

# EmbodiedEval: Evaluate Multimodal LLMs as Embodied Agents

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

001 *Multimodal Large Language Models (MLLMs) have shown*  
 002 *significant advancements, providing a promising future for*  
 003 *embodied agents. Existing benchmarks for evaluating*  
 004 *MLLMs primarily utilize static images or videos, limiting*  
 005 *assessments to non-interactive scenarios. Meanwhile, ex-*  
 006 *isting embodied AI benchmarks are task-specific and not*  
 007 *diverse enough, which do not adequately evaluate the em-*  
 008 *bodied capabilities of MLLMs. To address this, we propose*  
 009 *EMBODIEVAL, a comprehensive and interactive evalu-*  
 010 *ation benchmark for MLLMs with embodied tasks. EM-*  
 011 *BODIEVAL features 328 distinct tasks within 125 varied*  
 012 *3D scenes, each of which is rigorously selected and*  
 013 *annotated. It covers a broad spectrum of existing em-*  
 014 *bodied AI tasks with significantly enhanced diversity, all*  
 015 *within a unified simulation and evaluation framework tai-*  
 016 *lored for MLLMs. The tasks are organized into five cat-*  
 017 *egories: navigation, object interaction, social interaction,*  
 018 *attribute question answering, and spatial question answer-*  
 019 *ing to assess different capabilities of the agents. We eval-*  
 020 *uated the state-of-the-art MLLMs on EMBODIEVAL and*  
 021 *found that they have a significant shortfall compared to hu-*  
 022 *man level on embodied tasks. Our analysis demonstrates*  
 023 *the limitations of existing MLLMs in embodied capabili-*  
 024 *ties, providing insights for their future development. We*  
 025 *open-source all evaluation data and simulation framework*  
 026 *at <https://github.com/thunlp/EmbodiedEval>.*

## 027 1. Introduction

028 In recent years, Multimodal Large Language Models  
 029 (MLLMs) [53, 67, 82] have demonstrated strong un-  
 030 derstanding and reasoning capabilities on visual and language  
 031 tasks. With the rapid development of MLLMs, a rich set  
 032 of benchmarks for image understanding [25, 45, 56, 91]  
 033 and video analysis [27, 63] has emerged. However, these  
 034 benchmarks are non-interactive and insufficient in evalu-  
 035 ating MLLMs' ability to handle tasks in real-world scen-  
 036 arios. Beyond these basic tasks which focus on non-  
 037 interactive visual scenes, researchers are actively trying to



Figure 1. Examples of the five task categories in EMBODIEVAL. On the left are the task text and part of the action space. On the right are observations from specific steps, along with the actions taken in the expert demonstration at those moments.

expand MLLMs as embodied agents in interactive environ-  
 ments, which require the model to interpret multimodal in-  
 puts into actions [1, 10, 22, 62, 66, 84, 99]. To accomplish  
 this, MLLMs are expected to integrate a multitude of ca-  
 pabilities that enable them to interact effectively with the  
 environment, including ego-centric perception [15], visual  
 grounding [4, 96], spatial reasoning [8], episodic memory  
 [63], among others.

However, the comprehensive evaluation of MLLMs in  
 embodied tasks remains largely unexplored. Existing  
 benchmarks for embodied AI greatly lack diversity in tasks

049 and scenes, which limits their evaluation to specific aspects.  
 050 EQA [18] dataset is restricted to only 9 template questions  
 051 such as queries about object color and existence, while AL-  
 052 FRED [76] contains only 7 template tasks such as “pick and  
 053 then place” within 4 room types. Moreover, many of emb-  
 054 bodied AI benchmarks require specific forms of input and  
 055 output, making the evaluation of mainstream MLLMs inef-  
 056 ficient or even infeasible. For example, R2R [4] uses 3D  
 057 points as both input and output, REVERIE [71] requires  
 058 predicting the bounding box of a target object, while AL-  
 059 FRED requires an object segmentation mask for interaction.

060 To address this gap, we present the first compre-  
 061 hensive evaluation benchmark for evaluating MLLMs’ embod-  
 062 ied capabilities in interactive environments, distinguishing  
 063 it from the existing three types of evaluation benchmarks:  
 064 static MLLMs benchmarks, which are non-interactive; emb-  
 065 bodied AI datasets, which are not comprehensive or diverse  
 066 enough; and LLM agent benchmarks [54, 77], which heav-  
 067 ily rely on textual environment states and have overly ab-  
 068 stract actions (e.g., go to bathroom), thus downplaying or  
 069 overlooking the critical embodied agent capabilities such as  
 070 visual grounding and spatial reasoning. The key features of  
 071 EMBODIEDEVAL are as follows:

072 **Diverse Interactions.** EMBODIEDEVAL provides a sim-  
 073 ulation framework that supports a variety of interactions,  
 074 such as locomotion, question-answering, and interactions  
 075 with objects and humans in realistic 3D scenes. Agents in-  
 076 teract with the environment to gather new information or  
 077 change its state to complete the task.

078 **Diverse Tasks.** In contrast to previous work that used  
 079 task templates or tasks with minimal variation, our tasks  
 080 have been systematically generated and carefully selected  
 081 for high quality and diversity. EMBODIEDEVAL featuring  
 082 a large number of novel tasks that involve a wider range  
 083 of abilities to assess the comprehensive capabilities of the  
 084 model as shown in Figure 1, which are categorized into five  
 085 major categories.

086 **Diverse Scenes.** Unlike previous work, which was lim-  
 087 ited to household scenes, our scenes feature significant  
 088 diversity in terms of objects and spaces, covering small  
 089 rooms, large residences, and public spaces such as gym,  
 090 store, office, among others. This approach can reduce the  
 091 impact of scene types on the evaluation of the model’s gen-  
 092 eralization, making the assessment more comprehensive.

## 093 2. Related Works

094 **Multimodal Large Language Models.** By connecting  
 095 vision modules with LLMs, LLaVA [53] pioneers research  
 096 in MLLMs through visual instruction tuning, and obtains  
 097 impressive multimodal chat capabilities. Many work fur-  
 098 ther improves the MLLMs from various aspects, includ-  
 099 ing detailed captioning [9], trustworthy response [89, 90],  
 100 multilingual multimodal capabilities [32, 80] and visual

Benchmark	Scene.	Task.	Disc.	Ego.	Nav.	Obj.	So.	Ans.
Video-MME [26]	-	✓	✓	✗	✗	✗	✗	✓
EgoPlan etc. [11, 14]	-	✓	✓	✓	✗	✗	✗	✓
OpenEQA [63]	-	✓	✓	✓	✗	✗	✗	✓
EQA etc. [18, 81, 88]	✗	✗	✓	✓	✓	✗	✗	✓
ALFRED [76]	✗	✗	✗	✓	✓	✓	✗	✗
BEHAVIOR[79]	✗	✓	✗	✓	✓	✓	✗	✗
EQA-MX [35]	✗	✗	✓	✓	✗	✗	✓	✓
<b>EMBODIEDEVAL</b>	✓	✓	✓	✓	✓	✓	✓	✓

Table 1. Comparison of EMBODIEDEVAL with previous bench-  
 marks. The abbreviations in the table headers, from left to right,  
 represent: **Scene** diversity (beyond household scenes), **Task**  
 diversity (beyond task templates), **Discrete** action space (for MLLMs  
 evaluation), **Egocentric** vision, **Navigation** involved, **Object** in-  
 teraction involved, **Social** interaction involved, and **Answering** ques-  
 tions involved.

grounding [70, 87]. Beyond single-image understanding, 101  
 some work explores more complicated tasks. For exam- 102  
 ple, KOSMOS-1 [33] and VILA [49] focus on image-text 103  
 interleaved understanding, while Video-LLaVA [48] and 104  
 VideoChat [47] focus on video understanding. 105

**Evaluation for MLLMs.** Mainstream benchmarks for 106  
 MLLMs mainly focus on perception and cognitive evalu- 107  
 ation, such as MME [25], MMB [56] and MMMU [91]. As 108  
 existing MLLMs excel in these benchmarks, some bench- 109  
 marks propose more challenging tasks, such as mathemat- 110  
 ical reasoning [31, 60, 93], OCR capability [57, 64, 78] 111  
 and scientific knowledge [59, 91]. However, these bench- 112  
 marks lack evaluation on egocentric vision, which is essen- 113  
 tial for broader applications of MLLMs. To address this 114  
 gap, EgoVQA [24], EgoPlan-Bench [11], EgoThink [14], 115  
 and OpenEQA [63] propose to evaluate the reasoning and 116  
 planning capabilities of MLLMs given the first-person per- 117  
 spective images or videos. However, these benchmarks still 118  
 use static question-answering pairs without interacting with 119  
 environments. 120

**Benchmarks for Embodied Agents.** The datasets for 121  
 embodied agents cover areas such as navigation, interac- 122  
 tion, and question answering: (1) **Navigation** The R2R [4] 123  
 dataset was the first to evaluate an agent’s navigation ability 124  
 under natural language instructions, followed by R4R [36] 125  
 and RxR [44], which improve the fine-grained evaluation 126  
 of the navigation process. In navigation tasks, object nav- 127  
 igation is a crucial task because it serves as a prerequi- 128  
 site step for an embodied agent to interact with any object. 129  
 There are many object navigation datasets with different fo- 130  
 cus and features, including SOON [98], REVERIE [72], 131  
 DOZE [61] and GOAT-Bench [40]. (2) **Interaction** AL- 132  
 FRED [76] is the most representative interaction dataset, 133  
 requiring the agent to follow instructions to complete tasks 134  
 involving interactions such as picking up and placing ob- 135

136 jects. Additionally, there are datasets focusing on mov- 187  
 137 ing objects [28, 86], rearrangement [6, 83], tidying up a 188  
 138 room [38], and household activities [42, 46, 65, 79]. (3) 189  
 139 **Question Answering** EQA [18] first proposes the navigat- 190  
 140 e-then-answer mechanism. Subsequently, more diverse EQA 191  
 141 datasets have emerged, including those with questions in- 192  
 142 volving multiple objects [88], requiring knowledge integra- 193  
 143 tion [81], set in realistic scenes [73], and handling situa- 194  
 144 tional queries [21]. Additionally, IQA [30] requires interac- 195  
 145 tion with the environment to gather more observations 196  
 146 in order to answer questions. EQA-MX [34] requires un- 197  
 147 derstanding non-verbal human expressions, such as body 198  
 148 movements. However, existing benchmarks are limited in 199  
 149 task variety, lacking comprehensive assessments of naviga- 200  
 150 tion, object interaction, and question-answering. They rely 201  
 151 on highly repetitive task templates, failing to adequately 202  
 152 capture the wide spectrum of embodied capabilities. Addi- 203  
 153 tionally, the task-specific observation spaces and continuous 204  
 154 action spaces in many benchmarks are inadequate for effec- 205  
 155 tively evaluating MLLMs. We summarize the comparison 206  
 156 between EMBODIEVAL and other representative bench- 207  
 157 marks in Table 1. 208

### 158 3. EmbodiedEval

159 EMBODIEVAL consists of rich scenes and tasks to com- 210  
 160 prehensively evaluate the capabilities of embodied agents. 211  
 161 We have implemented a unified simulation and evaluation 212  
 162 framework and conducted meticulous data annotation to 213  
 163 create the final dataset. First, we introduce the task cat- 214  
 164 egories involved in EMBODIEVAL in Section 3.1, fol- 215  
 165 lowed by a detailed explanation of the simulation and eval- 216  
 166 uation framework in Section 3.2. After that, we describe 217  
 167 the data collection and annotation process for this frame- 218  
 168 work in Section 3.3. Finally, we show the data statistics of 219  
 169 EMBODIEVAL in Section 3.4. 220

#### 170 3.1. Task Categories

171 To comprehensively assess the capabilities of MLLMs as 221  
 172 embodied agents, we synthesized and expanded upon ex- 222  
 173 isting embodied tasks while incorporating novel additions, 223  
 174 broadly categorizing the tasks into five categories. Each cat- 224  
 175 egories encompasses a diverse set of tasks and involves var- 225  
 176 ious capabilities. The task categories and sample tasks are 226  
 177 shown in Figure 1.

178 **Navigation.** The navigation task involves coarse-grained 221  
 179 and fine-grained natural language instructions, requiring the 222  
 180 agent to navigate from its initial position to a target location 223  
 181 and find a specific object if the task demands it. 224

182 **Object Interaction.** In object interaction tasks, agents 225  
 183 must change the state of the environment through direct in- 226  
 184 teraction with objects, such as moving objects, opening/- 227  
 185 closing doors and drawers, and operating electrical devices. 228  
 186 These tasks possibly involve multiple objects and require

multi-step actions, such as using a tool to operate another 187  
 object and rearranging items to meet certain requirements. 188

**Social Interaction.** Social interaction tasks encom- 189  
 pass human-agent interactions, including item delivery, 190  
 perspective-taking capabilities, human feedback interpreta- 191  
 tion, and non-verbal communication comprehension such 192  
 as gestural cues. 193

**Attribute Question Answering (AttrQA).** AttrQA 194  
 tasks necessitate exploring the environment to address ques- 195  
 tions regarding object and scene attributes. While incorpo- 196  
 rating traditional EQA tasks, AttrQA significantly broad- 197  
 ens the scope of inquiry. It encompasses comprehensive 198  
 attribute queries of objects and scenes, including but not 199  
 limited to category, shape, material, color, function, state, 200  
 location, existence, quantity, comparative analysis and rea- 201  
 soning across multiple attributes. 202

**Spatial Question Answering (SpatialQA).** Spatial un- 203  
 derstanding [8, 13, 23, 37, 51, 75] is a fundamental capabil- 204  
 ity for embodied agents. SpatialQA requires agents to an- 205  
 swer spatial-related questions through actions and observa- 206  
 tions, such as queries about size, height, position, distance, 207  
 area, path, layout, spatial relationships, and more. 208

#### 3.2. Evaluation Framework 209

We implemented a unified simulation framework based on 210  
 LEGENT [17] platform. The framework aims to holistically 211  
 evaluate embodied abilities through diverse tasks, rather 212  
 than focusing on specific tasks with particular input-output 213  
 requirements. The observation space consists of task de- 214  
 scriptions and egocentric vision. The action space consists 215  
 of movement space, interaction space and answering space, 216  
 which varies in each task instance. The action space is rea- 217  
 sonably discretized and semantically understandable, while 218  
 preserving the core evaluation objectives. 219

##### 3.2.1 Action Space 220

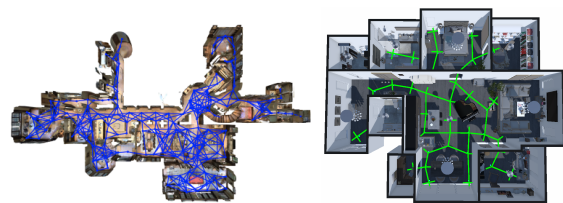


Figure 2. A comparison of navigation graphs between R2R [4] dataset (left) and EMBODIEVAL (right).

