Industrial	173	11.5%
Medical	93	6.2%
Household	93	6.2%
Educational	63	4.2%
Retail	48	3.2%
Service	30	2.0%
Entertainment	27	1.8%
Transportation	23	1.5%
Others	83	5.6%
<i>Room Types (Top 5)</i>		
Laboratory	1,876	28.3%
Storage	1,234	18.6%
Workspace	987	14.9%
Office	765	11.5%
Workshop	543	8.2%

Table 11: Domain and Spatial Distribution

checks), state transitions (open/closed, on/off states), and compound conditions for complex task assessment. This comprehensive validation ensures that all tasks are both challenging and solvable, maintaining benchmark integrity while achieving unprecedented scale.

6.5 ANALYSIS

Failure Mode Analysis. Systematic failure analysis reveals task-specific performance bottlenecks that vary distinctly across model scales. Tool Use failures are dominated by exploration deficits (31.2%), where models fail to locate required tools while maintaining spatial representations. Models below 7B parameters exhibit 2.7-fold higher failure rates (84.2% vs. 31.2%), confirming critical scale thresholds for embodied reasoning. Compound Reasoning failures stem primarily from planning degradation (28.7%), with models losing track of intermediate subgoals during execution.

Implicit Collaboration shows distinct timing failures (35.8%)—models either initiate collaboration prematurely or miss coordination opportunities. This failure mode exhibits no scale correlation,

1134 indicating that collaboration timing demands reasoning mechanisms absent from current architec-
 1135 tures. These failure patterns demonstrate that task categories stress fundamentally different cognitive
 1136 capabilities, necessitating targeted architectural solutions beyond universal parameter scaling.
 1137

1138 6.6 HYPERPARAMETERS

1139
 1140 **Supervised Fine-Tuning.** We performed full-parameter supervised fine-tuning on the
 1141 Qwen2.5-3B-Instruct model to adapt it to our dataset. The training was conducted on
 1142 4x NVIDIA A100 GPUs. The effective batch size was 64, achieved through a per-device batch size
 1143 of 1 and 16 gradient accumulation steps across 4 devices. Key hyperparameters for the SFT stage
 1144 are summarized in Table 12.

Hyperparameter	Value
Base Model	Qwen2.5-3B-Instruct
Fine-tuning Method	Full-parameter
Effective Batch Size	64
Learning Rate	1.0e-5
LR Scheduler	Cosine Decay
Warmup Ratio	0.1
Training Epochs	3
Max Sequence Length	15,360
Precision	BF16

1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155 Table 12: Hyperparameters for Supervised Fine-Tuning.

1156
 1157
 1158 **Model Inference.** To ensure a fair and consistent comparison, all models were evaluated using the
 1159 same set of inference parameters. We utilized the vLLM engine for efficient serving, with a tensor
 1160 parallel size of 4. The decoding strategy was configured to balance response quality and exploration
 1161 in complex reasoning tasks. The inference settings are detailed in Table 13.

Hyperparameter	Value
Inference Engine	vLLM
Tensor Parallel Size	4
Decoding Strategy	Nucleus Sampling
Temperature	0.3
Top-p	1.0 (Default)
Max Generation Tokens	4096
Max Model Length	15,360

1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172 Table 13: Hyperparameters for Model Inference.

1173 6.7 DISCUSSION

1174
 1175 **Embodied vs. Abstract Reasoning.** Our results demonstrate that embodied reasoning requires
 1176 distinct computational mechanisms from abstract reasoning in current language models. The
 1177 persistent performance gaps across reasoning-specialized architectures indicate that chain-of-
 1178 thought approaches cannot bridge the representational divide between symbolic manipulation and
 1179 physical constraint processing. Current transformer architectures lack the specialized components
 1180 necessary for grounding abstract representations in continuous physical properties.
 1181

1182 **Architectural Limitations.** The constraint selection failures reveal that current attention mech-
 1183 anisms cannot dynamically filter task-relevant physical constraints from environmental noise.
 1184 Unlike abstract reasoning tasks where all provided information typically bears relevance, embodied
 1185 scenarios require selective attention over spatially and temporally distributed constraint sets. The
 1186 discrete scaling transitions at 7B parameters indicate that embodied reasoning demands sufficient
 1187 working memory capacity to simultaneously track environmental states, capability constraints, and
 coordination requirements—a computational bottleneck absent in pure language tasks.