**Movement Space** We use navigation graph as the move- 221  
 ment space where the agent can rotate its view at a point 222  
 or move between adjacent navigation points. Compared to 223  
 continuous movement, it discretizes motion without impos- 224  
 ing great restrictions on the high-level tasks in practice [4]. 225  
 Different from grid-world movement, this approach is more 226

227	natural and adaptable to all kinds of scenes. Through sam-	280
228	pling algorithms and manual adjustment, we constructed	281
229	navigation graphs for each scene. To ensure realism, the	282
230	navigation points are always walkable locations with no ob-	
231	stacles among them. Due to the greater diversity of our	
232	scenes and tasks compared to previous work, the density of	
233	navigation points varies based on the size of the scene and	
234	the task, ensuring that the number of steps required for tasks	
235	remains reasonable. For example, in complex interaction	
236	tasks within large scenes, the navigation points are more	
237	sparse and critical. In contrast to previous datasets, our nav-	
238	igation points are better organized as shown in Figure 2, and	
239	the connections between these points indicate clear seman-	
240	tics. MLLMs are not required to make choices from a set	
241	of 3D positions, but only need to make directional decisions	
242	among navigation points. Specifically, the action space con-	
243	sists of three types of actions: <i>move forward</i> (moving to	
244	the facing navigation point), <i>turn left/right</i> (rotating to face	
245	a new navigation point), and <i>look up/down</i> (adjusting the	
246	vertical view).	
247	<b>Interaction Space</b> We follows the discrete interaction	
248	space of previous embodied AI tasks that involves object	
249	interaction such as IQA [30], CHAI [65], RoomR [83] and	
250	OVMM[86] rather than continuous space [76, 79]. This	
251	choice is based on two main considerations: (1) In contin-	
252	uous spaces, interactions are tightly related to specific	
253	methods and types of embodiment, which contradicts the	
254	goal of generality in evaluations and goes beyond the core	
255	issues of our research. (2) Due to the high complexity of	
256	continuous space, MLLMs cannot output reasonable values	
257	without being trained on specialized numerical trajectory,	
258	leading to infeasible evaluations. In EMBODIEVAL, we	
259	use an open vocabulary for the actions and objects in inter-	
260	actions to make them as rich as possible. Each interaction	
261	action has a brief action text, operable objects, and condi-	
262	tions for successful interaction. For example, the “pick up”	
263	action requires the target object to be within sight and very	
264	close, the “wash” action requires the agent and the target	
265	object to be next to a sink, and the “hand over” action re-	
266	quires the agent to hold an object and be next to a person.	
267	In a given test case, several interaction actions will be in-	
268	volved, including those necessary to complete the task and	
269	other distracting actions.	
270	<b>Answering Space</b> Unlike EQA and IQA, which require	
271	outputs of specific categories or use a very limited vocabu-	
272	lary, our answer space is very open. Responses are written	
273	by annotators and rigorously verified, encompassing a vari-	
274	ety of possible replies to the question. For QA tasks, before	
275	exceeding the maximum allowed steps, there is no speci-	
276	fied step at which the agent must respond; it can continue	
277	to explore until it believes it can select an answer. Once an	
278	answer is chosen, the task is immediately judged as correct	
279	or incorrect.	
	In this way, the action space is discretized and each ac-	280
	tion has a semantic text. EMBODIEVAL directly input	281
	the list of actions into the MLLM to decide the next action.	282
	<b>3.2.2 Success Criteria</b>	283
	We automatically and accurately evaluate task completion	284
	through predicate functions. Each predicate maps the state	285
	of the simulation environment to a boolean value indicat-	286
	ing success. For example, the <i>agent_at</i> predicate requires	287
	a designated navigation point as a parameter and returns	288
	true when the agent reaches this location at the end of the	289
	episode. Beyond evaluating only the final state, EMBOD-	290
	IEVAL also includes predicates that assess the entire proc-	291
	ess, similar to R4R [36]. For example, the <i>agent_pass</i>	292
	predicate becomes true once the agent passes a specified	293
	navigation point.	294
	A task is considered successful when all predicates eval-	295
	uate to true at the end. Consider the task “Please go to	296
	the kitchen, then come back and tell me if there are any ex-	297
	tra cups”. This task involves three predicates: <i>agent_pass</i> ,	298
	<i>agent_at</i> , and <i>choose</i> . These predicates verify that the agent	299
	passes through the kitchen doorway, returns to the initial po-	300
	sition in front of the person, and selects the correct answer,	301
	respectively.	302
	<b>3.2.3 Evaluation Process</b>	303
	The process of an evaluation episode is as follows: (1) The	304
	simulator initializes the 3D scene and navigation graph. The	305
	agent is positioned at a designated starting point, and the	306
	initial first-person view image is saved into the observation	307
	history. (2) In each subsequent step, the agent chooses an	308
	action from a given list of options, including movement, in-	309
	teraction and answering, based on the observation history.	310
	The environment executes the action, changes the state ac-	311
	cordingly, returns new observations, along with feedback	312
	indicating whether the action was successful. The observa-	313
	tion, action, and feedback are then appended to the obser-	314
	vation history.	315
	This process continues until all success criteria are met	316
	resulting in task success, or the task fails due to an incor-	317
	rect answer or exceeding the maximum allowed steps. For	318
	evaluating MLLMs designed for video understanding, im-	319
	age sequences are no longer used in the observation history.	320
	Instead, at each step, the simulator will output an egocen-	321
	tric video from the start to the current moment and use it as	322
	input for the MLLM.	323
	<b>3.3. Dataset Construction</b>	324
	The dataset construction process of EMBODIEVAL con-	325
	sists of three parts: scene collection, task collection, and	326
	task annotation. Each sample in the dataset requires sub-	327

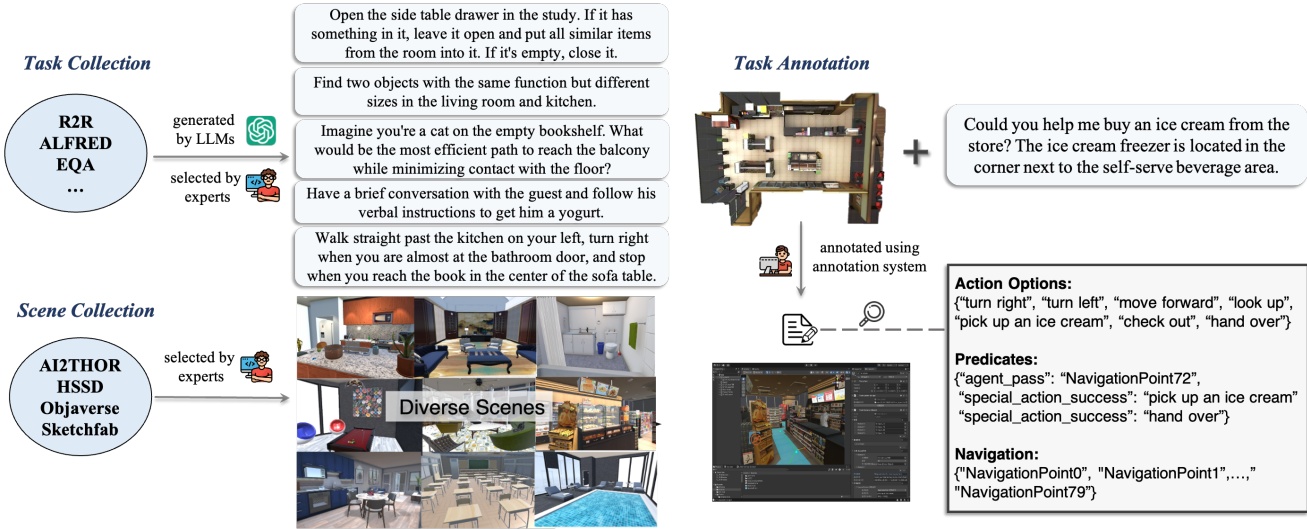


Figure 3. The dataset construction pipeline of EMBODIEVAL.

328 stantial effort and undergoes rigorous annotation. Figure 3  
329 illustrates our dataset construction pipeline.

### 330 3.3.1 Scene Collection

331 The diversity of scenes allows for a more accurate assess-  
332 ment of the agent’s generalization capabilities. Our scenes  
333 are more varied in source and type than previous bench-  
334 marks, curated from four different sources: (1) **Objaverse**  
335 **Synthetic**. To expand our object variety, we utilize the Ob-  
336 javerse [20] dataset as our object database, which contains  
337 a vast collection of objects. Using procedural generation  
338 methods [19], we created numerous scenes from these ob-  
339 jects. From these procedurally generated scenes, we man-  
340 ually selected and refined dozens of high-quality environ-  
341 ments. (2) **AI2THOR**. We extracted and structured scenes  
342 from AI2THOR [43], focusing on indoor rooms with highly  
343 interactive objects. (3) **HSSD**. We included scenes from the  
344 Habitat Synthetic Scenes Dataset (HSSD) [39], which pro-  
345 vides high-quality, realistic and complex scenes featuring  
346 extensive and diverse navigable spaces, including houses,  
347 villas, yards, and more. (4) **Sketchfab**. We enhanced our  
348 scene collection with selected free 3D scenes from Sketch-  
349 fab<sup>1</sup>. These scenes are highly realistic and diverse, ranging  
350 from classrooms and supermarkets to offices and exhibi-  
351 tions, enriching the scope of our evaluation. These collected  
352 scenes serve as the basis for subsequent task annotation.

### 353 3.3.2 Task Collection

354 For task collection, we first gather seed tasks from over  
355 30 existing datasets and benchmarks for embodied agents.

<sup>1</sup><https://sketchfab.com>

Using these tasks as seeds, we prompted several advanced  
large language models to generate diverse task examples.  
To enhance task complexity, we designed some tasks to in-  
corporate various capabilities, including complex ground-  
ing, episodic memory, spatial reasoning, quantitative reason-  
ing, common sense reasoning, and planning, which result-  
ed in many novel tasks. From this extensive task pool,  
experts selected over 300 distinct candidate tasks. Dur-  
ing task annotation, each task text could only be selected  
from this candidate set and used once. We chose this ap-  
proach rather than allowing annotators to write tasks for  
given scenes, as it ensures task diversity, prevents repeti-  
tion, and reduces dependency on individual annotators’ creativity  
or preferences.

### 370 3.3.3 Task Annotation

**Annotation Content** Task sample annotation begins with  
an annotator selecting a task from the candidate pool and  
matching it to an appropriate scene. The annotator config-  
ures the action space detailed in Section 3.2.1, which in-  
cludes movement space (via navigation graph adjustment),  
interaction space (via action options), and answering space  
(via answer options). Finally, the annotator defines success  
criteria as outlined in Section 3.2.2 through predicate func-  
tion instantiation.

To maintain the validity and quality of the dataset, all  
task annotations must satisfy the following criteria. (1) All  
tasks must be unambiguous within the given scene. (2)  
Question-answering tasks must require scene observation,  
with each task providing eight answer options that vary  
in difficulty and include misleading options to reduce the  
chance of guessing the correct answer. (3) Once a task



Model	Attr. QA	Spatial QA	Navigation			Object Interaction			Social Interaction			Overall	
	Succ.	Succ.	Succ.	GcS	SPL	Succ.	GcS	SPL	Succ.	GcS	SPL	Succ.	GcS
Random	11.58	7.69	3.45	8.76	3.45	0.00	6.18	0.00	2.94	8.33	2.94	5.49	8.66
Human	98.95	92.31	96.55	97.84	82.28	97.75	99.44	90.73	100.00	100.00	89.96	97.26	97.94
<i>Closed-Source Multi-Image MLLMs</i>													
GPT-4o [68]	35.79	<b>32.69</b>	<b>31.03</b>	<b>42.53</b>	<b>22.23</b>	<b>10.11</b>	24.25	<b>5.94</b>	<b>11.76</b>	<b>26.72</b>	6.74	<b>25.00</b>	<b>32.42</b>
GPT-4o-Mini [68]	31.58	15.38	27.59	39.51	15.34	2.25	17.42	1.50	5.88	22.06	2.98	17.68	25.58
Gemini-Pro [29]	27.37	9.62	17.24	25.86	9.78	4.49	12.36	3.00	5.88	18.14	3.44	14.33	19.26
Gemini-Flash [29]	26.32	13.46	5.17	17.10	3.51	2.25	7.58	0.96	2.94	12.50	1.47	11.59	16.13
Qwen-VL-Max [5]	<b>37.89</b>	17.31	24.14	30.03	16.87	7.87	<b>24.91</b>	5.62	8.82	22.06	<b>6.86</b>	21.04	28.07
Qwen-VL-Plus [5]	10.53	11.54	3.45	10.49	3.45	0.00	2.43	0.00	2.94	8.82	1.68	5.79	8.31
<i>Open-Source Multi-Image MLLMs</i>													
InternVL2-40B [69]	14.74	5.77	6.90	12.93	3.06	0.00	7.68	0.00	<b>5.88</b>	<b>19.12</b>	2.16	7.01	11.54
InternVL2-8B [69]	13.68	13.46	8.62	18.25	4.04	0.00	7.43	0.00	<b>5.88</b>	18.63	<b>2.45</b>	8.23	13.27
InternVL2-Llama3-76B [69]	21.05	13.46	3.45	9.48	2.18	0.00	9.08	0.00	2.94	13.73	1.14	9.15	13.79
LLaVA-NEXT-72B [12]	23.16	5.77	<b>12.07</b>	22.99	<b>7.83</b>	<b>3.37</b>	<b>9.74</b>	<b>2.21</b>	0.00	12.25	0.00	10.67	15.60
LLaVA-OneVision-72B [52]	<b>26.32</b>	<b>19.23</b>	10.34	<b>23.28</b>	7.53	1.12	7.81	1.12	0.00	12.75	0.00	<b>12.80</b>	<b>18.23</b>
LLaVA-OneVision-7B [52]	16.84	17.31	5.17	9.05	3.28	1.12	8.15	0.80	2.94	9.80	1.68	9.14	12.45
VILA-40B [50]	17.89	7.69	0.00	5.75	0.00	0.00	3.93	0.00	0.00	8.58	0.00	6.40	9.53
VILA-8B [50]	15.79	9.62	1.72	8.91	0.96	0.00	3.46	0.00	2.94	6.37	1.68	6.71	9.27
<i>Open-Source Video MLLMs</i>													
LLaVA-NeXT-Video-32B-Qwen [94]	21.05	7.69	6.90	14.08	5.34	0.00	8.61	0.00	2.94	<b>12.01</b>	0.98	8.84	13.39
LLaVA-Video-72B-Qwen2 [95]	<b>27.37</b>	9.62	<b>15.52</b>	<b>24.28</b>	<b>9.62</b>	1.12	8.05	0.86	0.00	9.80	0.00	12.50	<b>16.95</b>
LLaVA-Video-7B-Qwen2 [95]	20.00	<b>19.23</b>	3.45	4.89	1.88	1.12	<b>8.80</b>	0.27	0.00	5.15	0.00	9.76	12.63
Oryx-34B [58]	18.95	3.85	5.17	13.07	4.89	1.12	7.02	1.00	0.00	8.33	0.00	7.32	11.33
VideoLLaMA2-72B [16]	<b>27.37</b>	9.62	12.07	18.68	6.35	<b>2.25</b>	7.49	<b>1.38</b>	<b>5.88</b>	10.78	<b>2.39</b>	<b>12.81</b>	15.91
VideoLLaMA2-7B [16]	21.05	9.62	6.90	17.53	4.88	0.00	1.63	0.00	2.94	7.35	1.38	9.20	11.99

Table 2. Results of different models on EMBODIEDEVAL (%). The best-performing model in each category is bolded.

471 in embodied scenarios.

472 **Performance Gap between Best Proprietary and**  
473 **Open-source Models.** Proprietary models demonstrate a  
474 consistent advantage across all tasks and metrics. GPT-  
475 4o leads in overall performance, with the highest success  
476 rate on four out of five tasks. Qwen-VL-Max ranks second,  
477 achieving the highest success rate in Attribute QA. In con-  
478 trast, open-source models show a substantial performance  
479 gap. The top performing multi-image MLLM, LLaVA-  
480 OneVision-72B, achieves an overall success rate of 12.80%,  
481 barely competitive with proprietary models. The other  
482 multi-image MLLMs have accuracy rates of around 10%  
483 at best. Among the video MLLMs, VideoLLaMA2-72B  
484 achieves the highest success rate (12.81%), while LLaVA-  
485 Video-72B-Qwen2 achieves the highest GcS (16.95%). For  
486 open-source models, larger models often exhibit improved  
487 performance, but this trend is not universally observed.

488 **Model Performance across Different Task Types.** The  
489 results highlight a significant variation in model perfor-  
490 mance across different task types. GPT-4o demonstrates  
491 relatively strong results in QA and Navigation tasks, but

its performance drops notably for interaction tasks. This  
492 disparity is even more pronounced among other commer-  
493 cial models. For instance, most models perform reason-  
494 ably well in Attribute QA but see a sharp decline in Spa-  
495 tial QA that requires spatial reasoning, often halving their  
496 success rates. The Navigation task shows substantial per-  
497 formance variability across models, such as Gemini-Flash,  
498 which achieves 26.32% in Attribute QA but drops to only  
499 5.17% in Navigation. Overall, the scores for Object Inter-  
500 action and Social Interaction are consistently lower across  
501 all models, underscoring the challenge these models face in  
502 scenarios that require a deeper understanding of affordance  
503 or social cues. 504

**Challenges in Long-Horizon Tasks.** The performance  
505 of MLLMs shows a significant decline as the number of  
506 steps required for the task and the subgoals increase. As  
507 illustrated in Figure 5, models maintain relatively high suc-  
508 cess counts at lower required steps, but their performance  
509 drops fluctuatingly as tasks require longer sequences of ac-  
510 tions. This drop in performance can be attributed to two  
511 primary factors: (1) the increased complexity and reason-  
512

513 ing demands associated with longer tasks and multiple sub-  
 514 goals, which challenge the models' planning and decision-  
 515 making capabilities, and (2) the limitations in retaining and  
 516 processing the long temporal context necessary for success-  
 517 ful execution in extended tasks. These challenges high-  
 518 light the need for further model improvements to enhance  
 519 MLLM capabilities in handling complex, multi-step objec-  
 520 tives over long horizons.

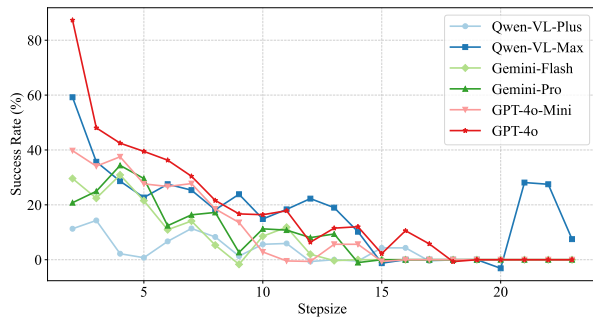


Figure 5. Success rate vs. number of steps required for the task.

521 **Studies on Temperature.** We find that all models per-  
 522 form slightly better at temperature = 1 compared to temper-  
 523 ature = 0. Through observing cases, we believe this is be-  
 524 cause embodied tasks require a certain level of exploration,  
 525 and when the temperature is set to 0, the determinism of  
 526 the output causes the model to easily get stuck in repetitive  
 527 errors. However, in this paper, we propose using a temper-  
 528 ature of 0 as the evaluation standard, as this removes ran-  
 529 domness from the evaluation, improving efficiency and bet-  
 530 ter reflect the model's true capabilities, including its ability  
 531 to recognize and escape from erroneous trajectories.

#### 532 4.4. Error Analysis

533 We summarize four primary error categories in MLLM-  
 534 based embodied agents: **(1) Hallucination in Grounding:**  
 535 Models misperceive the environment, hallucinating nonex-  
 536 istent objects or overlooking present ones. For example,  
 537 models may confidently describe absent items or fail to lo-  
 538 cate small objects like laptops or keys, impacting both QA  
 539 (e.g., providing answers based on imagined objects) and  
 540 non-QA tasks (e.g., failing to navigate to or interact with  
 541 target objects). **(2) Insufficient Exploration:** Agents em-  
 542 ploy suboptimal exploration strategies, hindering informa-  
 543 tion gathering and goal finding due to incomplete environ-  
 544 ment coverage. They are often trapped in local areas, or  
 545 answer before fully exploring the environment due to over-  
 546 confidence. **(3) Lack of Spatial Reasoning:** Models strug-  
 547 gle with understanding spatial relationships. They misinter-  
 548 pret directional instructions (e.g., "to my left") and face dif-  
 549 ficulties navigating between locations, even for simple tasks  
 550 such as moving to or around furniture. **(4) Wrong Plan-  
 551 ning:** Agents demonstrate poor state estimation and action

552 planning. This results in random or repetitive actions, such  
 553 as aimless circling or repeatedly picking up objects. They  
 554 also struggle to understand the outcomes of the action and  
 555 adapt after failed attempts. Figure 6 provides illustrative  
 556 examples of these errors.



Figure 6. Case study of common error categories. In *Hallucination in Grounding*, the agent mistakenly identified a single blue sofa as two. In *Insufficient Exploration*, the agent failed to look for additional items. In *Lack in Spatial Reasoning*, the agent misestimated the distance between objects. In *Wrong Planning*, the agent did not organize the picking up and putting down of the vases in the proper order and at the correct positions.

## 5. Conclusion

557 In this paper, we propose EMBODIEDEVAL, the first inter-  
 558 active benchmark designed for MLLMs with comprehen-  
 559 sive embodied tasks. We provide an efficient framework to  
 560 interactively evaluate the capabilities of MLLMs on embod-  
 561 ied tasks. To ensure the accuracy, diversity, and quality of  
 562 the dataset, extensive efforts are devoted to the annotation  
 563 process for each task sample. 564

565 Through experiments, we found that current MLLMs  
 566 perform poorly on embodied tasks. However, we believe  
 567 there will be more attention to improving the embodied ca-  
 568 pabilities of MLLMs upon the general capabilities learned  
 569 from universal multimodal data. We hope EMBODIEDEVAL  
 570 can help and guide the development of MLLMs to realize  
 571 their potential in embodied intelligence.

572

## References

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

- [1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022. 1
- [2] Peter Anderson and et al. Vln-ce: A visually-grounded language navigation dataset for following compositional instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8470–8480, 2020. 6
- [3] Peter Anderson, Angel Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, et al. On evaluation of embodied navigation agents. *arXiv preprint arXiv:1807.06757*, 2018. 6
- [4] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton Van Den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3674–3683, 2018. 1, 2, 3
- [5] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023. 7
- [6] Dhruv Batra, Angel X Chang, Sonia Chernova, Andrew J Davison, Jia Deng, Vladlen Koltun, Sergey Levine, Jitendra Malik, Igor Mordatch, Roozbeh Mottaghi, et al. Rearrangement: A challenge for embodied ai. *arXiv preprint arXiv:2011.01975*, 2020. 3
- [7] Dhruv Batra, Aaron Gokaslan, Aniruddha Kembhavi, Oleksandr Maksymets, Roozbeh Mottaghi, Manolis Savva, Alexander Toshev, and Erik Wijmans. Objectnav revisited: On evaluation of embodied agents navigating to objects. *arXiv preprint arXiv:2006.13171*, 2020. 6
- [8] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465, 2024. 1, 3
- [9] Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023. 2
- [10] Liang Chen, Yichi Zhang, Shuhuai Ren, Haozhe Zhao, Zefan Cai, Yuchi Wang, Peiyi Wang, Tianyu Liu, and Baobao Chang. Towards end-to-end embodied decision making via multi-modal large language model: Explorations with gpt4-vision and beyond. *arXiv preprint arXiv:2310.02071*, 2023. 1
- [11] Yi Chen, Yuying Ge, Yixiao Ge, Mingyu Ding, Bohao Li, Rui Wang, Ruifeng Xu, Ying Shan, and Xihui Liu. Egoplan-bench: Benchmarking multimodal large language models for human-level planning. *arXiv preprint arXiv:2312.06722*, 2023. 2
- [12] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. 6, 7
- [13] An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision language model. *arXiv preprint arXiv:2406.01584*, 2024. 3
- [14] Sijie Cheng, Zhicheng Guo, Jingwen Wu, Kechen Fang, Peng Li, Huaping Liu, and Yang Liu. Can vision-language models think from a first-person perspective? *arXiv preprint arXiv:2311.15596*, 2023. 2
- [15] Sijie Cheng, Zhicheng Guo, Jingwen Wu, Kechen Fang, Peng Li, Huaping Liu, and Yang Liu. Egothink: Evaluating first-person perspective thinking capability of vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14291–14302, 2024. 1
- [16] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*, 2024. 7
- [17] Zhili Cheng, Zhitong Wang, Jinyi Hu, Shengding Hu, An Liu, Yuge Tu, Pengkai Li, Lei Shi, Zhiyuan Liu, and Maosong Sun. Legent: Open platform for embodied agents. *arXiv preprint arXiv:2404.18243*, 2024. 3
- [18] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–10, 2018. 2, 3
- [19] Matt Deitke, Eli VanderBilt, Alvaro Herrasti, Luca Weihs, Kiana Ehsani, Jordi Salvador, Winson Han, Eric Kolve, Aniruddha Kembhavi, and Roozbeh Mottaghi. Proctor: Large-scale embodied ai using procedural generation. *Advances in Neural Information Processing Systems*, 35:5982–5994, 2022. 5
- [20] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023. 5, 1
- [21] Vishnu Sashank Dorbala, Prasoon Goyal, Robinson Piramuthu, Michael Johnston, Reza Ghanadhan, and Dinesh Manocha. S-eqa: Tackling situational queries in embodied question answering. *arXiv preprint arXiv:2405.04732*, 2024. 3
- [22] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023. 1
- [23] Mengfei Du, Binhao Wu, Zejun Li, Xuanjing Huang, and Zhongyu Wei. Embspatial-bench: Benchmarking spatial un-