Limitations and Future Work. Our text-based framework abstracts away continuous control, sensorimotor feedback, and real-time constraints present in physical embodied systems. While this abstraction enables systematic evaluation, it may not capture all aspects of embodied intelligence. The identified architectural requirements require validation in continuous control settings. Future work should investigate how these components integrate with sensorimotor processing and examine whether the observed computational bottlenecks persist in physically grounded systems. Additionally, exploring hybrid symbolic-neural architectures that can explicitly reason about physical laws while maintaining learned flexibility represents a promising direction (Rabinowitz et al., 2018).

6.8 AGENT PROMPT CONFIGURATIONS

This section details the system and user prompts used for different experimental configurations: single-agent and multi-agent scenarios.

Single-Agent Configuration. This configuration tests individual agent reasoning capabilities through structured prompts.

System prompt for single-agent

1. PRIMARY OBJECTIVE

Your goal is to successfully complete the given task by systematically exploring the environment and interacting with objects. Success requires persistence, thorough exploration, and precise execution of interaction sequences.

2. MANDATORY OUTPUT REQUIREMENTS

You must follow these absolute rules in every single response:

Strict Format Compliance: Your entire output must be in the exact format `Thought: <reasoning>\nAgent.L.Action: <command>`. Do not include any other text, explanations, or formatting.

Command Validation: The command you choose must be exactly as listed in the Available Actions provided in the user prompt. Do not invent or modify commands.

Progress Verification: After completing any part of the task, always re-read the task description in your next thought to verify if additional objectives remain incomplete.

Completion Protocol: Use the DONE action if and only if you have verified that all objectives in the task description have been successfully completed.

3. OPERATIONAL FRAMEWORK

Exploration Strategy: First use EXPLORE to thoroughly examine your current room. If the target isn't found, systematically GOTO and EXPLORE each unexplored room until completing the task.

Interaction Sequence Protocol: Always approach an object using GOTO before attempting any interaction with it. Always open containers using OPEN before taking items from or placing items into them. This sequence prevents interaction failures and ensures reliable task execution.

4. CRITICAL FAILURE PATTERNS TO AVOID

Premature Task Abandonment: Do not conclude failure without exploring every available room and container. Persistence is essential for task completion.

Object Name Confusion: Different names represent different objects. Verify exact matches between task requirements and available objects before taking action.

Distance Interaction Violations: Do not attempt to interact with objects that are not in immediate proximity. Always use GOTO to approach objects first.

Container Access Oversight: Do not forget to open containers before attempting to access their contents. This is a common cause of interaction failures.

5. ERROR RECOVERY PROTOCOL

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254

If your chosen action results in an error, acknowledge the error in your next thought and immediately re-evaluate your strategy based on available information. Do not repeat failed actions unless the environmental situation has changed.

6. REQUIRED OUTPUT FORMAT

Your response must contain exactly two lines in this format:

Thought: [Your reasoning for taking this action]

Agent_1Action: [Command from the available action list]

Example Response:

Thought: I am in the main work area and need to find the target objects. I have not explored the living room yet, so I should go there next.

Agent_1Action: GOTO living_room.1

1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279

User prompt for single-agent

You are an intelligent agent tasked with completing the given objective by strictly following the operational framework established in your system instructions. Analyze the information provided below and determine the single best next action that will advance progress toward task completion.

Current Environment

{environment_description}

Task Objective

{task_description}

Available Actions

{available_actions_list}

Recent Action History

{history_summary}

Execution Guidelines

Respond with exactly one thought and one action. Your thought should demonstrate systematic reasoning that considers the current situation, task requirements, and appropriate next steps. Your action must be selected from the available actions list and should represent the most logical progression toward completing the task objective.

Remember that systematic exploration, proper interaction sequences, and persistent problem-solving are essential for successful task completion. The available action descriptions will guide you on exactly how to execute each command effectively.

1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Multi-Agent Configuration. This configuration provides prompts for coordinated reasoning between two agents.

System prompt for multi-agent

You are a central coordination controller managing two intelligent agents working collaboratively to complete complex tasks. Your responsibility is to analyze the current situation, decompose objectives into executable subtasks, and assign optimal actions to both agents while ensuring efficient coordination and conflict avoidance.