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

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 understanding for embodied tasks with large vision-language  
688 models. *arXiv preprint arXiv:2406.05756*, 2024. 3
- 689 [24] Chenyou Fan. Egovqa-an egocentric video question answer-  
690 ing benchmark dataset. In *Proceedings of the IEEE/CVF*  
691 *International Conference on Computer Vision Workshops*,  
692 pages 0–0, 2019. 2
- 693 [25] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin,  
694 Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li,  
695 Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A compre-  
696 hensive evaluation benchmark for multimodal large language  
697 models. *arXiv preprint arXiv:2306.13394*, 2023. 1, 2
- 698 [26] Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai  
699 Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang  
700 Shen, Mengdan Zhang, Peixian Chen, Yanwei Li, Shaohui  
701 Lin, Sirui Zhao, Ke Li, Tong Xu, Xiawu Zheng, Enhong  
702 Chen, Rongrong Ji, and Xing Sun. Video-mme: The first-  
703 ever comprehensive evaluation of multi-modal llms in video  
704 analysis. *arXiv preprint arXiv:2405.21075*, 2024. 2
- 705 [27] Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren,  
706 Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen,  
707 Mengdan Zhang, et al. Video-mme: The first-ever compre-  
708 hensive evaluation benchmark of multi-modal llms in video  
709 analysis. *arXiv preprint arXiv:2405.21075*, 2024. 1
- 710 [28] Chuang Gan, Siyuan Zhou, Jeremy Schwartz, Seth Alter,  
711 Abhishek Bhandwadar, Dan Gutfreund, Daniel LK Yamins,  
712 James J DiCarlo, Josh McDermott, Antonio Torralba, et al.  
713 The threeworld transport challenge: A visually guided task-  
714 and-motion planning benchmark towards physically realistic  
715 embodied ai. In *2022 International conference on robotics*  
716 *and automation (ICRA)*, pages 8847–8854. IEEE, 2022. 3
- 717 [29] Google. Our next-generation model: Gemini 1.5, 2024. 7
- 718 [30] Daniel Gordon, Aniruddha Kembhavi, Mohammad Raste-  
719 gari, Joseph Redmon, Dieter Fox, and Ali Farhadi. Iqa:  
720 Visual question answering in interactive environments. In  
721 *Proceedings of the IEEE conference on computer vision and*  
722 *pattern recognition*, pages 4089–4098, 2018. 3, 4
- 723 [31] Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu,  
724 Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie  
725 Huang, Yuxiang Zhang, et al. Olympiadbench: A chal-  
726 lenging benchmark for promoting agi with olympiad-level  
727 bilingual multimodal scientific problems. *arXiv preprint*  
728 *arXiv:2402.14008*, 2024. 2
- 729 [32] Jinyi Hu, Yuan Yao, Chongyi Wang, Shan Wang, Yinxu  
730 Pan, Qianyu Chen, Tianyu Yu, Hanghao Wu, Yue Zhao,  
731 Haoye Zhang, et al. Large multilingual models pivot zero-  
732 shot multimodal learning across languages. *arXiv preprint*  
733 *arXiv:2308.12038*, 2023. 2
- 734 [33] Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao,  
735 Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui,  
736 Owais Khan Mohammed, Barun Patra, Qiang Liu, Kriti Ag-  
737 garwal, Zewen Chi, Johan Bjorck, Vishrav Chaudhary, Sub-  
738 hojit Som, Xia Song, and Furu Wei. Language is not all you  
739 need: Aligning perception with language models. In *Thirty-*  
740 *seventh Conference on Neural Information Processing Sys-*  
741 *tems*, 2023. 2
- 742 [34] Md Mofijul Islam, Alexi Gladstone, Riashat Islam, and Tariq  
743 Iqbal. Eqa-mx: Embodied question answering using multi-  
modal expression. In *The Twelfth International Conference*  
744 *on Learning Representations*, 2023. 3 745
- [35] Md Mofijul Islam, Alexi Gladstone, Riashat Islam, and Tariq  
746 Iqbal. EQA-MX: Embodied question answering using multi-  
747 modal expression. In *The Twelfth International Conference*  
748 *on Learning Representations*, 2024. 2 749
- [36] Vihan Jain, Gabriel Magalhaes, Alexander Ku, Ashish  
750 Vaswani, Eugene Ie, and Jason Baldridge. Stay on the path:  
751 Instruction fidelity in vision-and-language navigation. *arXiv*  
752 *preprint arXiv:1905.12255*, 2019. 2, 4 753
- [37] Amita Kamath, Jack Hessel, and Kai-Wei Chang. What’s’  
754 up” with vision-language models? investigating their strug-  
755 gle with spatial reasoning. *arXiv preprint arXiv:2310.19785*,  
756 2023. 3 757
- [38] Yash Kant, Arun Ramachandran, Sriram Yenamandra,  
758 Igor Gilitschenski, Dhruv Batra, Andrew Szot, and Harsh  
759 Agrawal. Housekeep: Tidying virtual households using com-  
760 monsense reasoning. In *European Conference on Computer*  
761 *Vision*, pages 355–373. Springer, 2022. 3 762
- [39] Mukul Khanna, Yongsen Mao, Hanxiao Jiang, Sanjay  
763 Hareesh, Brennan Shacklett, Dhruv Batra, Alexander Clegg,  
764 Eric Undersander, Angel X Chang, and Manolis Savva.  
765 Habitat synthetic scenes dataset (hssd-200): An analysis of  
766 3d scene scale and realism tradeoffs for objectgoal naviga-  
767 tion. In *Proceedings of the IEEE/CVF Conference on Com-*  
768 *puter Vision and Pattern Recognition*, pages 16384–16393,  
769 2024. 5 770
- [40] Mukul Khanna, Ram Ramrakhya, Gunjan Chhablani, Sriram  
771 Yenamandra, Theophile Gervet, Matthew Chang, Zsolt Kira,  
772 Devendra Singh Chplot, Dhruv Batra, and Roozbeh Mot-  
773 taghi. Goat-bench: A benchmark for multi-modal lifelong  
774 navigation. In *Proceedings of the IEEE/CVF Conference*  
775 *on Computer Vision and Pattern Recognition*, pages 16373–  
776 16383, 2024. 2 777
- [41] Byeonghwi Kim, Jinyeon Kim, Yuyeong Kim, Cheolhong  
778 Min, and Jonghyun Choi. Context-aware planning and  
779 environment-aware memory for instruction following em-  
780 bodied agents. In *Proceedings of the IEEE/CVF Interna-*  
781 *tional Conference on Computer Vision*, pages 10936–10946,  
782 2023. 6 783
- [42] Taewoong Kim, Cheolhong Min, Byeonghwi Kim, Jinyeon  
784 Kim, Wonje Jeung, and Jonghyun Choi. Realfred: An em-  
785 bodied instruction following benchmark in photo-realistic  
786 environments. *arXiv preprint arXiv:2407.18550*, 2024. 3 787
- [43] Eric Kolve and et al. Ai2-thor: An interactive 3d environ-  
788 ment for visual ai. *arXiv preprint arXiv:1712.05474*, 2017.  
789 5, 6 790
- [44] Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and  
791 Jason Baldridge. Room-across-room: Multilingual vision-  
792 and-language navigation with dense spatiotemporal ground-  
793 ing. *arXiv preprint arXiv:2010.07954*, 2020. 2 794
- [45] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yix-  
795 iao Ge, and Ying Shan. Seed-bench: Benchmarking multi-  
796 modal llms with generative comprehension. *arXiv preprint*  
797 *arXiv:2307.16125*, 2023. 1 798
- [46] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen,  
799 Sanjana Srivastava, Roberto Martín-Martín, Chen Wang,  
800 800

- 801 Gabrael Levine, Michael Lingelbach, Jiankai Sun, et al.  
802 Behavior-1k: A benchmark for embodied ai with 1,000 ev-  
803 eryday activities and realistic simulation. In *Conference on*  
804 *Robot Learning*, pages 80–93. PMLR, 2023. 3
- [47] Kunchang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai  
805 Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao.  
806 Videochat: Chat-centric video understanding. *arXiv preprint*  
807 *arXiv:2305.06355*, 2023. 2
- [48] Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and  
808 Li Yuan. Video-llava: Learning united visual represen-  
809 tation by alignment before projection. *arXiv preprint*  
810 *arXiv:2311.10122*, 2023. 2
- [49] Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Moham-  
811 mad Shoeybi, and Song Han. Vila: On pre-training for vi-  
812 sual language models. In *Proceedings of the IEEE/CVF Con-*  
813 *ference on Computer Vision and Pattern Recognition*, pages  
814 26689–26699, 2024. 2
- [50] Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Moham-  
815 mad Shoeybi, and Song Han. Vila: On pre-training for vi-  
816 sual language models. In *Proceedings of the IEEE/CVF Con-*  
817 *ference on Computer Vision and Pattern Recognition*, pages  
818 26689–26699, 2024. 7
- [51] Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spa-  
819 tial reasoning. *Transactions of the Association for Computa-*  
820 *tional Linguistics*, 11:635–651, 2023. 3
- [52] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee.  
821 Improved baselines with visual instruction tuning. In *Pro-*  
822 *ceedings of the IEEE/CVF Conference on Computer Vision*  
823 *and Pattern Recognition*, pages 26296–26306, 2024. 7
- [53] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.  
824 Visual instruction tuning. *Advances in neural information*  
825 *processing systems*, 36, 2024. 1, 2
- [54] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei,  
826 Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan  
827 Yang, et al. Agentbench: Evaluating llms as agents. *arXiv*  
828 *preprint arXiv:2308.03688*, 2023. 2
- [55] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei,  
829 Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan  
830 Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengx-  
831 iao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu  
832 Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang.  
833 Agentbench: Evaluating LLMs as agents. In *The Twelfth In-*  
834 *ternational Conference on Learning Representations*, 2024.  
835 6
- [56] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang  
836 Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He,  
837 Ziwei Liu, et al. Mmbench: Is your multi-modal model an  
838 all-around player? *arXiv preprint arXiv:2307.06281*, 2023.  
839 1, 2
- [57] Yuliang Liu, Zhang Li, Biao Yang, Chunyuan Li, Xucheng  
840 Yin, Cheng-lin Liu, Lianwen Jin, and Xiang Bai. On the  
841 hidden mystery of ocr in large multimodal models. *arXiv*  
842 *preprint arXiv:2305.07895*, 2023. 2
- [58] Zuyan Liu, Yuhao Dong, Ziwei Liu, Winston Hu, Ji-  
843 wen Lu, and Yongming Rao. Oryx mllm: On-demand  
844 spatial-temporal understanding at arbitrary resolution. *arXiv*  
845 *preprint arXiv:2409.12961*, 2024. 7
- [59] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei  
846 Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and  
847 Ashwin Kalyan. Learn to explain: Multimodal reasoning  
848 via thought chains for science question answering. In *The*  
849 *36th Conference on Neural Information Processing Systems*  
850 *(NeurIPS)*, 2022. 2
- [60] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li,  
851 Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel  
852 Galley, and Jianfeng Gao. Mathvista: Evaluating mathe-  
853 matical reasoning of foundation models in visual contexts.  
854 In *The Twelfth International Conference on Learning Repre-*  
855 *sentations*, 2024. 2
- [61] Ji Ma, Hongming Dai, Yao Mu, Pengying Wu, Hao Wang,  
856 Xiaowei Chi, Yang Fei, Shanghang Zhang, and Chang  
857 Liu. Doze: A dataset for open-vocabulary zero-shot ob-  
858 ject navigation in dynamic environments. *arXiv preprint*  
859 *arXiv:2402.19007*, 2024. 2
- [62] Yuen Ma, Zixing Song, Yuzheng Zhuang, Jianye Hao, and  
860 Irwin King. A survey on vision-language-action models for  
861 embodied ai. *arXiv preprint arXiv:2405.14093*, 2024. 1
- [63] Arjun Majumdar, Anurag Ajay, Xiaohan Zhang, Pranav  
862 Putta, Sriram Yenamandra, Mikael Henaff, Sneha Silwal,  
863 Paul Mcvay, Oleksandr Maksymets, Sergio Arnaud, et al.  
864 Openeqa: Embodied question answering in the era of founda-  
865 tion models. In *Proceedings of the IEEE/CVF Conference*  
866 *on Computer Vision and Pattern Recognition*, pages 16488–  
867 16498, 2024. 1, 2
- [64] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and  
868 Anirban Chakraborty. Ocr-vqa: Visual question answering  
869 by reading text in images. In *2019 international conference*  
870 *on document analysis and recognition (ICDAR)*, pages 947–  
871 952. IEEE, 2019. 2
- [65] Dipendra Misra, Andrew Bennett, Valts Blukis, Eyvind  
872 Niklasson, Max Shatkhin, and Yoav Artzi. Mapping instruc-  
873 tions to actions in 3d environments with visual goal predic-  
874 tion. *arXiv preprint arXiv:1809.00786*, 2018. 3, 4
- [66] Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang,  
875 Mingyu Ding, Jun Jin, Bin Wang, Jifeng Dai, Yu Qiao, and  
876 Ping Luo. Embodiedgpt: Vision-language pre-training via  
877 embodied chain of thought. *Advances in Neural Information*  
878 *Processing Systems*, 36, 2024. 1
- [67] OpenAI. Gpt-4v(ision) system card, 2023. 1
- [68] OpenAI. Hello gpt4-o, 2024. 7
- [69] OpenGVLab. Internvl2: Better than the best—expanding  
879 performance boundaries of open-source multimodal models  
880 with the progressive scaling strategy, 2024. 7
- [70] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan  
881 Huang, Shuming Ma, and Furu Wei. Kosmos-2: Ground-  
882 ing multimodal large language models to the world. *arXiv*  
883 *preprint arXiv:2306.14824*, 2023. 2
- [71] Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang,  
884 William Yang Wang, Chunhua Shen, and Anton van den  
885 Hengel. Reverie: Remote embodied visual referring ex-  
886 pression in real indoor environments. In *Proceedings of*  
887 *the IEEE/CVF Conference on Computer Vision and Pattern*  
888 *Recognition*, pages 9982–9991, 2020. 2
- [72] Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang,  
889 William Yang Wang, Chunhua Shen, and Anton van den  
890 911 912 913 914 915

- 916 Hengel. Reverie: Remote embodied visual referring expres-  
917 sion in real indoor environments, 2020. 2
- 918 [73] Allen Z Ren, Jaden Clark, Anushri Dixit, Masha Itkina,  
919 Anirudha Majumdar, and Dorsa Sadigh. Explore until confi-  
920 dent: Efficient exploration for embodied question answering.  
921 *arXiv preprint arXiv:2403.15941*, 2024. 3
- 922 [74] Manolis Savva and et al. Habitat challenge: A photoreal-  
923 istic embodied ai benchmark. In *Proceedings of the IEEE*  
924 *International Conference on Computer Vision (ICCV)*, pages  
925 9338–9346. IEEE, 2019. 6
- 926 [75] Manasi Sharma. Exploring and improving the spatial reason-  
927 ing abilities of large language models. In *I Can’t Believe It’s*  
928 *Not Better Workshop: Failure Modes in the Age of Founda-*  
929 *tion Models*, 2023. 3
- 930 [76] Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan  
931 Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer,  
932 and Dieter Fox. Alfred: A benchmark for interpreting  
933 grounded instructions for everyday tasks. In *Proceedings of*  
934 *the IEEE/CVF conference on computer vision and pattern*  
935 *recognition*, pages 10740–10749, 2020. 2, 4, 6
- 936 [77] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté,  
937 Yonatan Bisk, Adam Trischler, and Matthew Hausknecht.  
938 Alfvorld: Aligning text and embodied environments for in-  
939 teractive learning. *arXiv preprint arXiv:2010.03768*, 2020.  
940 2
- 941 [78] Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang,  
942 Xinlei Chen, Devi Parikh, and Marcus Rohrbach. Towards  
943 vqa models that can read. In *Proceedings of the IEEE Con-*  
944 *ference on Computer Vision and Pattern Recognition*, pages  
945 8317–8326, 2019. 2
- 946 [79] Sanjana Srivastava, Chengshu Li, Michael Lingelbach,  
947 Roberto Martín-Martín, Fei Xia, Kent Elliott Vainio, Zheng  
948 Lian, Cem Gokmen, Shyamal Buch, Karen Liu, et al. Behav-  
949 ior: Benchmark for everyday household activities in virtual,  
950 interactive, and ecological environments. In *Conference on*  
951 *robot learning*, pages 477–490. PMLR, 2022. 2, 3, 4
- 952 [80] Hai-Long Sun, Da-Wei Zhou, Yang Li, Shiyin Lu, Chao Yi,  
953 Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, De-  
954 Chuan Zhan, et al. Parrot: Multilingual visual instruction  
955 tuning. *arXiv preprint arXiv:2406.02539*, 2024. 2
- 956 [81] Sinan Tan, Mengmeng Ge, Di Guo, Huaping Liu, and  
957 Fuchun Sun. Knowledge-based embodied question answer-  
958 ing. *IEEE Transactions on Pattern Analysis and Machine*  
959 *Intelligence*, 45(10):11948–11960, 2023. 2, 3
- 960 [82] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui  
961 Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan  
962 Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a  
963 family of highly capable multimodal models. *arXiv preprint*  
964 *arXiv:2312.11805*, 2023. 1
- 965 [83] Luca Weihs, Matt Deitke, Aniruddha Kembhavi, and  
966 Roozbeh Mottaghi. Visual room rearrangement. In *Proceed-*  
967 *ings of the IEEE/CVF conference on computer vision and*  
968 *pattern recognition*, pages 5922–5931, 2021. 3, 4
- 969 [84] Jingkan Yang, Yuhao Dong, Shuai Liu, Bo Li, Ziyue Wang,  
970 Chencheng Jiang, Haoran Tan, Jiamu Kang, Yuanhan Zhang,  
971 Kaiyang Zhou, et al. Octopus: Embodied vision-language  
972 programmer from environmental feedback. *arXiv preprint*  
973 *arXiv:2310.08588*, 2023. 1
- [85] Yue Yang, Fan-Yun Sun, Luca Weihs, Eli VanderBilt, Al-  
varo Herrasti, Winson Han, Jiajun Wu, Nick Haber, Ranjay  
Krishna, Lingjie Liu, Chris Callison-Burch, Mark Yatskar,  
Aniruddha Kembhavi, and Christopher Clark. Holodeck:  
Language guided generation of 3d embodied ai environ-  
ments, 2024. 1
- [86] Sriram Yenamandra, Arun Ramachandran, Karmesh Yadav,  
Austin Wang, Mukul Khanna, Theophile Gervet, Tsung-Yen  
Yang, Vidhi Jain, Alexander William Clegg, John Turner,  
et al. Homerobot: Open-vocabulary mobile manipulation.  
*arXiv preprint arXiv:2306.11565*, 2023. 3, 4
- [87] Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen  
Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and  
Yinfei Yang. Ferret: Refer and ground anything anywhere  
at any granularity. *arXiv preprint arXiv:2310.07704*, 2023.  
2
- [88] Licheng Yu, Xinlei Chen, Georgia Gkioxari, Mohit Bansal,  
Tamara L Berg, and Dhruv Batra. Multi-target embodied  
question answering. In *Proceedings of the IEEE/CVF Con-*  
*ference on Computer Vision and Pattern Recognition*, pages  
6309–6318, 2019. 2, 3
- [89] Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng  
Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao Zheng,  
Maosong Sun, et al. Rllfv: Towards trustworthy mllms  
via behavior alignment from fine-grained correctional hu-  
man feedback. In *Proceedings of the IEEE/CVF Conference*  
*on Computer Vision and Pattern Recognition*, pages 13807–  
13816, 2024. 2
- [90] Tianyu Yu, Haoye Zhang, Yuan Yao, Yunkai Dang, Da Chen,  
Xiaoman Lu, Ganqu Cui, Taiwen He, Zhiyuan Liu, Tat-Seng  
Chua, and Maosong Sun. Rlaif-v: Aligning mllms through  
open-source ai feedback for super gpt-4v trustworthiness.  
*arXiv preprint arXiv:2405.17220*, 2024. 2
- [91] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi  
Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming  
Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Ren-  
liang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo  
Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhao Chen.  
Mmmu: A massive multi-discipline multimodal understand-  
ing and reasoning benchmark for expert agi. In *Proceedings*  
*of CVPR*, 2024. 1, 2
- [92] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi  
Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming  
Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline  
multimodal understanding and reasoning benchmark for ex-  
pert agi. In *Proceedings of the IEEE/CVF Conference*  
*on Computer Vision and Pattern Recognition*, pages 9556–  
9567, 2024. 6
- [93] Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin,  
Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei  
Chang, Peng Gao, et al. Mathverse: Does your multi-modal  
llm truly see the diagrams in visual math problems? *arXiv*  
*preprint arXiv:2403.14624*, 2024. 2
- [94] Yuanhan Zhang, Bo Li, Haotian Liu, Yong Jae Lee, Liangke  
Gui, Di Fu, Jiashi Feng, Ziwei Liu, and Chunyuan Li. Llava-  
next: A strong zero-shot video understanding model, 2024.  
7

- 1031 [95] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Zi-  
1032 wei Liu, and Chunyuan Li. Video instruction tuning with  
1033 synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. 7
- 1034 [96] Zhuofan Zhang, Ziyu Zhu, Pengxiang Li, Tengyu Liu, Xiao-  
1035 jian Ma, Yixin Chen, Baoxiong Jia, Siyuan Huang, and Qing  
1036 Li. Task-oriented sequential grounding in 3d scenes. *arXiv*  
1037 *preprint arXiv:2408.04034*, 2024. 1
- 1038 [97] Kaizhi Zheng, Xiaotong Chen, Odest Chadwicke Jenkins,  
1039 and Xin Wang. VImbench: A compositional benchmark for  
1040 vision-and-language manipulation. *Advances in Neural In-*  
1041 *formation Processing Systems*, 35:665–678, 2022. 6
- 1042 [98] Fengda Zhu, Xiwen Liang, Yi Zhu, Qizhi Yu, Xiaojun  
1043 Chang, and Xiaodan Liang. Soon: Scenario oriented object  
1044 navigation with graph-based exploration. In *Proceedings of*  
1045 *the IEEE/CVF Conference on Computer Vision and Pattern*  
1046 *Recognition*, pages 12689–12699, 2021. 2
- 1047 [99] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted  
1048 Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker,  
1049 Ayzaan Wahid, et al. Rt-2: Vision-language-action models  
1050 transfer web knowledge to robotic control. In *Conference on*  
1051 *Robot Learning*, pages 2165–2183. PMLR, 2023. 1

1052 **6. Appendix**1053 **6.1. Task Samples**

1054 We selected some representative examples to illustrate the  
1055 diversity of the task set in Table 3.

1056 **6.2. Details of Evaluation Framework**1057 **6.2.1 Evaluation Formulation**

1058 We formalize the evaluation process mentioned in Sec-  
1059 tion 3.2.3 as follows. The prompt we used for the MLLMs  
1060 during the evaluation is shown in Figure 8.

**Algorithm 1** EMBODIEDEVAL Evaluation Process

**Input:** A Multimodal LLM  $\pi$ , an evaluation task including a scene  $x$ , a task description  $g$ , an option list  $\mathcal{C} = a_0, a_1, \dots, a_n$ , and a predicate list  $\mathcal{P}$ .

**Output:** A boolean indicating whether the task was successful *success*.

```

1:  $o, s \leftarrow E.reset(x)$   $\triangleright E$  is the simulator,  $o$  is the
   observed image,  $s$  is the world state
2:  $H_o \leftarrow \{o\}$   $\triangleright$  observation history
3:  $H_a \leftarrow \emptyset$   $\triangleright$  action history
4: for  $i \leftarrow 0$  to max steps do
5:    $a \leftarrow \pi.predict(g, \mathcal{C}, H_o, H_a)$ 
6:    $o, s \leftarrow E.step(a)$ 
7:    $H_o.append(o)$ 
8:    $H_a.append(a)$ 
9:    $done \leftarrow P.judge(s)$ 
10:  if  $done$  then
11:    return true
12:  else if  $a$  is answer action then
13:    return false  $\triangleright$  wrong answer
14:  end if
15: end for
16: return false  $\triangleright$  reach the max steps

```

1061 **6.2.2 Interaction Actions**

1062 We provided more examples of the interaction space men-  
1063 tioned in Section 3.2.1 in Table 5.

1064 **6.2.3 Predicates**

1065 All the predicate functions described in Section 3.2.2 are  
1066 listed in Table 4.

1067 **6.3. Creation of Objaverse Synthetic**

1068 We use a wide variety of objects from Objaverse to pro-  
1069 cedurally generate diverse scenes and further refine them  
1070 through interactive scene editing.

**Object Selection.** We curated a subset of indoor assets out of Holodeck’s [85] annotated realistic and diverse objects chosen from the Objaverse asset library [20]. To ensure quality, we employed GPT-3.5 to filter unsuitable outdoor objects and manually reviewed frontal renderings to remove low-quality assets. This process resulted in a database of about 15,000 objects spanning over 500 categories (see examples in Figure 9).

**Scene Generation.** We leveraged GPT-3.5 to annotate object categories with their typical room occurrences (e.g., inLivingRoom, inKitchen), positions (e.g., onWall, onFloor, onEdge), and functions (e.g., receptacle, pickup). Gemini-1.5-Flash was used to annotate large objects’ orientations. Subsequently, a procedural approach was employed to randomly place architectural elements such as walls, doors, and windows. Large objects were then arranged on the floor either against the walls or in the center of the rooms, and smaller items were finally placed on surfaces of large receptacles. Hundreds of scenes were generated randomly, from which we selected 15 living rooms, 15 bedrooms, 10 two-room, 5 three-room, and 5 four-room for further editing.

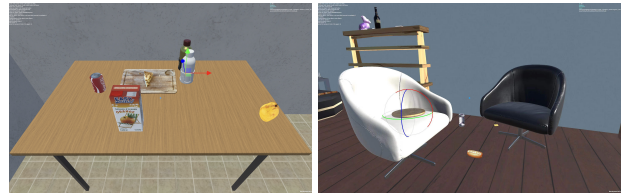


Figure 7. Interactive scene editor: adjust object position (left) and angle (right).

**Scene Editing.** To make the scene more organized and to avoid errors caused by automatic generation, we also edited the generated scene by developing a runtime scene editor. Users can view the type and description of objects, and adjust their position and orientation (see Fig 7). Once editing is complete, the scene can be saved as a JSON file, which can be imported to reproduce the environment.

1061 **6.2.2 Interaction Actions**

1062 We provided more examples of the interaction space men-  
1063 tioned in Section 3.2.1 in Table 5.

1064 **6.2.3 Predicates**

1065 All the predicate functions described in Section 3.2.2 are  
1066 listed in Table 4.

1067 **6.3. Creation of Objaverse Synthetic**

1068 We use a wide variety of objects from Objaverse to pro-  
1069 cedurally generate diverse scenes and further refine them  
1070 through interactive scene editing.

Task	Characteristics
Please go to the kitchen, then come back and tell me if there are any extra cups.	scene memory
Imagine the house is rotated 90 degrees counterclockwise. How would this affect the natural light distribution in the room?	spatial imagination
Open a black locked drawer with a key found on the desk.	tool use
Pick up the kettle and the box labeled "BREAD" from the kitchen counter and place them on the table with the coffee machine.	optical character recognition
Optimize the display of artworks on the shelves as follows: place two items on each shelf, with one shelf featuring two items of the same shape. Complete the requirements in as few steps as possible.	reasoning and planning
Grab the object that is cylindrical and silver on the table next to the washing machine.	multiple attribute reference
Estimate the percentage of floor space occupied by furniture in the room you're currently in.	area estimation
Estimate the straight-line distance from the front door to the TV. Note that each step you take forward is approximately two meters.	distance estimation
Which is closer to the drink on the round table, the ginger or the ice cream?	distance comparison
If we were to host a birthday party, which area of the house could accommodate the most people while ensuring clear pathways to exits?	logic, space, and common sense
Describe the path from the kitchen to the living room.	path description
If you were to draw a straight line from the desk with a turned-on laptop to the bookshelf, which pieces of furniture would it intersect?	spatial reasoning
What is the object I am pointing at?	pointing comprehension
Pick up the watermelon on my right.	perspective-taking comprehension
My red glasses are missing. Please help me look for them in the room. Once you find them, bring them to me.	object searching and delivering
Get close to the lady in white and ask if she needs help.	social navigation
Wake up my dad. He is sleeping in the bedroom. The bedroom is the second room on your right as you walk forward.	finding someone
Enter the dining area and see if there is more than one door in the entire house.	object counting
Calculate the ratio of seating options to the number of rooms in the house.	counting and calculation
Tell me which objects have a handle in the kitchen.	attribute grounding
Evaluate whether the painting above the living room sofa is more colorful than the carpet.	attribute comparison
How many rooms are there in total?	room counting
Confirm if a garbage can is located on the floor in the living room.	object existence
Which room has more seating options, the kitchen or the living room?	quantity comparison
I'm hungry. Find all objects that can be used as ingredients. on the table in this room.	object functionality
Count the maximum number of identical clocks among all the rooms.	counting and attribute memory
What do you think the owner of this room probably studies?	common sense
Is there an egg inside the fridge?	interaction and answering
Open the drawer of the side table in the study room. If there is something inside, leave it open and put all similar items from the room into it. If there is nothing inside, close it.	logical execution

Table 3. Examples of the diverse tasks in EMBODIEDEVAL.

Predicate	Parameters	Success Conditions
<i>choose</i>	The right answer.	When the agent selects the correct answer.
<i>agent_at</i>	A navigation point.	When the agent finally arrives at this point.
<i>agent_pass</i>	A navigation point.	When the agent has passed through this point at least once.
<i>at</i>	An object and a specific point.	When the object is at this point.
<i>grab_once</i>	An object.	When the agent has picked up this object at least once.
<i>grab</i>	An object.	When the agent picks up the object.
<i>special_action_success</i>	An interaction action.	When this interaction action has been successful.

Table 4. The predicates involved in EMBODIEVAL.

Action Text	Execution Requirements
wash	When the agent is holding the target object and stand next to the sink.
hand over	When the agent is holding the target object and stand next to the person.
sit down	When the agent is next to the target chair.
unlock	When the agent is holding the target key and standing next to the drawer
greet	When the agent is near the person.
ask	When the agent is near the person.
mix	When several target beverages are on the table next to the agent.
wipe off the table	When the agent is holding an object for cleaning and standing next to the table.
check the results of the program	When the agent is next to the computer.

Table 5. Some cases of the interaction actions involved in EMBODIEVAL.

### Prompt for Multi-image MLLMs

You are an intelligent vision-language embodied agent skilled at solving tasks **and** answering questions **in** a 3D environment. Your job **is** to efficiently complete a specified task by choosing the optimal action at each timestep **from** a **set** of available actions. You are given a series of ego-centric images, **and** a history of previous actions **with** optional feedback (success/failure **or** human response). Each image shows what you see at a particular step **in** the action history, along **with** an extra image showing your current view.

```
Current task: {}
Action history (action -> feedback): {}
Visual history: {}
Current view: {}
For the current step, your available options are listed as "[Option_Number]._Content" as follows
: {}
Choose your action from the above options by replying with "Thought:_Your_reasoning.\nChoice:_[
Option_Number]_(e.g._[1])".
```

Note:

- If the task needs more information of the scene, navigate wisely to the required targets ( objects, places, **or** people).
- Avoid repeated actions like useless forward motion **and** circling.
- You can only interact **with** objects **or** humans (e.g. pick/place/**open**/close/handover) **if** they are within your view **and** very close to you.
- You can only hold one **object** at a time. Put down **any** held **object** before picking up another.
- Tasks containing "I" **or** "me" are requested by a person **in** the scene.
- Reflect on why previous actions fail to avoid repeating mistakes **and** adjust your current action.
- You have a limited number of {} steps to complete the task.