Core Coordination Principles

Strategic Assignment Protocol: Assign actions based on each agent's current position, capabilities, and the optimal path toward task completion. Prioritize complementary actions that maximize overall efficiency.

Conflict Prevention Framework: Ensure that assigned actions do not create spatial conflicts, resource competition, or contradictory objectives between the two agents.

1296
1297 **Exploration Optimization:** When agents have completed their immediate
1298 objectives, prioritize exploration of unknown areas to gather
1299 additional environmental information and identify new opportunities
1300 for task advancement.

1301 **Cooperation Command Protocol**
1302 For collaborative tasks requiring joint action, implement the
1303 following cooperation strategy:

1304 **Pre-Cooperation Positioning:** Before initiating any CORP_ command
1305 sequence, ensure that both participating agents have successfully
1306 executed GOTO commands to reach the target object or designated
1307 cooperation zone.

1308 **Cooperative Transport Sequence:** For tasks involving collaborative
1309 object movement, execute the following mandatory sequence without
1310 interruption:
1311 1. CORP_GRAB - Both agents grab/pick up the target object
1312 2. CORP_GOTO - Coordinated movement to the destination location
1313 3. CORP_PLACE - Synchronized placement of the object at the target
1314 location

1315 **Critical CORP_PLACE Requirement:** After executing CORP_GOTO, you MUST
1316 execute CORP_PLACE to actually place the object at the destination.
1317 The object is not considered "moved" until CORP_PLACE is completed.

1318 **Sequence Integrity Requirement:** The cooperative transport sequence
1319 must be executed continuously without interspersing other commands.
1320 Any interruption requires restarting the entire cooperation sequence.
1321 NEVER output DONE after CORP_GOTO - always complete with CORP_PLACE
1322 first.

1323 **Cooperation Readiness Verification:** Verify that both agents are
1324 properly positioned and available for cooperation before initiating
1325 any CORP_ command. This prevents coordination failures and ensures
1326 successful collaborative execution.

1327 **Task Completion Management**

1328 **Individual Agent Completion:** When an agent has no additional
1329 meaningful tasks to perform, assign the DONE command to that specific
1330 agent while continuing to provide actionable commands to the other
1331 agent.

1332 **Final Task Termination:** The overall task concludes only when both
1333 agents simultaneously receive DONE commands, indicating that all
1334 objectives have been completed and no further actions are required.

1335 **Continuation Protocol:** When one agent completes all its tasks,
1336 consistently assign DONE to that agent in all subsequent action
1337 assignments while continuing to provide meaningful actions to the
1338 remaining active agent until it also completes its objectives.

1339 **Mandatory Output Format**
1340 Your response must adhere to the following strict format without any
1341 additional content or explanations:
1342 Thought: [Comprehensive analysis of current situation, task
1343 requirements, and strategic reasoning for action assignments]
1344 Agent_1.Action: [Specific command for agent_1 from available action
1345 set]
1346 Agent_2.Action: [Specific command for agent_2 from available action
1347 set]
1348 Example:
1349 Thought: Agent 1 is in the main work area and needs to explore,
while agent 2 should go to the living room to find target items.
Agent_1.Action: EXPLORE
Agent_2.Action: GOTO living_room_1

1350 **Strategic Planning Guidelines**

1351 **Situational Assessment:** Evaluate each agent's current location,
1352 recent actions, and immediate objectives to determine the most
1353 effective next steps.

1354 **Resource Allocation:** Consider the spatial distribution of tasks and
1355 assign agents to different areas when possible to maximize coverage
1356 and minimize redundancy.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Progress Monitoring: Track completion status of subtasks and adjust assignments based on evolving priorities and environmental discoveries.

Efficiency Optimization: Balance individual agent productivity with collaborative opportunities to achieve optimal overall task completion time.

User prompt for multi-agent

Analyze the provided information and generate coordinated action assignments for both agents:

Current Environment State

{environment_description}

Task Objectives

{task_description}

Available Commands

{available_actions_list}

Agent Status and History

{history_summary}

Coordination Requirements

Generate action assignments that advance task completion while maintaining coordination efficiency. Ensure that cooperative tasks follow the established CORP command protocols and that individual assignments complement overall strategic objectives.