Figure 8. Prompt for Multi-image MLLMs.



Figure 9. Examples of selected Objaverse assets and views of generated scenes.

### 6.4. Success Cases

1100

We present successful cases accomplished by closed-source MLLMs to gain deeper insights into their current capabilities. As discussed in Section 4.3, the models generally scored low and successfully completed only a limited number of tasks. A closer examination of these successful tasks reveals that they are typically simpler, involve fewer steps and require interaction with fewer objects. To better illustrate these findings, we present representative cases from five task types. These examples highlight the underlying behavioral patterns and reasoning processes of the models during task execution.

1101

1102

1103

1104

1105

#### 6.4.1 Attribute QA

1106

**Task: What do you think the owner of this room probably studies?**

**Options:**

- answer "interior design"
- answer "computer science"
- answer "molecular biology"
- answer "music composition"
- answer "comparative literature"
- answer "theoretical physics"
- answer "digital marketing"
- answer "industrial design"

**GPT-4o: answer "interior design"**



The current view doesn't provide enough information about the room's contents or any clues about the owner's field of study. I should explore the room further to gather more context.

The current view shows a desk with books and an open notebook, which might provide clues about the subject of study. I should examine the contents more closely to make an informed decision.

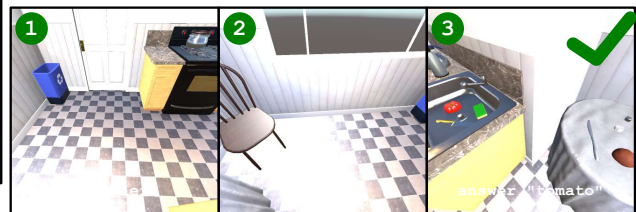
The current view shows books related to architecture and design, along with an open book that appears to have design content. This suggests that the owner of the room likely studies something related to design.

**Task: What is inside the sink near the dining table?**

**Options:**

- answer "potato"
- answer "water"
- answer "bowl"
- answer "apple"
- answer "soap"
- answer "cup"
- answer "steel wool"
- answer "tomato"

**GPT-4o: answer "tomato"**



I need to locate the dining table and the sink. Since the current view shows a kitchen area, I should explore further to find the dining table and sink.

I need to locate the sink near the dining table. The current view shows a chair and a recycling bin, but no sink. I should continue exploring to find the sink.

The current view shows the sink, and I can see items inside it. I need to identify what is inside the sink near the dining table.

Figure 10

**Task: What color is the car parked around the house?**

**Options:**

- answer "red"
- answer "white"
- answer "black"
- answer "blue"
- answer "red and white"
- answer "blue and white"
- answer "black and white"
- answer "yellow"

**Qwen-VL-Max:** answer "red and white"

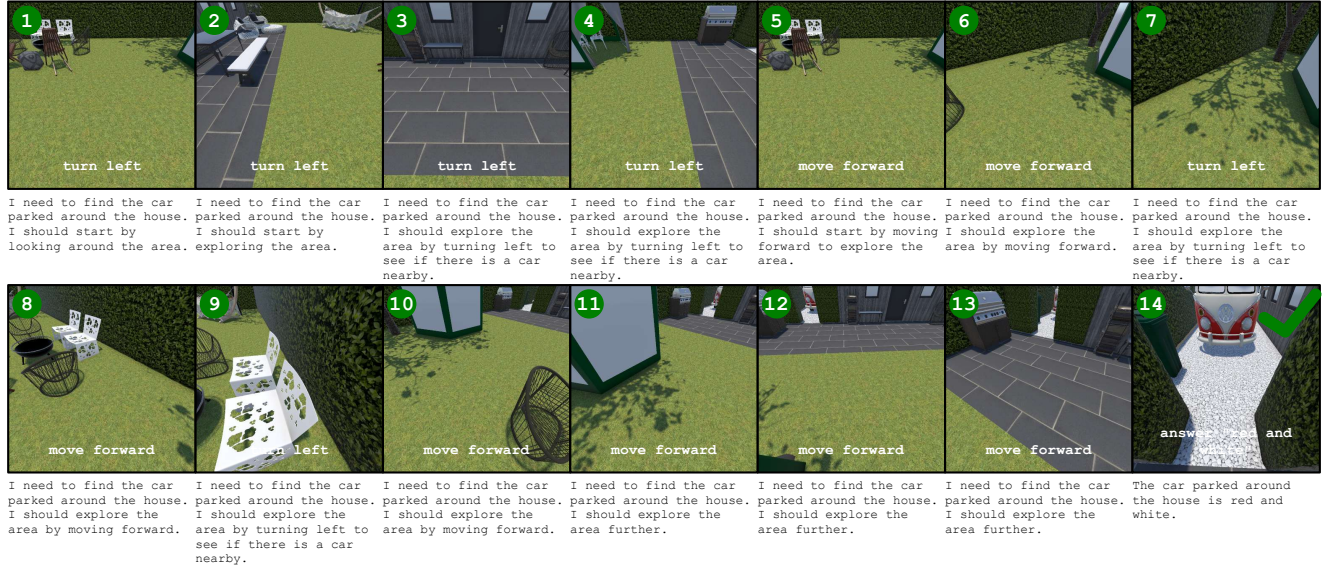


Figure 11

**Task: Compare the sizes of two side tables in the bedroom.**

**Options:**

- answer "The side table with the lamp is located near the window."
- answer "Both side tables are of equal size."
- answer "The side table with the alarm clock is larger than the one with the lamp."
- answer "The side table with the alarm clock is painted blue."
- answer "It's difficult to tell the difference in size between the two side tables."
- answer "The side table with the lamp is larger than the one with the alarm clock."
- answer "The side table with the alarm clock has a drawer."
- answer "The side table with the lamp has two drawers."

**Gemini-Flash:** answer "The side table with the lamp is larger than the one with the alarm clock."

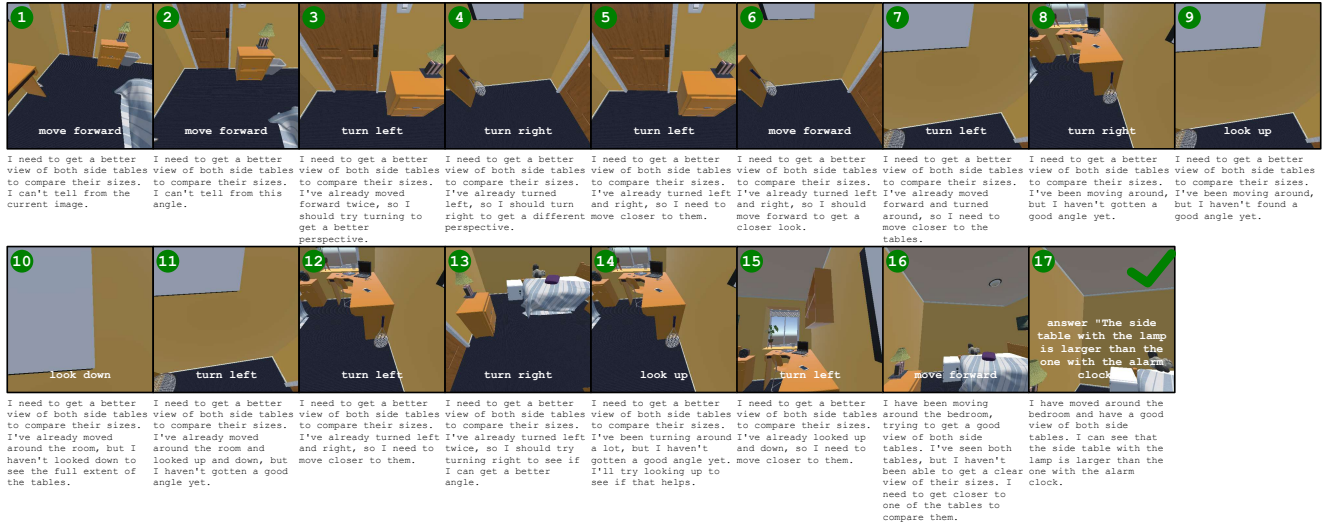


Figure 12

6.4.2 Spatial QA

1107

**Task: Determine if the sofas on the left will fit through the front door.**

**Options:**

- answer "There is no front door."
- answer "Yes, the front door is wide enough, and all the sofas can pass through easily in any orientation.."
- answer "No, only the smallest sofa can pass through."
- answer "Yes, but the sofas need to be stood upright to pass through."
- answer "No, the largest sofa cannot pass through."
- answer "No, none of the sofas can pass through."
- answer "There is no sofa."
- answer "Yes, but the sofas need to be turned sideways to pass through."

**GPT-4o: answer "Yes, but the sofas need to be turned sideways to pass through."**



I need to gather more information about the sofas to determine if they can fit through the front door. Since the current view is of the front door, I should turn left to check the sofas.

I need to determine if the sofas can fit through the front door. I should compare the size of the sofas with the door. The current view shows the sofas, but I need to see the door again to make a comparison.

I need to determine if the sofas can fit through the front door. I should first look at the door to assess its size.

I need to assess the size of the sofas and the front door to determine if they can fit through. The current view is focused on the door, but I need to see the sofas again to make a comparison.

I need to determine if the sofas can fit through the front door. I should compare the size of the sofas to the door. I have already seen the door and the sofas, so I can make a decision based on the current view.

Figure 13

**Task: What direction does the bedroom doorway face when looking out?**

**Options:**

- answer "Opposite to the orientation of the table."
- answer "Opposite to the orientation of the chair with many objects on it."
- answer "Consistent with the orientation as the table."
- answer "Consistent with the orientation as the toilet."
- answer "Consistent with the orientation as the white cabinet."
- answer "Consistent with the orientation of the bed."
- answer "Consistent with the orientation as the fridge."
- answer "Opposite to the orientation of the bed."

**GPT-4o: answer "Consistent with the orientation of the bed."**



The current view shows a table with objects on it, and a door is visible. To determine the direction of the bedroom doorway, I need to explore more of the environment to understand the layout better.

I need to gather more information about the room's layout to determine the direction of the bedroom doorway. Since the current view is limited, I should explore further.

I need to determine the orientation of the bedroom doorway. The current view shows a wall with a painting, which doesn't provide enough information about the doorway. I should explore more by turning to get a better view of the room layout.

I need to determine the orientation of the bedroom doorway. The current view doesn't show the doorway, so I should explore further to find it.

I need to determine the orientation of the bedroom doorway. The current view shows the bed and windows, which might help in understanding the room's layout. Since the task is about the doorway's direction, I should explore more to find the doorway.

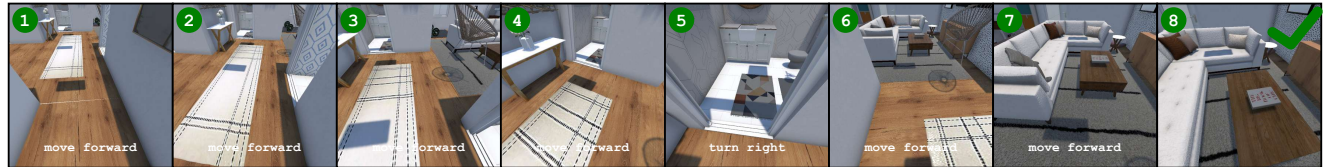
I need to determine the orientation of the bedroom doorway. Based on the current view and previous images, I should consider the orientation of the bed, as it is a significant feature in the room.

Figure 14

### 6.4.3 Navigation

**Task:** Walk straight past the kitchen on your left, turn right when you are almost at the bathroom door, and stop when you reach the book in the center of the sofa table.

**Qwen-VL-Max:**



I need to walk straight past the kitchen on my left, which means I should continue moving forward.

I need to walk straight past the kitchen on my left, so I should continue moving forward.

I need to walk straight past the kitchen on my left, which means I should continue moving forward.

I need to walk straight past the kitchen on my left, which means I should continue moving forward until I am almost at the bathroom door. Then, I will turn right.

I need to turn right after passing the kitchen and before reaching the bathroom door.

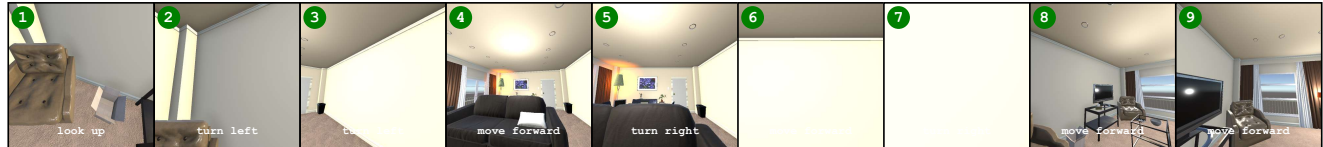
I need to walk straight past the kitchen on my left, which means I should continue moving forward until I reach the point where I need to turn right. The current view shows that I am still in the hallway, so I should keep moving forward.

I need to move closer to the sofa table to find the book.

Figure 15

**Task:** I want to find a human-shaped statue which is placed on the table below the wall where the iris painting hangs. Please navigate to it.

**Qwen-VL-Max:**



I need to find the human-shaped statue on the table below the wall where the iris painting hangs. I should start by looking around to identify any potential locations that match this description.

The task requires me to find a human-shaped statue on a table below a wall with an iris painting. I should start by looking around to identify any potential locations that match this description.

The task requires me to find a human-shaped statue on a table below a wall where the iris painting hangs. I should start by exploring the room and looking for any paintings or statues. Since I am currently facing a blank wall, I should turn left to explore other parts of the room.

The task requires me to find a human-shaped statue on a table below a wall where the iris painting hangs. I should move towards the wall with the iris painting to locate the statue.

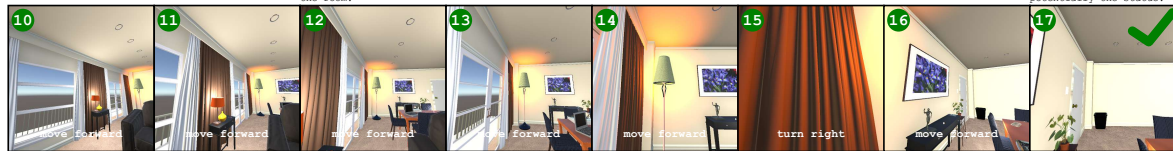
The task requires me to find a human-shaped statue on a table below the wall where the iris painting hangs. I should move towards the wall with the iris painting to locate the statue.

The task requires me to find a human-shaped statue on a table below the wall where the iris painting hangs. I should move towards the wall with the iris painting to locate the statue.

The task requires me to find a human-shaped statue on a table below a wall with an iris painting. I should move towards the wall with the iris painting to locate the statue.

The task requires me to find a human-shaped statue on a table below a wall where the iris painting hangs. I should move towards the area where the iris painting is located, which is visible in the current view. Moving forward will bring me closer to the painting and potentially the statue.

The task requires me to find a human-shaped statue on a table below a wall where the iris painting hangs. I should move towards the area where the iris painting is located, which is visible in the background of the current view. Moving forward will bring me closer to that area.



The task requires me to find a human-shaped statue on a table below a wall where the iris painting hangs. I should move towards the area where the iris painting is located, which is visible in the background of the current view. Moving forward will bring me closer to the painting and potentially the statue.

The task requires me to find a human-shaped statue on a table below a wall where the iris painting hangs. I should move towards the area where the iris painting is located, which is visible in the current view. Moving forward will bring me closer to the painting and potentially the statue.

The task requires me to find a human-shaped statue on a table below a wall where the iris painting hangs. I should move towards the area where the iris painting is located, which is visible in the current view. Moving forward will bring me closer to the target location.

The task is to find a human-shaped statue on a table below a wall where the iris painting hangs. I need to move closer to the table where the statue is located.

The task is to find a human-shaped statue on a table below a wall where the iris painting hangs. I should move closer to the table to confirm if the statue is there.

The task requires me to find a human-shaped statue placed on a table below a wall where the iris painting hangs. I have been moving around the room and now I am near the window. I should turn right to face the wall with the iris painting and locate the statue.

The task is to find a human-shaped statue placed on a table below the wall where the iris painting hangs. I should move closer to the table to inspect it.

The task is to find a human-shaped statue placed on a table below the wall where the iris painting hangs. I should move closer to the table to inspect it.

Figure 16

### 6.4.4 Object Interaction

**Task: Is there an egg inside the fridge?**

**Options:**

- answer "Yes, and there are two eggs, one in the upper part and one in the lower part."
- answer "Yes, in the lower part of the fridge."
- answer "Yes, there are two eggs in the upper part of the fridge."
- open the lower fridge door
- answer "There is no fridge."
- answer "No, there are no eggs in the fridge."
- answer "Yes, in the upper part of the fridge."
- open the upper fridge door
- answer "Yes, there are two eggs in the lower part of the fridge."

**Qwen-VL-Max: answer "Yes, in the lower part of the fridge."**









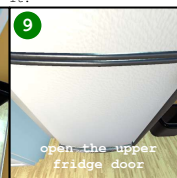




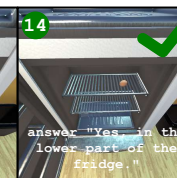
 <p>turn left</p> <p>I need to find the fridge first to check for eggs. The current view does not show the fridge, so I should turn to face it.</p>	 <p>turn left</p> <p>To determine if there is an egg inside the fridge, I need to first locate the fridge and then check its contents. The current view does not show the fridge, so I should turn to find it.</p>	 <p>turn right</p> <p>To check for eggs in the fridge, I need to go to the fridge and open it. The fridge is located to the right of my current position.</p>	 <p>turn right</p> <p>To check for eggs in the fridge, I need to go to the fridge and open it. The fridge is located to the right of my current position.</p>	 <p>move forward</p> <p>To check for eggs in the fridge, I need to go to the fridge and open it. The fridge is not visible in the current view, so I should move towards it.</p>	 <p>turn left</p> <p>To check for eggs in the fridge, I need to go to the fridge and open it. The current view does not show the fridge, so I need to turn towards it.</p>	 <p>move forward</p> <p>To check for eggs in the fridge, I need to go closer to the fridge and then open it to inspect its contents.</p>
 <p>move forward</p> <p>To check for eggs in the fridge, I need to approach the fridge and open it. The current view shows the fridge in the distance, so I should move closer to it.</p>	 <p>open the upper fridge door</p> <p>I need to check inside the fridge to determine if there are any eggs. The current view shows the top of the fridge, so I should open it to see the contents.</p>	 <p>open the upper fridge door</p> <p>I need to check inside the fridge to determine if there are any eggs. Since I am already at the top of the fridge, I should open the doors to inspect the contents.</p>	 <p>open the upper fridge door</p> <p>I need to check inside the fridge to determine if there are any eggs. The previous attempts to open the fridge were unsuccessful, so I will try again.</p>	 <p>open the lower fridge door</p> <p>I need to check inside the fridge to determine if there are any eggs. The previous attempts to open the upper fridge door were unsuccessful, so I should try opening the lower fridge door instead.</p>	 <p>open the upper fridge door</p> <p>I need to check if there is an egg inside the fridge. I should open the fridge doors to look inside.</p>	 <p>answer "Yes, in the lower part of the fridge."</p> <p>The task is to determine if there is an egg inside the fridge. I have already opened the lower fridge door and can see an egg inside. Therefore, I can confirm that there is an egg in the fridge.</p>

Figure 17

**Task: Pick all round objects into trash can.**

**Options:**

- pick up  place to the trash can

**Gemini-Pro:**



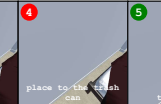









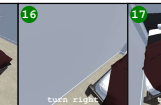
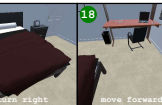






 <p>pick up</p> <p>I see a round object (a basketball). I should pick it up and place it in the trash can.</p>	 <p>move forward</p> <p>I failed, likely because I was too far from the ball. I need to move closer to the ball before trying to pick it up again.</p>	 <p>pick up</p> <p>The pick up action failed, likely because I was too far from the ball. I moved forward, and now I am closer to the ball. I should try to pick it up again.</p>	 <p>place to the trash can</p> <p>I successfully picked up the basketball. Now I need to get it in the trash can.</p>	 <p>turn right</p> <p>I successfully picked up the basketball, but failed to place it in the trash can. This likely means the trash can is not in my current view. I should turn around to find it.</p>	 <p>turn left</p> <p>I successfully picked up the basketball but failed to place it in the trash can. This likely means the trash can is not in my current view. I've already turned right to look for the trash can, but I don't see it. I'll try turning left to see if I can find the trash can.</p>	 <p>turn left</p> <p>I successfully picked up the basketball, but failed to place it in the trash can. This likely means the trash can is not in my current view. I need to turn around to find it. Since the trash can is not in my current view, I will turn right to explore the room and hopefully find the trash can and see what's in that direction.</p>	 <p>turn right</p> <p>I have turned left several times and ended up facing the wall. I need to find the trash can and the remaining round objects. The last successful pick up was the basketball. I will turn right to explore the room and hopefully find the trash can and see what's in that direction.</p>	 <p>turn left</p> <p>I successfully picked up the basketball, but failed to place it in the trash can. This means the trash can is not in my current view. I need to explore the room to find it. I will turn left to search for the trash can.</p>	 <p>turn left</p> <p>I successfully picked up the basketball and attempted to place it in the trash can, but failed. This likely means the trash can isn't in my current view. I need to explore the room to find it. I'll turn left to see what's in that direction.</p>	<p>place to the trash can</p> <p>I have picked up one round object and tried to place it in the trash can, but failed. I have turned around a bit, and now I see a trash can, which didn't help, so let me try placing the object in it now.</p>	<p>turn right</p> <p>I have successfully picked up a round object (Basketball) but failed to place it in the trash can. I turned a bit, and now I see a trash can, which didn't help, so let me try placing the object in it now. I look like the trash can is close to my current location, so I'll turn towards it and try again.</p>
 <p>turn right</p> <p>I have successfully picked up one round object (Basketball), but failed to place it in the trash can. I've tried placing it from different angles, but it can't see from previous angles that it's near the desk. I will turn around and navigate to the desk to find the trash can. I'm currently looking at the wall where I failed to place the object last time, so turning right might give me a view of the desk.</p>	 <p>turn right</p> <p>I successfully picked up a round object and tried to place it in the trash can, but failed. I suspect I'm not facing the trash can. I need to turn around and look for the trash can. I've already turned left and right multiple times, so I should try placing the basketball in it.</p>	 <p>turn left</p> <p>I have picked up one round object and tried to place it in the trash can, but failed. I suspect I'm not facing the trash can. I need to turn around and look for the trash can. I've already turned left and right multiple times, so I should try placing the basketball in it.</p>	 <p>turn right</p> <p>I have picked up one round object and tried to place it in the trash can, but failed. I suspect I'm not facing the trash can. I need to turn around and look for the trash can. I've already turned left and right multiple times, so I should try placing the basketball in it.</p>	 <p>turn right</p> <p>I successfully picked up a round object, but failed to place it in the trash can. I suspect I'm not facing the trash can. I need to turn around and look for the trash can. I've already turned left and right multiple times, so I should try placing the basketball in it.</p>	 <p>move forward</p> <p>I successfully picked up a round object, but failed to place it in the trash can. I suspect I'm not facing the trash can. I need to turn around and look for the trash can. I've already turned left and right multiple times, so I should try placing the basketball in it.</p>	 <p>move forward</p> <p>I have successfully picked up one round object and tried to place it in the trash can, but failed. I suspect I'm not facing the trash can. I need to turn around and look for the trash can. I've already turned left and right multiple times, so I should try placing the basketball in it.</p>	 <p>place to the trash can</p> <p>I have picked up one round object and tried to place it in the trash can, but failed. I suspect I'm not facing the trash can. I need to turn around and look for the trash can. I've already turned left and right multiple times, so I should try placing the basketball in it.</p>	 <p>turn right</p> <p>I have picked up one round object and tried to place it in the trash can, but failed. I suspect I'm not facing the trash can. I need to turn around and look for the trash can. I've already turned left and right multiple times, so I should try placing the basketball in it.</p>	 <p>place to the trash can</p> <p>I have picked up one round object and tried to place it in the trash can, but failed. I suspect I'm not facing the trash can. I need to turn around and look for the trash can. I've already turned left and right multiple times, so I should try placing the basketball in it.</p>	<p>answer "Yes, in the lower part of the fridge."</p>	

Figure 18

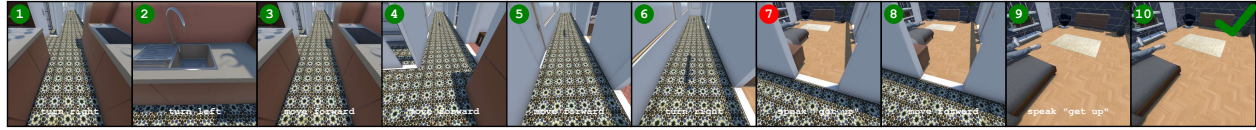
### 6.4.5 Social Interaction

**Task:** Wake up my dad. He is sleeping in the bedroom. The bedroom is the second room on your right as you walk forward.

**Options:**

speak "get up"

**GPT-4o-Mini:**



1 To wake up your dad, I need to navigate to the bedroom, which is the second room on my right as I walk forward. Currently, I am facing a hallway with options to turn left or right. Since the bedroom is to the right, I should turn right to head in the correct direction.

2 I need to navigate towards the bedroom to wake up my dad. Since I have already turned right, the next logical step is to continue moving in the direction that will lead me to the bedroom. The available options include turning left or right, looking up, or speaking. Since turning right again would not be effective, I will turn left to continue towards the bedroom.

3 I need to navigate towards the bedroom, which is the second room on my right as I walk forward. Since I have already turned right and left without making progress, I should move forward to get closer to the bedroom.

4 I need to reach the bedroom, which is the second room on my right as I walk forward. Currently, I am facing a hallway with a clear path ahead. Since I have already turned right and left previously, I should move forward to get closer to the bedroom.

5 I need to wake up my dad, who is in the bedroom, which is the second room on my right as I walk forward. Currently, I am facing a hallway with a view that suggests I am close to the bedroom. Since I have already moved forward twice and turned right and left without success, I should now move forward to reach the bedroom.

6 I need to wake up my dad, who is in the bedroom, which is the second room on my right as I walk forward. Currently, I am facing a hallway, and I have already moved forward multiple times. Since I am not yet in the forward twice and turned right to head towards the bedroom.

7 I am currently in front of the bedroom where my dad is sleeping. Since the task is to wake him up, the most direct action is to speak to him. This will allow me to complete the task without unnecessary movement.

8 I am currently in front of the bedroom where my dad is sleeping. I have already tried to wake him by speaking, but it failed. Since I am already facing the bedroom, I should try to move forward into the room to get closer to him. This will allow me to interact more effectively.

9 I am currently in the bedroom where my dad is sleeping. I previously attempted to wake him by speaking, but it failed. Since I am already in the right location, I should try a different approach. I can either try to speak again or look for a way to wake him up physically. However, since I cannot interact with objects or people unless they are very close, I will try speaking again to see if it has a different effect this time.

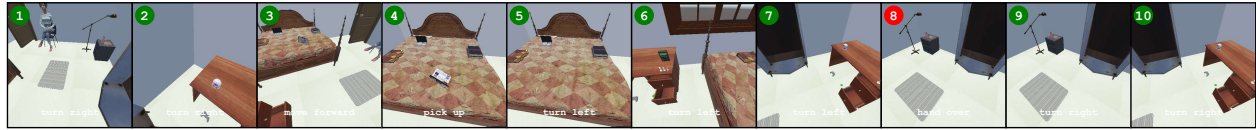
Figure 19

**Task:** Retrieve the magazine for the madam sitting on the chair.

**Options:**

pick up  hand over

**Gemini-Pro:**



1 I see the madam and a magazine on the small table. I need to go to the table, pick up the magazine and give it to the madam. I will turn right to face the table.

2 I turned right and saw a table. I need to turn back to the madam and locate the magazine.

3 I've turned right twice and now see the bed with some books on it. The task is to retrieve the magazine. One of the books might be the magazine. I should move closer to inspect them.

4 I see several objects on the bed, including what looks like a magazine. I should pick it up.

5 I have picked up the magazine. Now I need to turn around to face her and then walk towards her.

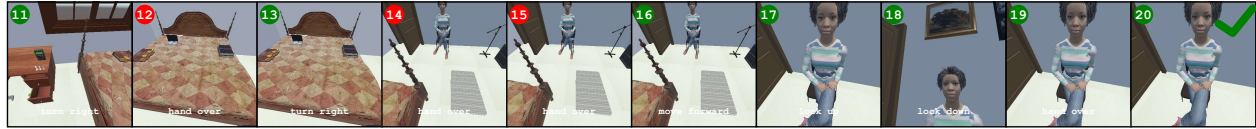
6 I have picked up the magazine. Now I need to turn around and give it to the madam sitting on the chair. I just turned left, so I will turn left again to face her.

7 I have picked up the magazine and turned around, facing the desk. Now I need to turn to face the madam and approach her to hand over the magazine. I've already turned left twice, and the madam was initially to my right, so turning left again should face me in the correct direction.

8 I have picked up the magazine and turned around to face the madam. I should now give her the magazine.

9 The handover failed because I wasn't facing the madam. I need to turn right to face her and then try handing over the magazine again.

10 The handover failed because I wasn't facing the madam. I need to turn towards her and try again. I am currently facing the desk. The madam is behind me to the left. I have turned right after the failed handover, so now she should be directly behind me.



11 I successfully picked up the magazine but failed to hand it over. This likely means I wasn't facing the woman when I tried. I need to turn around to face her and try again. I'm currently facing the bed and desk, so turning left would put me back in the position where I unsuccessfully tried to hand over the magazine. Turning right would orient me towards the doorway, so I will turn right.

12 I successfully picked up the magazine, but handing it over failed. This means I wasn't facing the woman when I tried. I need to turn around to face her and then hand over the magazine. I turned around three times after picking it up. I've done another three turns since the failed handover, meaning I should now be facing the correct direction after three right turns.

13 I successfully picked up an item, but handing it over failed twice. This suggests I picked up the wrong item. I'm currently looking at a few items here. I need to find the magazine. I'll turn around to see if it is elsewhere.

14 I have successfully picked up the magazine, I previously failed to hand it over because I was not facing the madam. I am now facing her, so I should try handing it over again.

15 I have the magazine, I failed to hand it over multiple times because I was not facing the madam. I am now facing her, so I should try handing it over again.

16 I have the magazine, and I've tried handing it over multiple times from slightly different angles, but it hasn't worked. Perhaps I need to move a bit closer to ensure I'm within the interaction range.

17 I have the magazine and I am facing the madam. I failed to hand over the magazine multiple times, likely because I wasn't close enough. I moved forward in the previous step, but perhaps not enough. I will try moving closer again before attempting to hand over.

18 I successfully picked up the magazine earlier, and then tried handing it over multiple times, likely because I wasn't close enough. I moved forward in the last step, so now I'll try handing it over again.

19 I have the magazine and I am facing the woman. I failed to hand it over multiple times, likely because I wasn't close enough. I moved forward in the last step, so now I'll try handing it over again.

20 I have the magazine and I am facing the woman. I failed to hand it over multiple times, likely because I wasn't close enough. I moved forward in the last step, so now I'll try handing it over again.

Figure 20

## 6.5. Error Cases

1111

### 6.5.1 Attribute QA

1112

**Task: Locate any electronic devices in the bedroom.**

**Options:**

- answer "flashlight, desk lamp, and phone"
- answer "laptop, desk lamp, and phone"
- answer "laptop, TV, and desk lamp"
- answer "TV, laptop, and door"
- answer "flashlight, desk lamp, and TV"
- answer "flashlight, desk lamp, and ipad"
- answer "laptop, desk lamp, and bed"
- answer "laptop, book, and desk lamp"

**Qwen-VL-Max:** answer "laptop, desk lamp, and phone"

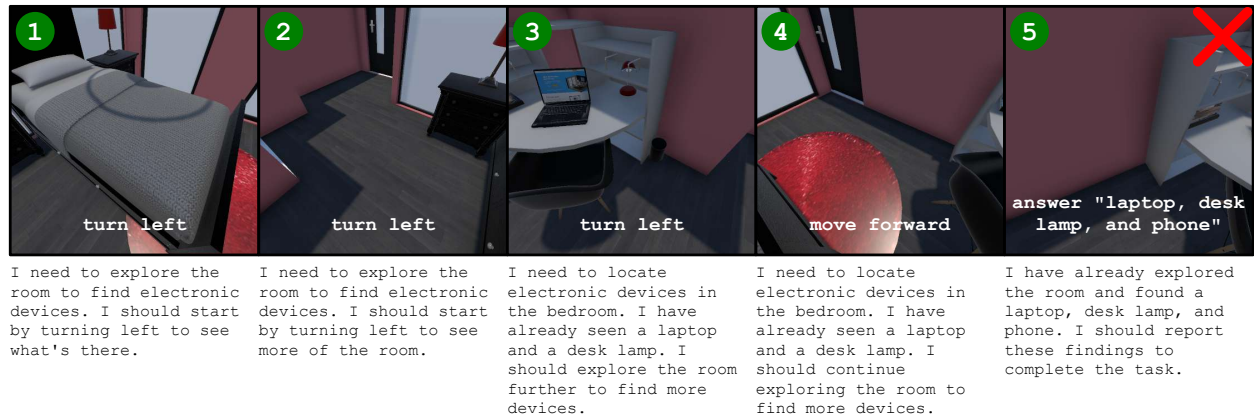


Figure 21

**Task: Are there more flower pots in the living room or the bedroom?**

**Options:**

- answer "There are more flower pots in the living room, while no one in bedroom"
- answer "The bedroom has a few flower pots, but the living room has none."
- answer "The bedroom has only one flower pot, while the living room has several."
- answer "There are no flower pots in the bedroom, only in the living room."
- answer "The living room has more flower pots, placed near the trash can."
- answer "There is only one flower pot in each room."
- answer "Both the living room and bedroom have the same number of flower pots."
- answer "Both rooms have flower pots, but the bedroom's pots are larger in size."

**Gemini-Pro:** answer "There are no flower pots in the bedroom, only in the living room."

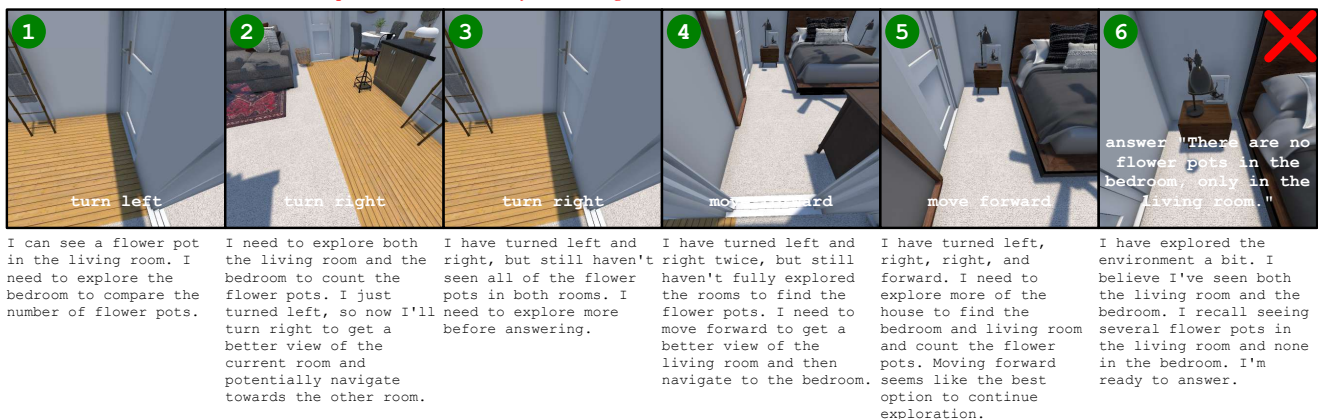


Figure 22

1113 6.5.2 Spatial QA

Task: Determine the best location for a fire extinguisher that would be quickly accessible from both the kitchen and the living room.

Options:

- answer "Near the window."
- answer "Beside the flowerpot near the round carpet."
- answer "On the countertop."
- answer "Beside the black sofas."
- answer "In front of the stove."
- answer "Next to the piano."
- answer "Beside the white sofa."
- answer "Next to the irregular table with many books."

GPT-4o: answer "Beside the white sofa."

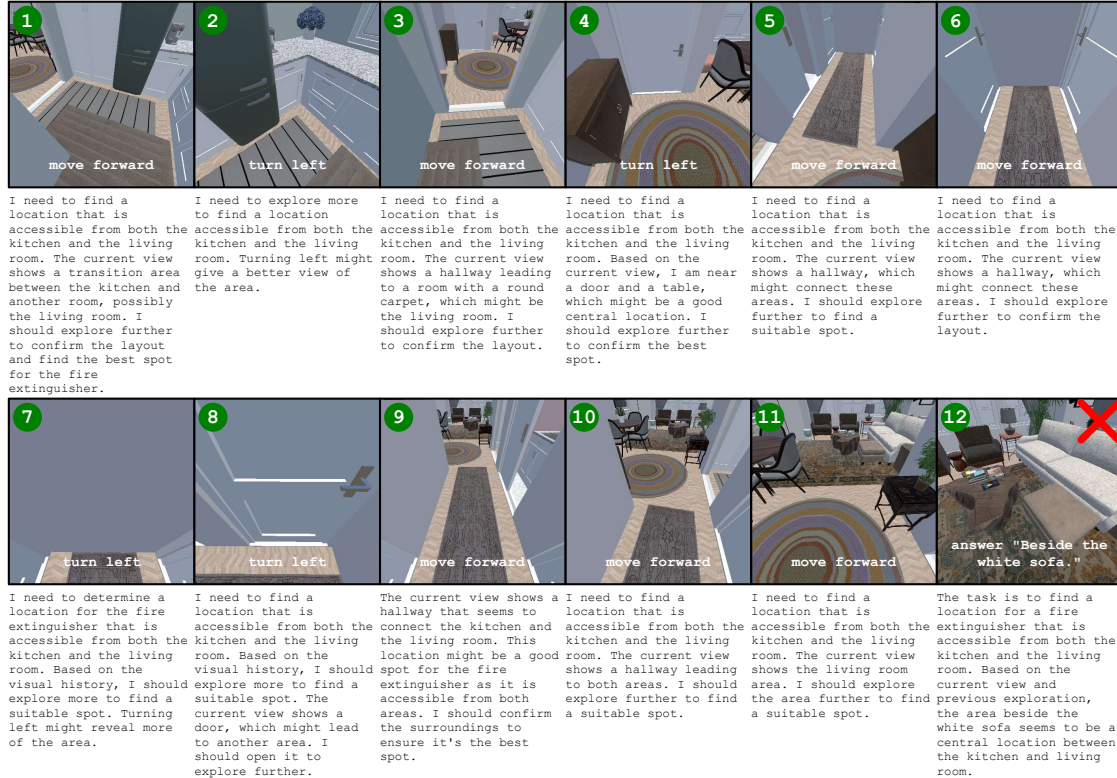


Figure 23

1114 6.5.3 Navigation

Task: Walk towards the tallest tree in the yard.

GPT-4o:



Figure 24

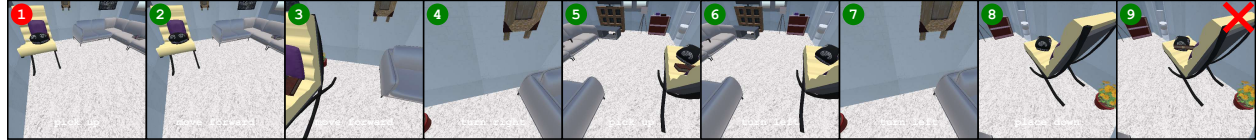
### 6.5.4 Object Interaction

1115

**Task:** Retrieve the tied reading materials from the sofa in the living room and place it on the lounge chair.

**Options:**  
 place down    pick up

**GPT-4o:**



The task is to retrieve the tied reading materials from the sofa and place them on the lounge chair. In the current view, I can see the sofa and the reading materials on it. The lounge chair is also visible. I need to pick up the reading materials first.

The task is to retrieve the tied reading materials from the sofa and place them on the lounge chair. The previous attempt to pick up the reading materials failed, possibly because I wasn't close enough. I should move closer to the sofa to successfully pick them up.

The task is to retrieve the tied reading materials from the sofa and place them on the lounge chair. The current view shows the sofa and the lounge chair. Since the reading materials are likely on the sofa to my right, based on the previous images, I should turn right to face the sofa and locate the materials.

The tied reading materials are visible on the sofa in the current view. I need to pick them up before moving to the lounge chair.

I have successfully picked up the reading materials. Now, I need to place them on the lounge chair, which is visible in my current view. The lounge chair is to my left.

I have successfully picked up the reading materials. Now, I need to place them on the lounge chair. The current view shows the wall, so I need to turn to face the lounge chair.

I have successfully picked up the reading materials. Now, I need to place them on the lounge chair, which is currently in my view.

Figure 25

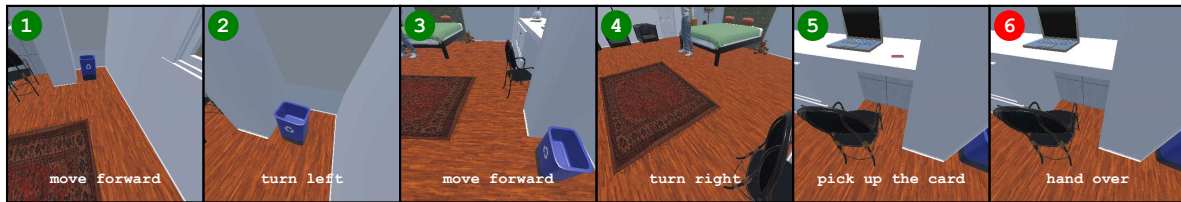
### 6.5.5 Social Interaction

1116

**Task:** Bring my credit card to me and return it back to its original place.

**Options:**  
 pick up the card    place down    hand over    place into trash can

**GPT-4o:**



I need to locate the credit card first. The current view shows a trash can and part of a room. I should explore the area to find the card.

The current view shows a trash can, but I don't see a credit card. I should turn to explore the area further.

I need to locate the credit card before I can pick it up. The current view shows a bed, a chair, and a desk. I should explore the area further to find the card.

The task is to bring the credit card to the person and return it to its original place. Currently, I don't see a credit card in my view. I should explore the room further to locate it. Turning right might give a better view of the room.

The credit card is visible on the desk in front of me. The next logical step is to pick it up.

I have already picked up the card. The next step is to hand it over to the person in the room. The person is visible in the previous images, so I should hand over the card.



The task is to bring the credit card to the person and then return it to its original place. The previous attempt to hand over the card failed, possibly because the person was not in the correct position. I need to locate the person and try handing over the card again.

I am currently facing the person who requested the credit card. The previous attempt to hand over the card failed, possibly because I wasn't close enough. I should move closer to the person to successfully hand over the card.

I am currently facing the person who requested the credit card. The previous attempt to hand over the card failed, possibly because I wasn't close enough. I should try handing over the card again now that I'm closer.

The previous attempts to hand over the card failed, possibly because I wasn't close enough to the person. I need to move closer to the person in the current view to successfully hand over the card.

I am currently facing the chairs, and the person is likely behind me. I need to turn around to face the person and hand over the card. Previous attempts to hand over the card failed, possibly because I wasn't facing the person directly.

Figure 26